

---

# Generalization Analysis for Multi-Label Learning

---

Yi-Fan Zhang<sup>1,2</sup> Min-Ling Zhang<sup>2,3</sup>

## Abstract

Despite great advances in algorithms for multi-label learning, research on the theoretical analysis of generalization is still in the early stage. Some recent theoretical results have investigated the generalization performance of multi-label learning under several evaluation metrics, however, how to reduce the dependency on the number of labels, explicitly introduce label correlations, and quantitatively analyze the impact of various inductive biases in the generalization analysis of multi-label learning is still a crucial and open problem. In an attempt to make up for the gap in the generalization theory of multi-label learning, we develop several novel vector-contraction inequalities, which exploit the Lipschitz continuity of loss functions, and derive generalization bounds with a weaker dependency on the number of labels than the state of the art in the case of decoupling the relationship among different components, which serves as theoretical guarantees for the generalization of multi-label learning. In addition, we derive the generalization bound for Macro-Averaged AUC and analyze its relationship with class-imbalance. The mild bounds without strong assumptions explain the good generalization ability of multi-label learning with first-order label correlations and high-order label correlations induced by norm regularizers.

## 1. Introduction

Multi-label learning is one of the most studied and important machine learning paradigms in practice, in which each object is represented by a single instance while being asso-

ciated with a set of labels instead of a single label. The goal of multi-label learning is to learn a hypothesis which can predict the proper sets of labels for unseen instances. It has made important advances in text categorization (Schapire & Singer, 2000; Rubin et al., 2012), multimedia content annotation (Boutell et al., 2004; Cabral et al., 2011), bioinformatics (Barutcuoglu et al., 2006; Cesa-Bianchi et al., 2012) and other fields (Yu et al., 2005). Although multi-label learning has achieved impressive empirical advances across a wide range of tasks (Zhang & Zhou, 2014), the problem of understanding multi-label learning theoretically remains relatively under-explored.

As we all know, the generalization ability, i.e., the performance of learning machines trained on certain datasets on unseen data, of learning machines is an important question of theoretical research in machine learning, and it is no exception for multi-label learning. Efforts to explain why multi-label models generalize well is an important open problem in multi-label learning community. Uniform convergence is a powerful tool in learning theory for understanding the generalization ability of learners, and it is also used in the generalization analysis of multi-label learning (Yu et al., 2014; Xu et al., 2016; Wu & Zhu, 2020; Wu et al., 2021b;a). However, the progress on the generalization analysis of multi-label learning appears to be severely scarce. A satisfactory and complete study of the generalization analysis for multi-label learning should include three aspects: 1) the reduction of the dependency on the number of labels of the generalization bounds, 2) the explicit introduction of label correlations in the generalization analysis, and 3) the impact of various inductive biases on the generalization performance. First of all, the generalization analysis of multi-label learning is more difficult than that of traditional supervised learning since their difference in problem settings. In particular, the vector-valued output of multi-label learning makes the typical theoretical results not applicable to multi-label learning (Maurer, 2016; Wu & Zhu, 2020). Hence, how to reduce the dependency on the number of labels of the generalization bounds is a very critical problem. Secondly, the consideration of label correlations can often effectively improve the generalization performance of multi-label learning (Zhang & Zhou, 2014), so it is very necessary to explicitly and formally introduce label correlations in the generalization analysis. Finally, the impressive em-

---

<sup>1</sup>School of Cyber Science and Engineering, Southeast University, Nanjing 210096, China <sup>2</sup>Key Laboratory of Computer Network and Information Integration (Southeast University), Ministry of Education, China <sup>3</sup>School of Computer Science and Engineering, Southeast University, Nanjing 210096, China. Correspondence to: Min-Ling Zhang <zhangml@seu.edu.cn>.

empirical success of multi-label learning algorithms motivates us to further investigate the inductive biases induced by these algorithms, e.g., the inductive bias induced by various norm regularizers (Huang et al., 2016) or label-specific features (Zhang & Wu, 2015; Hang & Zhang, 2022; Jia et al., 2023). Therefore, theoretical research on generalization can promote a better understanding of multi-label learning.

In this paper, we derive novel and tighter bounds based on the Rademacher complexity for multi-label learning. Specifically, for  $\ell_2$  Lipschitz loss, we improve the basic linear dependent bound to be independent on the number of labels, which decouples the relationship among different components. For  $\ell_\infty$  Lipschitz loss, we improve the square-root dependent bound to be independent on the number of labels, which also decouples the relationship among different components. These bounds are tighter than the state of the art. We also give several tight bounds for the coupling case to study the impact of different types of label correlations on the generalization analysis. Finally, we derive the generalization bound based on the label-based ranking multi-label Rademacher complexity for Macro-Averaged AUC, and analyze the relationship between Macro-Averaged AUC and class-imbalance.

Our generalization bounds reduce the dependency on the number of labels and account for different types of label correlations. Major contributions of the paper include:

- We prove several novel vector-contraction inequalities for the generalization analysis of multi-label learning, which exploits the Lipschitz continuity of the loss function with respect to the  $\ell_2$  and  $\ell_\infty$  norm and decouples the relationship among different components.
- We derive generalization bounds for general function classes with a weaker dependency on the number of labels than the state of the art, which provides general theoretical guarantees for multi-label learning with different types of label correlations.
- We introduce the label-based ranking multi-label Rademacher complexity and analyze the relationship between Macro-Averaged AUC and class-imbalance according to the generalization bound.

We structure our work as follows. We first introduce an overview of the problem setting for multi-label learning, the definitions of the related evaluation metrics and complexity measures in Section 2. We then present our main results in Sections 3, where we develop several novel vector-contraction inequalities and derive bounds with a weaker dependency on the number of labels than the state of the art for  $\ell_2$  and  $\ell_\infty$  Lipschitz loss, and study the impact of different types of label correlations. In Section 4, we derive

the bound for Macro-Averaged AUC and analyze its relationship with class-imbalance. In Section 5, we provide a comparison of our theoretical results with the related works. Finally, we give a conclusion of our work in Section 6.

## 2. Preliminaries

In this section, we first present the problem setting for multi-label learning. Secondly, we give the definitions of commonly used surrogate losses. Finally, we introduce the related complexity measures involved in the main results.

### 2.1. Multi-Label Learning

In the context of multi-label learning, given a dataset  $D = \{(\mathbf{x}_1, Y_1), \dots, (\mathbf{x}_n, Y_n)\}$  with  $n$  examples which are identically and independently distributed (i.i.d.) from a probability distribution  $P$  on  $\mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X} \subset \mathbb{R}^d$  denotes the  $d$ -dimensional input space and  $\mathcal{Y} = \{y_1, \dots, y_c\}$  denotes the label space with  $c$  class labels,  $\mathbf{x}_i \in \mathcal{X}$ ,  $Y_i \subseteq \mathcal{Y}$ . Let  $[n] := \{1, \dots, n\}$  for any natural number  $n$  and  $i \in [n]$ .

Let  $\mathcal{Y} = \{-1, +1\}^c$ , i.e., each  $\mathbf{y} = (y_1, \dots, y_c)$  is a binary vector and  $y_j = 1$  ( $y_j = -1$ ) denotes that the  $j$ -th label is relevant (irrelevant),  $j \in [c]$ . The task of multi-label learning is to learn a multi-label classifier  $\mathbf{h} \in \mathcal{H} : \mathcal{X} \mapsto \{-1, +1\}^c$  which assigns each instance with a set of relevant labels. A common strategy is to learn a vector-valued function  $\mathbf{f} = (f_1, \dots, f_c) : \mathcal{X} \mapsto \mathbb{R}^c$  and derive the classifier by a thresholding function which dichotomizes the label space into relevant and irrelevant label sets.

We consider the prediction function for each label of the general form  $f_j(\mathbf{x}) = \langle \mathbf{w}_j, \phi(\mathbf{x}) \rangle$ , where  $\phi$  represents a nonlinear mapping. We define a vector-valued function class of the multi-label learning as follows:

$$\begin{aligned} \mathcal{F} = \{ \mathbf{x} \mapsto \mathbf{f}(\mathbf{x}) : \mathbf{f}(\mathbf{x}) &= (f_1(\mathbf{x}), \dots, f_c(\mathbf{x})), \\ f_j(\mathbf{x}) &= \langle \mathbf{w}_j, \phi(\mathbf{x}) \rangle, \mathbf{x} \in \mathcal{X}, j \in [c] \\ \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_c) &\in \mathbb{R}^{d \times c}, \alpha(\mathbf{w}) \leq \Lambda, \\ \beta(\phi(\mathbf{x})) &\leq A, \Lambda > 0, A > 0 \}, \quad (1) \end{aligned}$$

where  $\alpha$  represents a functional that constrains weights,  $\beta$  represents a functional that constrains nonlinear mappings.

For any function  $\mathbf{f} : \mathcal{X} \mapsto \mathbb{R}^c$ , the quality of a prediction on a single example  $(\mathbf{x}, \mathbf{y})$  is measured by a loss function  $L : \mathbb{R}^c \times \{-1, +1\}^c \mapsto \mathbb{R}_+$ . The goal is to learn a hypothesis  $\mathbf{f} \in \mathcal{F}$  with good generalization performance from the dataset  $D$  by optimizing the loss  $L$ . The generalization performance is measured by the expected risk:  $R(\mathbf{f}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} [L(\mathbf{f}(\mathbf{x}), \mathbf{y})]$ . We denote the empirical risk w.r.t. the training dataset  $D$  as  $\widehat{R}_D(\mathbf{f}) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i)$ . In addition, we denote the optimal risk as  $R^* = \inf_{\mathbf{f} \in \mathcal{F}} R(\mathbf{f})$  and denote the minimizer of the empirical risk as  $\widehat{\mathbf{f}}^* = \arg \min_{\mathbf{f} \in \mathcal{F}} \widehat{R}_D(\mathbf{f})$ .

The above definitions apply to Hamming loss, Subset loss and Ranking loss.

For Macro-Averaged AUC, it involves the pairwise loss, so we additionally define the corresponding risks. Maximizing Macro-Averaged AUC is equivalent to minimizing the following empirical risk w.r.t. Macro-Averaged AUC:

$$\begin{aligned} \widehat{R}_D(\mathbf{f}) & \\ &= \frac{1}{c} \sum_{j=1}^c \frac{1}{|X_j^+||X_j^-|} \sum_{\mathbf{x}_i \in X_j^+} \sum_{\mathbf{x}'_i \in X_j^-} \ell_{0/1}(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)), \end{aligned} \quad (2)$$

where  $X_j^+ = \{\mathbf{x}_i \mid y_j = +1, i \in [n]\}$  ( $X_j^- = \{\mathbf{x}'_i \mid y_j = -1, i \in [n]\}$ ) corresponds to the set of test instances that are relevant (irrelevant) to the  $j$ -th label. The expected risk w.r.t. Macro-Averaged AUC is defined as  $R(\mathbf{f}) = \mathbb{E}_D[\widehat{R}_D(\mathbf{f})]$ .

However, the above mentioned loss is typically the 0 – 1 loss, which is hard to handle in practice. Hence, one usually consider its surrogate losses.

## 2.2. Related Evaluation Metrics

A number of evaluation metrics are proposed to measure the generalization performance of different approaches for multi-label learning. Here we focus on commonly used evaluation metrics, i.e., Hamming loss, Subset loss, Ranking loss and Macro-Averaged AUC, and their surrogate losses are defined as follows:

**Hamming Loss:**  $L_H(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \frac{1}{c} \sum_{j=1}^c \ell(y_j f_j(\mathbf{x}))$ ,

where the base convex surrogate loss  $\ell$  can be various popular forms, such as the hinge, logistic and exponential loss.

**Subset Loss:**  $L_S(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \max_{j \in [c]} \{\ell(y_j f_j(\mathbf{x}))\}$ .

**Ranking Loss:**

$$L_R(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \frac{1}{|Y^+||Y^-|} \sum_{p \in Y^+} \sum_{q \in Y^-} \ell(f_p(\mathbf{x}) - f_q(\mathbf{x})),$$

where  $Y^+$  ( $Y^-$ ) denotes the relevant (irrelevant) label index set induced by  $\mathbf{y}$ , and  $|\cdot|$  denotes the cardinality of a set.

The surrogate loss for **Macro-Averaged AUC**:

$$L_M(\mathbf{f}(\mathbf{x}_i, \mathbf{x}'_i), \mathbf{y}) = \frac{1}{c} \sum_{j=1}^c \ell(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)), \quad (3)$$

where  $\mathbf{x}_i$  ( $\mathbf{x}'_i$ ) corresponds to the instances that are relevant (irrelevant) to the  $j$ -th label.

## 2.3. Related Complexity Measures

Here we use the Rademacher complexity to perform generalization analysis for multi-label learning.

**Definition 2.1** (Rademacher complexity). Let  $\mathcal{G}$  be a class of real-valued functions mapping from  $\mathcal{X}$  to  $\mathbb{R}$ . Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set with  $n$  i.i.d. samples. The empirical **Rademacher complexity** over  $\mathcal{G}$  is defined by

$$\widehat{\mathfrak{R}}_D(\mathcal{G}) = \mathbb{E}_\epsilon \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \epsilon_i g(\mathbf{x}_i) \right],$$

where  $\epsilon_1, \dots, \epsilon_n$  are i.i.d. Rademacher random variables. In addition, we define the worst-case Rademacher complexity as  $\mathfrak{R}_n(\mathcal{G}) = \sup_{D \in \mathcal{X}^n} \widehat{\mathfrak{R}}_D(\mathcal{G})$ .

In multi-label learning,  $\mathcal{F}$  is a class of vector-valued functions, which makes traditional Rademacher complexity analysis methods invalid. A common practice is to use the multi-label Rademacher complexity to bound the Rademacher complexity of a loss function space associated with the vector-valued function class  $\mathcal{F}$  according to the vector-contraction inequality in (Maurer, 2016).

**Definition 2.2** (Multi-label Rademacher complexity). Let  $\mathcal{F}$  be a class of vector-valued functions mapping from  $\mathcal{X}$  to  $\mathbb{R}^c$ . Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set with  $n$  i.i.d. samples. The empirical **multi-label Rademacher complexity** over  $\mathcal{F}$  is defined by

$$\widehat{\mathfrak{R}}_D(\mathcal{F}) = \mathbb{E}_\epsilon \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c \epsilon_{ij} f_j(\mathbf{x}_i) \right],$$

where each  $\epsilon_{ij}$  is an independent doubly indexed Rademacher random variable, and  $f_j(\mathbf{x}_i)$  is the  $j$ -th component of  $\mathbf{f}(\mathbf{x}_i)$ .

Here we use the covering number to bound the Rademacher complexity for multi-label learning. The covering number can be bounded by the fat-shattering dimension (Srebro et al., 2010; Lei et al., 2019; Zhang & Zhang, 2023):

**Definition 2.3** (Covering number). Let  $\mathcal{F}$  be a class of real-valued functions mapping from  $\mathcal{X}$  to  $\mathbb{R}$ . Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set with  $n$  i.i.d. samples. For any  $\epsilon > 0$ , the empirical  $\ell_2$  (or  $\ell_\infty$ ) norm covering number  $\mathcal{N}_2(\epsilon, \mathcal{F}, D)$  (or  $\mathcal{N}_\infty(\epsilon, \mathcal{F}, D)$ ) w.r.t.  $D$  is defined as the minimal number  $m$  of a collection of vectors  $\mathbf{v}^1, \dots, \mathbf{v}^m \in \mathbb{R}^n$  such that ( $\mathbf{v}_i^j$  is the  $i$ -th component of the vector  $\mathbf{v}^j$ )

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - \mathbf{v}_i^j)^2} \leq \epsilon. \quad \left( \text{or } \max_{i=1, \dots, n} |f(\mathbf{x}_i) - \mathbf{v}_i^j| \leq \epsilon \right)$$

In this case, we call  $\{\mathbf{v}^1, \dots, \mathbf{v}^m\}$  an  $(\epsilon, \ell_2)$ -cover (or  $(\epsilon, \ell_\infty)$ -cover) of  $\mathcal{F}$  with respect to  $D$ . We also denote  $\mathcal{N}_2(\epsilon, \mathcal{F}, n) = \sup_D \mathcal{N}_2(\epsilon, \mathcal{F}, D)$  (or  $\mathcal{N}_\infty(\epsilon, \mathcal{F}, n) = \sup_D \mathcal{N}_\infty(\epsilon, \mathcal{F}, D)$ ).

**Definition 2.4** (Fat-shattering dimension). Let  $\mathcal{F}$  be a class of real-valued functions mapping from  $\mathcal{X}$  to  $\mathbb{R}$ . We

define the fat-shattering dimension  $\text{fat}_\epsilon(\mathcal{F})$  at scale  $\epsilon > 0$  as the largest  $p \in \mathbb{N}$  such that there exist  $p$  points  $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathcal{X}$  and witnesses  $s_1, \dots, s_p \in \mathbb{R}$  satisfying: for any  $\delta_1, \dots, \delta_p \in \{-1, +1\}$  there exists  $f \in \mathcal{F}$  with

$$\delta_i (f(\mathbf{x}_i) - s_i) \geq \frac{\epsilon}{2}, \quad \forall i = 1, \dots, p.$$

### 3. Generalization Bounds Based on the Rademacher Complexity

In this section, we derive several novel bounds under the assumption that the loss is Lipschitz continuous w.r.t.  $\ell_2$  norm or  $\ell_\infty$  norm. For  $\ell_2$  Lipschitz loss, we first introduce a basic bound with a linear dependency on  $c$ , then we develop a novel vector-contraction inequality and improve the bound with no dependency on  $c$ , which is tighter than the state of the art. For  $\ell_\infty$  Lipschitz loss, we first introduce a basic bound with a square-root dependency on  $c$ , then we develop a novel vector-contraction inequality and improve the bound with no dependency on  $c$ , which is tighter than the state of the art in the decoupling case. We also give several tight bounds for the coupling case to study the impact of different types of label correlations on the generalization analysis. The theoretical results in this section hold for surrogate Hamming loss, surrogate Subset loss and surrogate Ranking loss. The detailed proofs of the theoretical results in this paper are provided in the appendix.

#### 3.1. Generalization Bounds for $\ell_2$ Lipschitz Loss

We first introduce the assumptions used, and we show that for general function classes, Lipschitz continuity of the loss function w.r.t. the  $\ell_2$  norm combined with the multi-label Rademacher complexity yields basic generalization bounds with a linear dependency on the number of labels, which exploits the typical vector-contraction inequality (Maurer, 2016). Second, we develop a novel vector-contraction inequality and derive a tighter bound with no dependency on the number of labels for  $\ell_2$  Lipschitz loss.

**Assumption 3.1.** Assume that the loss function and the components of the vector-valued function are bounded:  $L(\cdot, \cdot) \leq M$ ,  $|f_j(\cdot)| \leq B$  for  $j \in [c]$  where  $M > 0$  and  $B > 0$  are constants.

**Assumption 3.2.** Assume that the loss function is  $\mu$ -Lipschitz continuous w.r.t. the  $\ell_2$  norm.

Assumption 3.1 is a relatively mild assumption. In fact, when we consider the function class (1) for multi-label learning, we often use the assumptions  $\|\mathbf{w}_j\|_2 \leq \Lambda$ ,  $\phi_j$  is  $\rho$ -Lipschitz w.r.t. the  $\ell_2$  norm and  $\|\mathbf{x}_i\|_2 \leq A$  for any  $j \in [c]$ ,  $i \in [n]$  to replace the boundedness of the components of the vector-valued function. The following Proposition 3.3 also illustrates that Assumption 3.2 is very mild.

**Proposition 3.3** (Lemma 1 in (Wu & Zhu, 2020)). *Assume that the base loss function  $\ell$  defined in Subsection 2.2 is  $\mu$ -Lipschitz continuous, then the surrogate Hamming Loss is  $\frac{\mu}{\sqrt{c}}$ -Lipschitz w.r.t. the  $\ell_2$  norm, the surrogate Subset Loss is  $\mu$ -Lipschitz w.r.t. the  $\ell_2$  norm, and the surrogate Ranking Loss is  $\mu$ -Lipschitz w.r.t. the  $\ell_2$  norm.*

#### 3.1.1. A BASIC BOUND FOR $\ell_2$ LIPSCHITZ LOSS

Using the  $\ell_2$  Lipschitz continuity of loss and the multi-label Rademacher complexity, we have the following theorem:

**Theorem 3.4.** *Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.2 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the following holds for any  $\mathbf{f} \in \mathcal{F}$ :*

$$R(\mathbf{f}) \leq \widehat{R}_D(\mathbf{f}) + \frac{2\sqrt{2}\mu c B}{\sqrt{n}} + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

*Proof Sketch.* We first use the multi-label Rademacher complexity to bound the Rademacher complexity of the loss function space associated with the vector-valued function class  $\mathcal{F}$  according to the typical vector-contraction inequality (Maurer, 2016), and then complete the proof with the McDiarmid's inequality and the symmetrization technique.  $\square$

*Remark 3.5.* The above bound with a linear dependency on the number of labels indicates that good generalization performance will be obtained when the number of examples ( $\sqrt{n}$ ) is larger than the number of labels ( $c$ ), but in practice, it is often encountered in the case of an extremely large number of labels, that is, extreme multi-label learning (Yu et al., 2014; Prabhu & Varma, 2014; Yen et al., 2016; Liu & Shen, 2019). At this time, the number of labels will probably be more than the number of examples, so the bound in Theorem 3.4 will not be able to provide theoretical guarantees, thus prompting us to develop the bound that is tighter on the number of labels. Our analysis in Theorem 3.4 implies the following inequality:  $\mathbb{E}_\epsilon \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c \epsilon_{ij} f_j(\mathbf{x}_i) \right] \leq c \max_j \mathbb{E}_\epsilon \left[ \sup_{f_j \in \mathcal{F}_j} \frac{1}{n} \sum_{i=1}^n \epsilon_{ij} f_j(\mathbf{x}_i) \right]$ , which shows that decoupling the relationship among different components (since the maximization over  $j \in [c]$  is outside of the expectation operator) will lead to bounds with a linear dependency on  $c$  for general function classes. When considering kernel function classes, the dependency of the bounds on  $c$  can be improved to square-root (Maurer, 2016; Wu & Zhu, 2020; Wu et al., 2021a). Such improvements essentially come from preserving the coupling among different components reflected by the constraint, i.e.,  $\|\mathbf{w}\| \leq \Lambda$ . As a comparison, when  $\|\mathbf{w}_j\|_2 \leq \Lambda$  for any  $j \in [c]$ , if we consider the group norm  $\|\cdot\|_{2,2}$ , we have  $\|\mathbf{w}\|_{2,2} \leq \sqrt{c}\Lambda$ ,

which means that these improved bounds still suffer from a linear dependency on the number of labels. Hence, the improvement in the preservation of coupling by a factor of  $\sqrt{c}$  benefits from replacing  $\Lambda$  with  $\sqrt{c}\Lambda$  in the constraint to some extent. This reveals that in order to improve the existing bounds, we need to improve the dependency on the number of labels in the decoupling case.

### 3.1.2. TIGHTER BOUNDS FOR $\ell_2$ LIPSCHITZ LOSS

We develop a novel vector-contraction inequality for  $\ell_2$  Lipschitz loss, which decouples the relationship among different components and guarantees that the derived generalization bounds are tighter than the state of the art.

We show that the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded by the worst-case Rademacher complexity of the projection function class  $\mathcal{P}(\mathcal{F})$ . We first define a function class  $\mathcal{P}$  consisting of projection operators  $p_j : \mathbb{R}^c \mapsto \mathbb{R}$  for any  $j \in [c]$  which project the  $c$ -dimensional vector onto the  $j$ -th coordinate. Then, we have  $\mathcal{P}(\mathcal{F}) = \{(j, \mathbf{x}) \mapsto p_j(\mathbf{f}(\mathbf{x})) : p_j(\mathbf{f}(\mathbf{x})) = f_j(\mathbf{x}), \mathbf{f} \in \mathcal{F}, (j, \mathbf{x}) \in [c] \times \mathcal{X}\}$ . With the above definitions, we develop the following vector-contraction inequality:

**Lemma 3.6.** *Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.2 hold. Given a dataset  $D$  of size  $n$ . Then, we have*

$$\hat{\mathfrak{R}}_D(\mathcal{F}) \leq 48\mu\sqrt{c}\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \left(1 + \log^{\frac{1}{2}}(nc) \cdot \log \frac{M\sqrt{n}}{\mu B}\right),$$

where  $\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))$  is the worst-case Rademacher complexity of the projection function class.

*Proof Sketch.* First, the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded by the empirical  $\ell_2$  norm covering number with the refined Dudley’s entropy integral inequality. Second, according to the Lipschitz continuity w.r.t the  $\ell_2$  norm, the empirical  $\ell_2$  norm covering number of  $\mathcal{F}$  can be bounded by the empirical  $\ell_2$  norm covering number of  $\mathcal{P}(\mathcal{F})$ . Third, the empirical  $\ell_2$  norm covering number of  $\mathcal{P}(\mathcal{F})$  can be bounded by using Sudakov’s minoration (Wainwright, 2019), which bounds the  $\ell_2$  norm covering number of a function class by the expectation of a Gaussian process indexed by the function class, and the expectation of the Gaussian process can be bounded by the worst-case Rademacher complexity of the projection function class  $\mathcal{P}(\mathcal{F})$ . Hence, the problem is transferred to the estimation of the worst-case Rademacher complexity. Finally, we estimate the lower bound of the worst-case Rademacher complexity of  $\mathcal{P}(\mathcal{F})$ , and then combined with the above steps, the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded.  $\square$

With the vector-contraction inequality above, we can derive the following tight bound for  $\ell_2$  Lipschitz loss:

**Theorem 3.7.** *Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.2 and 3.1 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the following holds for any  $\mathbf{f} \in \mathcal{F}$ :*

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + 3M\sqrt{\frac{\log \frac{2}{\delta}}{2n}} + \frac{96B\mu \left(1 + \log^{\frac{1}{2}}(nc) \cdot \log \frac{M\sqrt{n}}{\mu B}\right)}{\sqrt{n}}.$$

*Proof Sketch.* We first upper bound the worst-case Rademacher complexity  $\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))$ , and then combined with Lemma 3.6, the desired bound can be derived.  $\square$

*Remark 3.8.* Although Lemma 3.6 shows a factor of  $\sqrt{c}$ , the term  $\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \leq \frac{B}{\sqrt{nc}}$ , which makes the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  (i.e.,  $\hat{\mathfrak{R}}_D(\mathcal{F})$ ) actually independent on  $c$ , and results in a tighter bound than the  $O(c/\sqrt{n})$  bound in Theorem 3.4 with a faster convergence rate  $\tilde{O}(1/\sqrt{n})$ . Lemma 3.6 decouples the relationship among different components (since the supremum over  $j \in [c]$  is outside of the expectation operator by the definition of  $\mathcal{P}(\mathcal{F})$ ). Hence, the bound with no dependency on  $c$  in Theorem 3.7 is clearly tighter than the state of the art with square-root dependency on  $c$  in (Lei et al., 2015; Maurer, 2016; Wu & Zhu, 2020; Wu et al., 2021a) for  $\ell_2$  Lipschitz loss assumption, where the analyzes all preserve the coupling among different components, not to mention, as discussed in Remark 3.5, the constraint on preserving the coupling ( $\|\mathbf{w}\| \leq \Lambda$ ) directly implies the improvement by a factor of  $\sqrt{c}$ . In fact, decoupling the relationship or preserving the coupling among different components corresponds to different types of label correlations in multi-label learning (Zhang & Zhou, 2014). The former corresponds to first-order label correlations (which tackle multi-label learning problem by decomposing it into a number of independent binary classification problems, i.e., ignorance of label correlations), and the latter corresponds to high-order label correlations (which tackle multi-label learning problem by exploiting high-order relationships among labels) induced by norm regularizers. The assumption of the coupling case that holds for many learning scenarios (e.g., multi-class learning) does not hold in multi-label learning methods with first-order label correlations. This means that we need to develop new vector-contraction inequalities to handle the assumption of the decoupling case. The projection function class is used to help handle the generalization analysis in the decoupling case. The above bounds provide theoretical guarantees for first-order label correlations methods, e.g., Binary Relevance methods (Boutell et al., 2004; Zhang & Zhou, 2014; Zhang & Wu, 2015; Hang & Zhang, 2022). Furthermore, according to the Proposition 3.3, we

can obtain a  $\tilde{O}(1/\sqrt{nc})$  bound for  $\ell_2$  Lipschitz surrogate Hamming Loss in the decoupling case.

### 3.2. Generalization Bounds for $\ell_\infty$ Lipschitz Loss

We first introduce the assumption of Lipschitz continuity w.r.t. the  $\ell_\infty$  norm, and we derive a basic bound with a square-root dependency on  $c$  for general function classes by refining Theorem 1 in (Foster & Rakhlin, 2019). Second, we develop a novel vector-contraction inequality and derive a tighter bound with no dependency on  $c$  for  $\ell_\infty$  Lipschitz loss in the decoupling case, up to logarithmic terms, and we also give several bounds with no dependency on  $c$  for  $\ell_\infty$  Lipschitz loss in the coupling case.

**Assumption 3.9.** Assume that the loss function  $L$  is  $\mu$ -Lipschitz continuous w.r.t. the  $\ell_\infty$  norm, that is:

$$|L(\mathbf{f}(\mathbf{x}), \cdot) - L(\mathbf{f}'(\mathbf{x}), \cdot)| \leq \mu \|\mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\|_\infty,$$

where  $\mu > 0$ ,  $\|\mathbf{t}\|_\infty = \max_{j \in [c]} |t_j|$  for  $\mathbf{t} = (t_1, \dots, t_c)$ .

In fact, the commonly used loss functions in multi-label learning actually satisfy the Lipschitz continuity w.r.t. the  $\ell_\infty$  norm, and it has been considered in some literature (Lei et al., 2019; Wu et al., 2021b). The following Proposition 3.10 further illustrates that Assumption 3.9 is very mild.

**Proposition 3.10.** Assume that the base loss  $\ell$  defined in Subsection 2.2 is  $\mu$ -Lipschitz continuous, then the surrogate Hamming Loss is  $\mu$ -Lipschitz w.r.t. the  $\ell_\infty$  norm, the surrogate Subset Loss is  $\mu$ -Lipschitz w.r.t. the  $\ell_\infty$  norm, and the surrogate Ranking Loss is  $2\mu$ -Lipschitz w.r.t. the  $\ell_\infty$  norm.

#### 3.2.1. A BASIC BOUND FOR $\ell_\infty$ LIPSCHITZ LOSS

Using the Lipschitz continuity w.r.t. the  $\ell_\infty$  norm, we first show that the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded by the worst-case Rademacher complexity of the restriction of the function class along each coordinate with timing a factor of  $\tilde{O}(\sqrt{c})$ . Then, we derive the basic bound with a square-root dependency on the number of labels in the decoupling case.

**Lemma 3.11.** Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \eta < 1$ ,  $0 < a < 1$ , we have

$$\hat{\mathfrak{R}}_D(\mathcal{F}) \leq \frac{48\sqrt{2E}\mu}{a} \sqrt{c} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right),$$

where  $\tilde{\mathfrak{R}}_n(\mathcal{F}_j)$  is the worst-case Rademacher complexity,  $\mathcal{F}_j$  is the restriction of the function class along the  $j$ -th coordinate,  $f_j \in \mathcal{F}_j$  and  $E > 0$  is an absolute constant.

*Proof Sketch.* We obtain the lower bound of the worst-case Rademacher complexity  $\tilde{\mathfrak{R}}_n(\mathcal{F}_j)$  through the Khintchine-Kahane inequality, combined with the refined Theorem 1 in (Foster & Rakhlin, 2019), the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded.  $\square$

With the vector-contraction inequality above, we can derive the following basic bound for  $\ell_\infty$  Lipschitz loss:

**Theorem 3.12.** Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \delta, \eta, a < 1$ , there exists an absolute constant  $E > 0$  such that with probability at least  $1 - \delta$ , the following holds for any  $\mathbf{f} \in \mathcal{F}$ :

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}} + \frac{96\sqrt{2E}\mu B}{a} \sqrt{\frac{c}{n}} \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right).$$

*Proof Sketch.* We first upper bound the worst-case Rademacher complexity  $\tilde{\mathfrak{R}}_n(\mathcal{F}_j)$ , and then combined with Lemma 3.11, the desired bound can be derived.  $\square$

*Remark 3.13.* Lemma 3.11 is not a direct application of Theorem 1 in (Foster & Rakhlin, 2019) which absorbs all terms independent of  $c$  and  $n$  into an unspecified numerical constant and yields a logarithmic term of order  $O(\log^{\frac{3+\eta}{2}}(\sqrt{n}))$ . We refine the proof of Theorem 1 in (Foster & Rakhlin, 2019), bound the relevant constant terms and improve the order of the logarithmic term to  $O(\log^{\frac{1+\eta}{2}}(\sqrt{n}) \cdot \log(\sqrt{\frac{n}{c}}))$ . The term  $\max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \leq \frac{B}{\sqrt{n}}$  means that Lemma 3.11 decouples the relationship among different components, which results in a tighter bound than the  $O(c/\sqrt{n})$  bound in Theorem 3.4 with a faster convergence rate  $\tilde{O}(\sqrt{c/n})$ .

#### 3.2.2. TIGHTER BOUNDS FOR $\ell_\infty$ LIPSCHITZ LOSS

We first give a corollary with no dependency on  $c$ , which preserves the coupling among different components.

**Corollary 3.14.** Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the following holds for any  $\mathbf{f} \in \mathcal{F}$ :

**I.)** If  $\alpha(\mathbf{w}) := \|\mathbf{w}\|$ ,  $\beta(\phi(\mathbf{x})) := \sup_{\mathbf{x} \in \mathcal{X}, j \in [c]} \|\tilde{\phi}_j(\mathbf{x})\|_*$ , and  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ , we have:

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{36\mu\Lambda A \log^{\frac{3}{2}}(\sqrt{2n}^{\frac{3}{2}}c)}{\sqrt{n}} + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

2.) If  $\alpha(\mathbf{w}) := \|\mathbf{w}\|_{S_p}$ ,  $\beta(\phi(\mathbf{x})) := \max_{i \in [n]} \|\phi(\mathbf{x}_i)\|_2$ , when  $1 \leq p \leq 2$ , we have:

$$R(\mathbf{f}) \leq \widehat{R}_D(\mathbf{f}) + \frac{36\mu\Lambda A \log_2^{\frac{3}{2}}(\sqrt{2n}^{\frac{3}{2}}c)}{\sqrt{n}} + 3M\sqrt{\frac{\log \frac{2}{\delta}}{2n}},$$

when  $p > 2$ , we have:

$$R(\mathbf{f}) \leq \widehat{R}_D(\mathbf{f}) + 3M\sqrt{\frac{\log \frac{2}{\delta}}{2n}} + \frac{36\mu\Lambda A \min\{c, d\}^{\frac{1}{2} - \frac{1}{p}} \log_2^{\frac{3}{2}}(\sqrt{2n}^{\frac{3}{2}}c)}{\sqrt{n}}.$$

The Schatten- $p$  norm is defined as the  $\ell_p$  norm of the singular value vector of a matrix, i.e.,  $\|\mathbf{w}\|_{S_p} = \|\sigma(\mathbf{w})\|_p$ , where the singular values are sorted in non-increasing order. For any  $\mathbf{x} \in \mathcal{X}$  and  $j \in [c]$ , the notation  $\check{\phi}_j(\mathbf{x})$  is defined by  $\check{\phi}_j(\mathbf{x}) := (\underbrace{0, \dots, 0}_{j-1}, \phi(\mathbf{x}), \underbrace{0, \dots, 0}_{c-j}) \in \mathbb{R}^{d \times c}$ .

*Proof Sketch.* We first show that the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded by the worst-case Rademacher complexity of  $\widetilde{\mathcal{F}}$ , which refines Theorem 5 in (Lei et al., 2019) (i.e.,  $\widehat{\mathfrak{R}}_D(\mathcal{F}) \leq 16\sqrt{\log 2}\mu\sqrt{c}\widetilde{\mathfrak{R}}_{nc}(\widetilde{\mathcal{F}})(1 + \log_2^{\frac{3}{2}}\frac{\Lambda An\sqrt{c}}{\widehat{\mathfrak{R}}_{nc}(\widetilde{\mathcal{F}})})$ , where  $\widetilde{\mathcal{F}} := \{\mathbf{v} \mapsto \langle \mathbf{w}, \mathbf{v} \rangle : \mathbf{w}, \mathbf{v} \in \mathbb{R}^{d \times c}, \alpha(\mathbf{w}) \leq \Lambda, \beta(\mathbf{v}) \leq A, \mathbf{v} \in \widetilde{S}\}$  and  $\widetilde{S} = \{\check{\phi}_j(\mathbf{x}_i) : j \in [c], i \in [n]\}$ . Then, we upper bound  $\widetilde{\mathfrak{R}}_{nc}(\widetilde{\mathcal{F}})$  for different norm regularizers. Combining these results, the desired bounds can be derived.  $\square$

*Remark 3.15.* Corollary 3.14 for preserving the coupling case has less novelty, which has the same order  $\widetilde{O}(1/\sqrt{n})$  as the results in (Lei et al., 2019; Wu et al., 2021b), since these results all use the same vector-contraction inequality, i.e., Theorem 5 in (Lei et al., 2019). However, here we are more concerned with investigating the impact of the label correlation induced by the norm regularizer on the generalization analysis. The bounds here involve a constraint on the overall weight  $\mathbf{w}$ , which consider that the components share some constraint properties with each other, that is, the label correlation induced by the norm regularizer. The above bounds involve label correlations induced by the general norm regularizer (case 1) and the Schatten- $p$  norm regularizer (case 2), respectively. Trace norm regularizer, corresponding to Schatten- $p$  norm regularizer with  $p = 1$ , is a common practice to consider the label correlation in multi-label learning, which imposes a low-rank constraint on the spectrum of  $\mathbf{w}$ . The bounds here also explain the good generalization ability of multi-label learning with the label correlation induced by the trace norm regularizer.

Then, we develop a novel vector-contraction inequality, which guarantees that the derived bounds are tighter than the state of the art in the decoupling case.

**Lemma 3.16.** Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \eta < 1$ ,  $0 < a < 1$ ,  $E > 0$ , we have

$$\widehat{\mathfrak{R}}_D(\mathcal{F}) \leq \frac{96\sqrt{E}\mu}{a}\sqrt{c}\widetilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \times \left(1 + \log^{\frac{1+\eta}{2}}\left(\frac{8}{5B}\sqrt{nc}\right) \cdot \log\left(\frac{M\sqrt{n}}{\sqrt{E}\mu B}\right)\right),$$

where  $\widetilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))$  is the worst-case Rademacher complexity of the projection function class.

*Proof Sketch.* The overall proof idea is similar to Lemma 3.6, but there are two obvious differences. First of all, the covering numbers involved in the proof here are all  $\ell_\infty$  norm covering numbers instead of  $\ell_2$  norm covering numbers, which makes the proof techniques involved completely different. Secondly, in the third step, instead of using Sudakov's minoration inequality, we use the fat-shattering dimension to bound the empirical  $\ell_\infty$  norm covering number of  $\mathcal{P}(\mathcal{F})$ , and the fat-shattering dimension can be bounded by the worst-case Rademacher complexity of  $\mathcal{P}(\mathcal{F})$ .  $\square$

**Theorem 3.17.** Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \delta, \eta, a < 1$ , there exists an absolute constant  $E > 0$  such that with probability at least  $1 - \delta$ , the following holds for any  $\mathbf{f} \in \mathcal{F}$ :

$$R(\mathbf{f}) \leq \widehat{R}_D(\mathbf{f}) + 3M\sqrt{\frac{\log \frac{2}{\delta}}{2n}} + \frac{192\sqrt{E}\mu B}{a}\frac{1}{\sqrt{n}} \times \left(1 + \log^{\frac{1+\eta}{2}}\left(\frac{8}{5B}\sqrt{nc}\right) \cdot \log\left(\frac{M\sqrt{n}}{\sqrt{E}\mu B}\right)\right).$$

*Proof Sketch.* We first upper bound  $\widetilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))$ , then combined with Lemma 3.16, the bound is immediate.  $\square$

*Remark 3.18.* Lemma 3.16 decouples the relationship among different components, and the term  $\widetilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \leq \frac{B}{\sqrt{nc}}$ , which makes the Rademacher complexity  $\widehat{\mathfrak{R}}_D(\mathcal{F})$  actually independent on  $c$ , and results in a tighter bound than the  $\widetilde{O}(\sqrt{c/n})$  bound in Theorem 3.12 with a faster convergence rate  $\widetilde{O}(1/\sqrt{n})$ . How to develop novel vector-contraction inequalities that can induce  $\widetilde{O}(1/\sqrt{n})$  bounds and deal with the assumption of the decoupling case are the two most critical difficulties in deriving tighter bounds. The introduction of the projection function class plays an important role in solving these two difficulties. It improves the vector-contraction inequalities by a factor of  $\sqrt{c}$  and handles the assumption of the decoupling case indirectly. Theorem 3.17 improves the bounds from  $\widetilde{O}(\sqrt{c/n})$  to  $\widetilde{O}(1/\sqrt{n})$  for

$\ell_\infty$  Lipschitz Loss in the decoupling case, and explains the good generalization ability of multi-label learning with first-order label correlations.

#### 4. Generalization Analysis for Macro-Averaged AUC

Macro-Averaged AUC is a typical label-based ranking metric (Zhang & Zhou, 2014). Here we analyze the relationship between Macro-Averaged AUC and class-imbalance according to the generalization bound. We use Lemma 3.11 for the generalization analysis, since improving the dependency on the number of labels is not our main concern here.

Rademacher complexity has proved to be a powerful data-dependent measure of hypothesis space complexity (Bartlett & Mendelson, 2002; Koltchinskii & Panchenko, 2002). However, a sequence of pairs of i.i.d. individual observation in (2) is no longer independent, which makes standard techniques in the i.i.d case for traditional Rademacher complexity inapplicable for Macro-Averaged AUC. We convert the non-sum-of-i.i.d pairwise function to a sum-of-i.i.d form by using permutations in U-process (Cl  men  on et al., 2008). We denote  $|X_j^+|$  and  $|X_j^-|$  in (2) as  $s_j$  and  $t_j$ ,  $r_j = \min\{s_j, t_j\}$ ,  $r_0 = \min_{j \in [c]} \{r_j\}$ , and  $s_j + t_j = n$  for any  $j \in [c]$ . Then, we have the following definition:

**Definition 4.1.** Let  $\mathcal{F}$  be a class of vector-valued functions mapping from  $\mathcal{X}$  to  $\mathbb{R}^c$  and  $\ell$  be the base loss defined in (3). Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set with  $n$  i.i.d. samples. The empirical **label-based ranking multi-label Rademacher complexity** of a loss function space associated with the vector-valued function class  $\mathcal{F}$  is defined by

$$\hat{\mathfrak{R}}_D(\mathcal{F}) = \mathbb{E}_\epsilon \left[ \sup_{f \in \mathcal{F}} \frac{1}{c} \sum_{j=1}^c \frac{1}{r_j} \sum_{i=1}^{r_j} \epsilon_{ij} \ell(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)) \right],$$

where each  $\epsilon_{ij}$  is an independent doubly indexed Rademacher random variable, and  $f_j(\mathbf{x}_i)$  is the  $j$ -th component of  $\mathbf{f}(\mathbf{x}_i)$ . We refer to the expectation  $\mathfrak{R}(\mathcal{F}) = \mathbb{E}_D[\hat{\mathfrak{R}}_D(\mathcal{F})]$  as the label-based ranking multi-label Rademacher complexity of  $\mathcal{F}$ .

With this definition, we then derive the bound as follows:

**Theorem 4.2.** *Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1) and the loss be Macro-Averaged AUC. Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the following holds for any  $\mathbf{f} \in \mathcal{F}$ :*

$$R(\hat{\mathbf{f}}^*) - R^* = \tilde{O}(\sqrt{c/r_0}).$$

*Proof Sketch.* First, by using the U-process technique, we bound  $\mathbb{E}[R(\hat{\mathbf{f}}^*)] - R^*$  with the label-based ranking multi-label Rademacher complexity. Second, combining with

Lemma 3.11 and bounding the worst-case Rademacher complexity, we can upper bound the label-based ranking multi-label Rademacher complexity. Finally, with the McDiarmid’s inequality, the desired bound can be derived.  $\square$

*Remark 4.3.* Theorem 4.2 shows that when class-imbalance occurs,  $r_0$  will be smaller than  $\frac{n}{2}$ . When class-imbalance is more serious,  $r_0$  will be smaller, which will lead to a looser bound for Macro-Averaged AUC. This means that when class-imbalance becomes more and more serious, if the learned classifier cannot handle the problem of class imbalance well, then its performance on Macro-Averaged AUC will be worse. Wu et al. (2023) also obtained similar conclusions, but the methods used were completely different. Wu et al. (2023) transformed the macro-averaged maximization problem in multi-label learning into the problem of learning multiple tasks with graph-dependent examples-which is hard to verify in practice, then proposed a new McDiarmid-type inequality to develop  $O(\frac{1}{\sqrt{n}})$  bound in the balanced case. Our method is simpler and can also yield bounds with no dependency on  $c$  by combining Lemma 3.16.

#### 5. Comparison with Related Work

The generalization analysis of multi-label learning originated from (Dembczynski et al., 2010), which performed regret analysis on Hamming and Subset loss, and derived the relationship between the expectations of Hamming and Subset loss. Dembczynski et al. (2012) performed regret analysis on Ranking loss. These analyses laid the foundation for research in (Wu & Zhu, 2020; Wu et al., 2021a).

Theorem 3.4 shows that when using the typical vector-contraction inequality (Maurer, 2016) (i.e.,  $\ell_2$  norm Lipschitz loss) and the multi-label Rademacher complexity, one can only derive bounds of order  $O(c/\sqrt{n})$  for general function classes in the decoupling case. Wu & Zhu (2020); Wu et al. (2021a) also obtained similar results for surrogate Hamming, Subset and Ranking losses, which mainly exploited the relationship between losses. Wu & Zhu (2020); Wu et al. (2021a) also showed that for kernel function classes, the order of the bounds for some losses can be improved to  $O(\sqrt{c/n})$  when preserving the coupling. Lei et al. (2015) first derived a  $O(\sqrt{c/n})$  bound for multi-class SVM with  $\ell_p$  norm regularized kernel function classes under the assumption of  $\ell_2$  Lipschitz loss. Li et al. (2018) used the local Rademacher complexity to derive a  $O(\log^2 c/n)$  bound for multi-class classification with  $\ell_p$  norm regularized kernel function classes under the assumption of  $\ell_2$  Lipschitz and smooth loss. Theorem 3.7 improves the bounds to  $O(1/\sqrt{n})$  for  $\ell_2$  Lipschitz loss even in the decoupling case.

Theorem 3.12 shows that when using  $\ell_\infty$  norm Lipschitz loss, one can derive bounds of order  $O(\sqrt{c/n})$  for general



function classes. Liu et al. (2018) also obtained a bound of order  $O(\sqrt{c/n})$  for the dual set multi-label learning for margin loss and kernel function classes. Lei et al. (2019); Wu et al. (2021b) improved the dependency on  $c$  in the coupling case, where Lei et al. (2019) derived a  $\tilde{O}(1/\sqrt{n})$  bound for multi-class classification with norm regularized function classes under  $\ell_\infty$  Lipschitz loss, and Wu et al. (2021b) derived a  $O(\log^3(nc)/n\sigma)$  bound for vector-valued learning with norm regularized kernel function classes under the assumption of  $\ell_\infty$  Lipschitz and  $\sigma$ -strongly convex loss. Here Theorem 3.17 shows a  $\tilde{O}(1/\sqrt{n})$  bound for general function classes with  $\ell_\infty$  Lipschitz loss in the decoupling case. Yu et al. (2014) obtained a  $O(1/\sqrt{n})$  bound for trace norm regularized linear function classes with the decomposable loss. In addition, Xu et al. (2016) used the local Rademacher complexity to derive a  $\tilde{O}(1/n)$  bound for trace norm regularized linear function classes with the assumption that the singular values of  $w$  decay exponentially. Compared with these works, we obtain tighter bounds with the state-of-the-art dependency on  $c$  for general function classes under the mild assumptions in both decoupling and coupling cases.

## 6. Discussion

The main goal of our capacity-based generalization bounds is to provide general and efficient theoretical guarantees for empirically successful multi-label learning methods, especially regarding the dependency on the number of labels. The generalization differences of different multi-label learning models or algorithms are mainly reflected in two aspects. On the one hand, the differences are reflected in the Lipschitz constant of the loss functions, as we showed in Proposition 3.3 and 3.10, different loss functions have different  $\mu$  values in our bounds. On the other hand, the differences are reflected in the nonlinear mappings corresponding to the specific models used. In fact, when we analyze the generalization of the specific models, the constraint on nonlinear mappings  $\beta(\phi(\mathbf{x})) \leq A$  is actually refined as  $\|\phi(\mathbf{x})\| \leq A$  in our analysis, and we will further have  $\|\phi(\mathbf{x})\| \leq \rho\|\mathbf{x}\|$  ( $\rho$  is the Lipschitz constant of the nonlinear mappings) to take into account the differences or characteristics of different models. The generalization differences are further reflected in the corresponding Lipschitz constants  $\rho$ . Compared with the Lipschitz constants of deep models and shallow models, their differences are particularly obvious (Bartlett et al., 2017; Golowich et al., 2018; Bartlett et al., 2019; Zhang & Liao, 2020; Ledent et al., 2021; Zhang & Zhang, 2023). In order to provide general theoretical guarantees for multi-label learning, here we only make the most general assumptions (e.g., for nonlinear mappings) and do not specify specific models, so the generalization differences of different multi-label learning models or algorithms are not explicitly shown. And if a specific model or algorithm is specified or refined, the results

obtained will lose their generality and will not be able to provide theoretical guarantees for all or most multi-label learning models or algorithms.

The analysis of the lower bound will greatly promote our theoretical understanding of multi-label learning and is the most powerful criterion for testing whether the given upper bound is the tightest. However, the effective lower bound remains relatively under-explored. We believe that such a lower bound  $\Omega(1/\sqrt{n})$  is still not tight enough. The main evidence is as follows: if we regard the classification problem corresponding to each label as a task, then multi-label learning can be regarded as multi-task learning. Since tasks share some constraint or generative properties with each other, the typical bound for multi-task learning is  $O(1/\sqrt{nt})$ , where  $t$  is the number of tasks. Hence, if we consider that the components share some constraint properties or information with each other, i.e., the label correlations, then appropriate label correlations will improve the generalization performance of multi-label learning. This means that some constraint properties shared between labels (i.e., label correlations) will facilitate the learning effect of multiple labels. In such a case, we should have the lower bound  $\Omega(1/\sqrt{nc})$ . A  $\tilde{O}(1/\sqrt{nc})$  bound for  $\ell_2$  Lipschitz surrogate Hamming loss in the decoupling case (as we discussed in Remark 3.8) provides evidence for this analysis. Hence, although our bounds are tighter than the state of the art, the above analysis actually motivates us to develop new theories to investigate the lower bound and break through the theoretical limitations of current methodologies.

## 7. Conclusion

In this paper, we propose several novel vector-contraction inequalities and derive bounds with a weaker dependency on  $c$ , and study the impact of different label correlations on the generalization analysis. In addition, with the label-based ranking multi-label Rademacher complexity, we derive the bound for Macro-Averaged AUC and analyze the relationship between Macro-Averaged AUC and class-imbalance.

In future work, we will extend our bounds to more general settings (especially for label correlations induced by label-specific features and other norm regularizers), and derive tighter bounds for multi-label learning with a faster convergence rate with respect to the number of observations, and further design efficient algorithms to construct multi-label models with good generalization performance.

## Acknowledgements

The authors wish to thank the anonymous reviewers for their helpful comments and suggestions. This work was supported by the National Science Foundation of China (62225602), the Fundamental Research Funds for the Cen-

tral Universities (2242024K30035), and the Big Data Computing Center of Southeast University.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

- Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. Spectrally-normalized margin bounds for neural networks. *Advances in Neural Information Processing Systems*, 30 (NIPS 2017):6240–6249, 2017.
- Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17, 2019.
- Barutcuoglu, Z., Schapire, R. E., and Troyanskaya, O. G. Hierarchical multi-label prediction of gene function. *Bioinformatics*, 22(7):830–836, 2006.
- Boutell, M. R., Luo, J., Shen, X., and Brown, C. M. Learning multi-label scene classification. *Pattern Recognition*, 37(9):1757–1771, 2004.
- Cabral, R. S., la Torre, F. D., Costeira, J. P., and Bernardino, A. Matrix completion for multi-label image classification. *Advances in Neural Information Processing Systems*, 24 (NIPS 2011):190–198, 2011.
- Cesa-Bianchi, N., Re, M., and Valentini, G. Synergy of multi-label hierarchical ensembles, data fusion, and cost-sensitive methods for gene functional inference. *Machine Learning*, 88(1-2):209–241, 2012.
- Cléménçon, S., Lugosi, G., and Vayatis, N. Ranking and empirical minimization of u-statistics. *The Annals of Statistics*, 36(2):844–874, 2008.
- Dembczynski, K., Waegeman, W., Cheng, W., and Hüllermeier, E. Regret analysis for performance metrics in multi-label classification: The case of hamming and subset zero-one loss. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, volume 6321, pp. 280–295, 2010.
- Dembczynski, K., Kotlowski, W., and Hüllermeier, E. Consistent multilabel ranking through univariate losses. *arXiv:1206.6401*, 2012.
- Foster, D. J. and Rakhlin, A.  $\ell_\infty$  vector contraction for rademacher complexity. *arXiv:1911.06468v1*, 2019.
- Golowich, N., Rakhlin, A., and Shamir, O. Size-independent sample complexity of neural networks. *International Conference on Computational Learning Theory*, 75(COLT 2018):297–299, 2018.
- Hang, J.-Y. and Zhang, M.-L. Collaborative learning of label semantics and deep label-specific features for multi-label classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):9860–9871, 2022.
- Huang, J., Li, G., Huang, Q., and Wu, X. Learning label-specific features and class-dependent labels for multi-label classification. *IEEE Transactions on Knowledge and Data Engineering*, 28(12):3309–3323, 2016.
- Jia, B., Liu, J., Hang, J., and Zhang, M. Learning label-specific features for decomposition-based multi-class classification. *Frontiers of Computer Science*, 17(6):176348, 2023.
- Koltchinskii, V. and Panchenko, D. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1):1–50, 2002.
- Ledent, A., Mustafa, W., Lei, Y., and Kloft, M. Norm-based generalisation bounds for deep multi-class convolutional neural networks. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, volume 35, pp. 8279–8287, 2021.
- Lei, Y., Dogan, Ü., Binder, A., and Kloft, M. Multi-class svms: From tighter data-dependent generalization bounds to novel algorithms. In *Advances in Neural Information Processing Systems*, volume 28, pp. 2035–2043, 2015.
- Lei, Y., Dogan, Ü., Zhou, D., and Kloft, M. Data-dependent generalization bounds for multi-class classification. *IEEE Transactions on Information Theory*, 65(5):2995–3021, 2019.
- Li, J., Liu, Y., Yin, R., Zhang, H., Ding, L., and Wang, W. Multi-class learning: From theory to algorithm. *Advances in Neural Information Processing Systems*, 31(NeurIPS 2018):1593–1602, 2018.
- Liu, C., Zhao, P., Huang, S., Jiang, Y., and Zhou, Z. Dual set multi-label learning. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, number AAAI 2018, pp. 3635–3642, 2018.
- Liu, W. and Shen, X. Sparse extreme multi-label learning with oracle property. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 4032–4041, 2019.

- Lust-Piquard, F. and Pisier, G. Non commutative khintchine and paley inequalities. *Arkiv för matematik*, 29:241–260, 1991.
- Maurer, A. A vector-contraction inequality for rademacher complexities. In *Proceedings of the 27th International Conference on Algorithmic Learning Theory*, volume 9925, pp. 3–17, 2016.
- Prabhu, Y. and Varma, M. Fastxml: a fast, accurate and stable tree-classifier for extreme multi-label learning. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, number KDD 2014, pp. 263–272, 2014.
- Rubin, T. N., Chambers, A., Smyth, P., and Steyvers, M. Statistical topic models for multi-label document classification. *Machine Learning*, 88(1-2):157–208, 2012.
- Rudelson, M. and Vershynin, R. Combinatorics of random processes and sections of convex bodies. *Annals of Mathematics*, pp. 603–648, 2006.
- Schapire, R. E. and Singer, Y. Boostexter: A boosting-based system for text categorization. *Machine Learning*, 39(2/3):135–168, 2000.
- Srebro, N., Sridharan, K., and Tewari, A. Smoothness, low noise and fast rates. In *Advances in Neural Information Processing Systems*, volume 23, pp. 2199–2207, 2010.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Wu, G. and Zhu, J. Multi-label classification: do hamming loss and subset accuracy really conflict with each other? *Advances in Neural Information Processing Systems*, 33 (NeurIPS 2020), 2020.
- Wu, G., Li, C., Xu, K., and Zhu, J. Rethinking and reweighting the univariate losses for multi-label ranking: Consistency and generalization. *Advances in Neural Information Processing Systems*, 34(NeurIPS 2021):14332–14344, 2021a.
- Wu, G., Li, C., and Yin, Y. Towards understanding generalization of macro-auc in multi-label learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 37540–37570, 2023.
- Wu, L., Ledent, A., Lei, Y., and Kloft, M. Fine-grained generalization analysis of vector-valued learning. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, number AAAI 2021, pp. 10338–10346, 2021b.
- Xu, C., Liu, T., Tao, D., and Xu, C. Local rademacher complexity for multi-label learning. *IEEE Transactions on Image Processing*, 25(3):1495–1507, 2016.
- Yen, I. E., Huang, X., Ravikumar, P., Zhong, K., and Dhillon, I. S. Pd-sparse : A primal and dual sparse approach to extreme multiclass and multilabel classification. In *Proceedings of the 33rd International Conference on Machine Learning*, volume 48, pp. 3069–3077, 2016.
- Yu, H., Jain, P., Kar, P., and Dhillon, I. S. Large-scale multi-label learning with missing labels. In *Proceedings of the 31th International Conference on Machine Learning*, volume 32, pp. 593–601, 2014.
- Yu, K., Yu, S., and Tresp, V. Multi-label informed latent semantic indexing. In Baeza-Yates, R. A., Ziviani, N., Marchionini, G., Moffat, A., and Tait, J. (eds.), *Proceedings of the 28th Annual International Conference on Research and Development in Information Retrieval*, number SIGIR 2005, pp. 258–265, 2005.
- Zhang, M.-L. and Wu, L. Lift: Multi-label learning with label-specific features. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(1):107–120, 2015.
- Zhang, M.-L. and Zhou, Z.-H. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 26(8):1819–1837, 2014.
- Zhang, Y. and Liao, S. A kernel perspective for the decision boundary of deep neural networks. In *Proceedings of the 32nd IEEE International Conference on Tools with Artificial Intelligence*, number ICTAI 2020, pp. 653–660, 2020.
- Zhang, Y. and Zhang, M. Nearly-tight bounds for deep kernel learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 41861–41879, 2023.

## A. Appendix Outline

In the appendix, we give the detailed proofs of those theoretical results which we only provide proof sketches in the main paper. Our main proofs include:

- The basic generalization bound for  $\ell_2$  Lipschitz loss (Theorem 3.4).
- The novel vector-contraction inequality for  $\ell_2$  Lipschitz loss (Lemma 3.6).
- The generalization bound with a square-root dependency on  $c$  for  $\ell_2$  Lipschitz loss (Theorem 3.7).
- The  $\ell_\infty$  Lipschitz continuity of the commonly used loss for multi-label learning (Proposition 3.10).
- The novel vector-contraction inequality for  $\ell_\infty$  Lipschitz loss with a square-root dependency on  $c$  (Lemma 3.11).
- The generalization bound with a square-root dependency on  $c$  for  $\ell_\infty$  Lipschitz loss (Theorem 3.12).
- The generalization bounds with no dependency on  $c$  for  $\ell_\infty$  Lipschitz loss in the coupling case (Corollary 3.14).
- The novel vector-contraction inequality for  $\ell_\infty$  Lipschitz loss with no dependency on  $c$  (Lemma 3.16).
- The generalization bound with no dependency on  $c$  for  $\ell_\infty$  Lipschitz loss in the decoupling case (Theorem 3.17).
- The generalization bound for Macro-Averaged AUC w.r.t.  $\ell_\infty$  Lipschitz loss (Theorem 4.2).

## B. Preliminaries

First, we define the loss function space as follows:

$$\mathcal{L} = \{L(\mathbf{f}(\mathbf{x}), \mathbf{y}) : \mathbf{f} \in \mathcal{F}\},$$

where  $\mathcal{F}$  is the vector-valued function class (1) of multi-label learning defined in the main paper:

$$\begin{aligned} \mathcal{F} &= \{\mathbf{x} \mapsto \mathbf{f}(\mathbf{x}) : \mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_c(\mathbf{x})), f_j(\mathbf{x}) = \langle \mathbf{w}_j, \phi(\mathbf{x}) \rangle, \\ &\quad \mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_c) \in \mathbb{R}^{d \times c}, \alpha(\mathbf{w}) \leq \Lambda, \beta(\phi(\mathbf{x})) \leq A, \mathbf{x} \in \mathcal{X}, j \in [c]\}, \end{aligned}$$

For any training dataset  $D = \{(\mathbf{x}_i, \mathbf{y}_i) : i \in [n]\}$ , let  $D' = \{(\mathbf{x}_i, \mathbf{y}_i) : i \in [n]\}$  be the training dataset with only one sample different from  $D$ , where the  $k$ -th sample is replaced by  $(\mathbf{x}'_k, \mathbf{y}'_k)$ . Let  $\Phi(D) = \sup_{\mathbf{f} \in \mathcal{F}} [\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} [L(\mathbf{f}(\mathbf{x}), \mathbf{y})] - \frac{1}{n} \sum_{i=1}^n L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i)] = \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \widehat{R}_D(\mathbf{f})]$ , then

$$\begin{aligned} &\Phi(D') - \Phi(D) \\ &= \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \widehat{R}_{D'}(\mathbf{f})] - \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \widehat{R}_D(\mathbf{f})] \\ &\leq \sup_{\mathbf{f} \in \mathcal{F}} [\widehat{R}_D(\mathbf{f}) - \widehat{R}_{D'}(\mathbf{f})] \\ &= \sup_{\mathbf{f} \in \mathcal{F}} \frac{[L(\mathbf{f}(\mathbf{x}_k), \mathbf{y}_k) - L(\mathbf{f}(\mathbf{x}'_k), \mathbf{y}'_k)]}{n} \\ &\leq \frac{M}{n}. \end{aligned}$$

According to McDiarmid's inequality, for any  $0 < \delta < 1$ , with probability at least  $1 - \frac{\delta}{2}$  over the training dataset  $D$ , the following holds:

$$\Phi(D) \leq \mathbb{E}_D[\Phi(D)] + M \sqrt{\frac{\ln(2/\delta)}{2n}}. \quad (4)$$

Then, we will estimate the upper bound of  $\mathbb{E}_D[\Phi(D)]$ .

$$\begin{aligned}
 & \mathbb{E}_D[\Phi(D)] \\
 &= \mathbb{E}_D \left[ \sup_{\mathbf{f} \in \mathcal{F}} \left[ \mathbb{E}_{D'} \left[ \widehat{R}_{D'}(\mathbf{f}) - \widehat{R}_D(\mathbf{f}) \right] \right] \right] \\
 &\leq \mathbb{E}_{D, D'} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \left[ \widehat{R}_{D'}(\mathbf{f}) - \widehat{R}_D(\mathbf{f}) \right] \right] \\
 &= \mathbb{E}_{D, D'} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \left[ \sum_{i=1}^n L(\mathbf{f}(\mathbf{x}'_i), \mathbf{y}'_i) - L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right] \right] \\
 &= \mathbb{E}_{\epsilon, D, D'} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \left[ \sum_{i=1}^n \epsilon_i (L(\mathbf{f}(\mathbf{x}'_i), \mathbf{y}'_i)) - L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right] \right] \\
 &\leq \mathbb{E}_{\epsilon, D'} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i L(\mathbf{f}(\mathbf{x}'_i), \mathbf{y}'_i) \right] \\
 &+ \mathbb{E}_{\epsilon, D} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n -\epsilon_i L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right] \\
 &= 2\mathbb{E}_{\epsilon, D} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right]. \tag{5}
 \end{aligned}$$

Then apply McDiarmid's inequality to  $\mathbb{E}_{\epsilon} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right]$ , we have

$$\mathbb{E}_{\epsilon, D} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right] \leq \mathbb{E}_{\epsilon} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) \right] + M \sqrt{\frac{\ln(2/\delta)}{2n}},$$

i.e.,

$$\mathfrak{R}(\mathcal{L}) \leq \widehat{\mathfrak{R}}_D(\mathcal{L}) + M \sqrt{\frac{\ln(2/\delta)}{2n}}. \tag{6}$$

Combining with (4), (5) and (6), then

$$R(\mathbf{f}) \leq \widehat{R}_D(\mathbf{f}) + 2\widehat{\mathfrak{R}}_D(\mathcal{L}) + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}. \tag{7}$$

## C. Generalization Bounds for $\ell_2$ Lipschitz Loss

### C.1. Proof of Theorem 3.4

We first introduce the following lemmas:

**Lemma C.1** (Corollary 1 in (Maurer, 2016)). *Let  $\mathcal{F}$  be a vector-valued function class. Given a dataset  $D$  of size  $n$ . Assume that the loss function is  $\mu$ -Lipschitz continuous with respect to the  $\ell_2$  norm. Then*

$$\mathbb{E}_{\epsilon} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i L(\mathbf{f}(\mathbf{x}_i)) \right] \leq \sqrt{2}\mu \mathbb{E}_{\epsilon} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c \epsilon_{ij} f_j(\mathbf{x}_i) \right],$$

where each  $\epsilon_{ij}$  is an independent doubly indexed Rademacher random variable, and  $f_j(\mathbf{x}_i)$  is the  $j$ -th component of  $\mathbf{f}(\mathbf{x}_i)$ .

**Lemma C.2** (Khintchine-Kahane inequality (Lust-Piquard & Pisier, 1991)). *Let  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathcal{H}$ , where  $\mathcal{H}$  is a Hilbert space with  $\|\cdot\|$  being the associated  $p$ -th norm. Let  $\epsilon_1, \dots, \epsilon_n$  be a sequence of independent Rademacher variables. Then,*

for any  $p \geq 1$  there holds

$$\min(\sqrt{p-1}, 1) \left[ \sum_{i=1}^n \|\mathbf{v}_i\|^2 \right]^{\frac{1}{2}} \leq \left[ \mathbb{E}_\epsilon \left\| \sum_{i=1}^n \epsilon_i \mathbf{v}_i \right\|^p \right]^{\frac{1}{p}} \leq \max(\sqrt{p-1}, 1) \left[ \sum_{i=1}^n \|\mathbf{v}_i\|^2 \right]^{\frac{1}{2}},$$

and

$$\mathbb{E}_\epsilon \left\| \sum_{i=1}^n \epsilon_i \mathbf{v}_i \right\| \geq 2^{-\frac{1}{2}} \left[ \sum_{i=1}^n \|\mathbf{v}_i\|^2 \right]^{\frac{1}{2}}.$$

According to Lemma C.1, we then have

$$\begin{aligned} & \hat{\mathfrak{R}}_D(\mathcal{L}) \\ & \leq \sqrt{2\mu} \mathbb{E}_\epsilon \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^c \epsilon_{ij} f_j(\mathbf{x}_i) \right] \\ & \leq \sqrt{2\mu} c \max_j \mathbb{E}_\epsilon \left[ \sup_{f_j \in \mathcal{F}_j} \frac{1}{n} \sum_{i=1}^n \epsilon_{ij} f_j(\mathbf{x}_i) \right] \\ & \leq \sqrt{2\mu} c \max_j \sup_{f_j \in \mathcal{F}_j} \frac{1}{n} \left( \sum_{i=1}^n (f_j(\mathbf{x}_i))^2 \right)^{\frac{1}{2}} \quad (\text{Use Lemma C.2}) \\ & \leq \frac{\sqrt{2\mu} c B}{\sqrt{n}}. \quad (\text{Use Assumption 3.1 in the main paper}) \end{aligned}$$

Combining with (7), then

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{2\sqrt{2\mu} c B}{\sqrt{n}} + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

## C.2. Proof of Lemma 3.6

**Proof Sketch:** First, the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded by the empirical  $\ell_2$  norm covering number with the refined Dudley's entropy integral inequality. Second, according to the Lipschitz continuity w.r.t the  $\ell_2$  norm, the empirical  $\ell_2$  norm covering number of  $\mathcal{F}$  can be bounded by the empirical  $\ell_2$  norm covering number of  $\mathcal{P}(\mathcal{F})$ . Third, the empirical  $\ell_2$  norm covering number of  $\mathcal{P}(\mathcal{F})$  can be bounded by using Sudakov's minoration (Wainwright, 2019), which bounds the  $\ell_2$  norm covering number of a function class by the expectation of a Gaussian process indexed by the function class, and the expectation of the Gaussian process can be bounded by the worst-case Rademacher complexity of the projection function class  $\mathcal{P}(\mathcal{F})$ . Hence, the problem is transferred to the estimation of the worst-case Rademacher complexity. Finally, we estimate the lower bound of the worst-case Rademacher complexity of  $\mathcal{P}(\mathcal{F})$ , and then combined with the above steps, the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded.

We first introduce the following lemmas:

**Lemma C.3** (Sudakov's minoration (Wainwright, 2019)). *Let  $\{Z_f, f \in \mathcal{F}\}$  be a zero-mean Gaussian process. Then*

$$\mathbb{E} \left[ \sup_{f \in \mathcal{F}} Z_f \right] \geq \sup_{\epsilon > 0} \frac{\epsilon}{2} \sqrt{\log \mathcal{M}_2(\epsilon, \mathcal{F}, D)} \geq \sup_{\epsilon > 0} \frac{\epsilon}{2} \sqrt{\log \mathcal{N}_2(\epsilon, \mathcal{F}, D)}.$$

**Lemma C.4** (Relationship between Rademacher and Gaussian complexity (Wainwright, 2019)). *Let  $\mathcal{G}$  be a class of real-valued functions mapping from  $\mathcal{X}$  to  $\mathbb{R}$ . Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set with  $n$  i.i.d. samples. Then*

$$\sqrt{\frac{2}{\pi}} \hat{\mathfrak{R}}_D(\mathcal{G}) \leq \hat{\mathfrak{G}}_D(\mathcal{G}) \leq 2 \hat{\mathfrak{R}}_D(\mathcal{G}) \sqrt{\log n},$$

where  $\hat{\mathfrak{G}}_D(\mathcal{G}) = \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{x}_i) \right]$ , and  $\sigma_1, \dots, \sigma_n$  are i.i.d. random variables obeying the normal distribution.

The following lemma is a refined result of Proposition 22 in (Ledent et al., 2021), where we replace the function class taking values in  $[0, 1]$  with the  $b$ -bounded function class, the refinement is obvious.

**Lemma C.5** (Refined Dudley's entropy integral inequality). *Let  $\mathcal{F}$  be a real-valued function class with  $f \leq b$ ,  $f \in \mathcal{F}$ ,  $b > 0$ , and assume that  $0 \in \mathcal{F}$ . Let  $S$  be a finite sample of size  $n$ . For any  $2 \leq p \leq \infty$ , we have the following relationship between the Rademacher complexity  $\hat{\mathfrak{R}}_S(\mathcal{F})$  and the covering number  $\mathcal{N}_p(\epsilon, \mathcal{F}, S)$ .*

$$\hat{\mathfrak{R}}_S(\mathcal{F}) \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^b \sqrt{\log \mathcal{N}_p(\epsilon, \mathcal{F}, S)} d\epsilon \right).$$

**Step 1:** We first derive the relationship between the empirical  $\ell_2$  norm covering number  $\mathcal{N}_2(\epsilon, \mathcal{L}, D)$  and the empirical  $\ell_2$  norm covering number  $\mathcal{N}_2(\epsilon, \mathcal{P}(\mathcal{F}), [c] \times D)$ .

For the dataset  $D = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$  with  $n$  i.i.d. examples:

$$\begin{aligned} & \sqrt{\frac{1}{n} \sum_{i=1}^n (L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) - L(\mathbf{f}'(\mathbf{x}_i), \mathbf{y}_i))^2} \\ & \leq \sqrt{\frac{1}{n} \sum_{i=1}^n \mu^2 \|\mathbf{f}(\mathbf{x}_i) - \mathbf{f}'(\mathbf{x}_i)\|_2^2} \quad (\text{Use Assumption 3.2}) \\ & \leq \mu \sqrt{c \frac{1}{n} \sum_{i=1}^n \frac{1}{c} \sum_{j=1}^c (f_j(\mathbf{x}_i) - f'_j(\mathbf{x}_i))^2} \\ & \leq \mu \sqrt{c} \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{c} \sum_{j=1}^c (p_j(\mathbf{f}(\mathbf{x}_i)) - p_j(\mathbf{f}'(\mathbf{x}_i)))^2}. \quad (\text{The definition of the projection function class } \mathcal{P}(\mathcal{F})) \end{aligned}$$

Then, according to the definition of the empirical  $\ell_2$  covering number, we have that an empirical  $\ell_2$  cover of  $\mathcal{P}(\mathcal{F})$  at radius  $\epsilon/\mu\sqrt{c}$  is also an empirical  $\ell_2$  cover of the loss function space associated with  $\mathcal{F}$  at radius  $\epsilon$ , and we can conclude that:

$$\mathcal{N}_2(\epsilon, \mathcal{L}, D) \leq \mathcal{N}_2\left(\frac{\epsilon}{\mu\sqrt{c}}, \mathcal{P}(\mathcal{F}), [c] \times D\right). \quad (8)$$

**Step 2:** We show that the empirical  $\ell_2$  norm covering number of  $\mathcal{P}(\mathcal{F})$  can be bounded by using Sudakov's minoration (Wainwright, 2019), which bounds the  $\ell_2$  norm covering number of a function class by the expectation of a Gaussian process indexed by the function class, and the expectation of the Gaussian process can be bounded by the worst-case Rademacher complexity of the projection function class  $\mathcal{P}(\mathcal{F})$ .

Let  $\mathcal{G}$  be a class of real-valued functions. Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a set with  $n$  i.i.d. samples. We consider the Gaussian process  $Z_f = \frac{1}{\sqrt{n}} \sum_{i=1}^n \sigma_i f(\mathbf{x}_i)$ , then according to Lemma C.3, we have  $\mathbb{E}[\sup_{f \in \mathcal{F}} Z_f] = \sqrt{n} \hat{\mathfrak{G}}_D(\mathcal{F})$ . Then, combining Lemma C.3 and Lemma C.4, we can get  $\sqrt{\log \mathcal{N}_2(\epsilon, \mathcal{F}, D)} \leq \frac{4}{\epsilon} \hat{\mathfrak{R}}_D(\mathcal{F}) \sqrt{n \log n}$ .

Hence, for the projection function class  $\mathcal{P}(\mathcal{F})$ , we have

$$\sqrt{\log \mathcal{N}_2(\epsilon, \mathcal{P}(\mathcal{F}), [c] \times D)} \leq \frac{4}{\epsilon} \hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F})) \sqrt{n \log n}. \quad (9)$$

**Step 3:** According to Assumption 3.1 in the main paper, we can obtain the lower bound of the worst-case Rademacher complexity  $\hat{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))$  by the Khintchine-Kahane inequality with  $p = 1$ :

$$\begin{aligned} \hat{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) &= \sup_{[c] \times D \in [c] \times \mathcal{X}^n} \hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F})) = \sup_{[c] \times D \in [c] \times \mathcal{X}^n} \mathbb{E}_{\epsilon} \left[ \sup_{p_j(\mathbf{f}(\mathbf{x}_i)) \in \mathcal{P}(\mathcal{F})} \frac{1}{nc} \sum_{i=1}^n \sum_{j=1}^c \epsilon_i p_j(\mathbf{f}(\mathbf{x}_i)) \right] \\ &= \sup_{D \in \mathcal{X}^n} \mathbb{E}_{\epsilon} \left[ \sup_{f_j \in \mathcal{F}_j} \frac{1}{nc} \sum_{i=1}^n \sum_{j=1}^c \epsilon_i f_j(\mathbf{x}_i) \right] \geq \sup_{D \in \mathcal{X}^n} \frac{1}{nc} \sup_{f_j \in \mathcal{F}_j} \frac{1}{\sqrt{2}} \left[ \sum_{i=1}^n \sum_{j=1}^c (f_j(\mathbf{x}_i))^2 \right]^{\frac{1}{2}}. \end{aligned}$$

Since  $|f_j(\cdot)| \leq B$ , we set  $\sup_{D \in \mathcal{X}^n} \frac{1}{nc} \sup_{f_j \in \mathcal{F}_j} \left[ \sum_{i=1}^n \sum_{j=1}^c (f_j(\mathbf{x}_i))^2 \right]^{\frac{1}{2}} = \frac{B}{\sqrt{nc}}$ . So,

$$\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \geq \frac{B}{\sqrt{2nc}}. \quad (10)$$

Then, according to Lemma C.5 and combined with the above steps, we have

$$\begin{aligned} & \hat{\mathfrak{R}}_D(\mathcal{L}) \\ & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^M \sqrt{\log \mathcal{N}_2(\epsilon, \mathcal{L}, D)} d\epsilon \right) \\ & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_{\alpha}^M \sqrt{\log \mathcal{N}_2\left(\frac{\epsilon}{\mu\sqrt{c}}, \mathcal{P}(\mathcal{F}), [c] \times D\right)} d\epsilon \right) \quad (\text{Use inequality (8)}) \\ & \leq \inf_{\alpha > 0} \left( 4\alpha + 48\sqrt{c}\mu\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \log^{\frac{1}{2}}(nc) \int_{\alpha}^M \epsilon^{-1} d\epsilon \right) \\ & \quad (\text{Use inequality (9) and the definition of the worst-case Rademacher complexity}) \\ & \leq 48\mu\sqrt{c}\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) + 48\mu\sqrt{c}\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \log^{\frac{1}{2}}(nc) \cdot \log \frac{M}{12\sqrt{c}\mu\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))} \quad (\text{Choose } \alpha = 12\sqrt{c}\mu\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))) \\ & \leq 48\mu\sqrt{c}\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \left( 1 + \log^{\frac{1}{2}}(nc) \cdot \log \frac{M\sqrt{n}}{\mu B} \right). \quad (\text{Use inequality (10)}) \end{aligned}$$

### C.3. Proof of Theorem 3.7

We upper bound the worst-case Rademacher complexity  $\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))$  as the following:

$$\begin{aligned} & \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \\ & = \sup_{[c] \times D \in [c] \times \mathcal{X}^n} \hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F})) \\ & = \sup_{[c] \times D \in [c] \times \mathcal{X}^n} \mathbb{E}_{\epsilon} \left[ \sup_{p_j(\mathbf{f}(\mathbf{x}_i)) \in \mathcal{P}(\mathcal{F})} \frac{1}{nc} \sum_{i=1}^n \sum_{j=1}^c \epsilon_i p_j(\mathbf{f}(\mathbf{x}_i)) \right] \\ & = \sup_{D \in \mathcal{X}^n} \mathbb{E}_{\epsilon} \left[ \sup_{f_j \in \mathcal{F}_j} \frac{1}{nc} \sum_{i=1}^n \sum_{j=1}^c \epsilon_i f_j(\mathbf{x}_i) \right] \\ & \leq \sup_{D \in \mathcal{X}^n} \sup_{f_j \in \mathcal{F}_j} \frac{1}{nc} \left( \sum_{i=1}^n \sum_{j=1}^c (f_j(\mathbf{x}_i))^2 \right)^{\frac{1}{2}} \quad (\text{Use Lemma C.2}) \\ & \leq \frac{B}{\sqrt{nc}}. \quad (\text{Use Assumption 3.1 in the main paper}) \end{aligned} \quad (11)$$

Then, we have

$$\begin{aligned} \hat{\mathfrak{R}}_D(\mathcal{F}) & \leq 48\mu\sqrt{c}\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \left( 1 + \log^{\frac{1}{2}}(nc) \cdot \log \frac{M\sqrt{n}}{\mu B} \right) \\ & \leq \frac{48B\mu \left( 1 + \log^{\frac{1}{2}}(nc) \cdot \log \frac{M\sqrt{n}}{\mu B} \right)}{\sqrt{n}}. \end{aligned}$$

Combining with (7), then

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{96B\mu \left( 1 + \log^{\frac{1}{2}}(nc) \cdot \log \frac{M\sqrt{n}}{\mu B} \right)}{\sqrt{n}} + 3M\sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$



## D. Generalization Bounds for $\ell_\infty$ Lipschitz Loss

### D.1. Proof of Proposition 3.10

We first prove that the surrogate Hamming Loss is  $\mu$ -Lipschitz continuous with respect to the  $\ell_\infty$  norm.

$$\begin{aligned}
 & |L_H(\mathbf{f}(\mathbf{x}), \mathbf{y}) - L_H(\mathbf{f}'(\mathbf{x}), \mathbf{y})| \\
 &= \left| \frac{1}{c} \sum_{j=1}^c \ell(y_j f_j(\mathbf{x})) - \frac{1}{c} \sum_{j=1}^c \ell(y_j f'_j(\mathbf{x})) \right| \\
 &= \frac{1}{c} \sum_{j=1}^c |\ell(y_j f_j(\mathbf{x})) - \ell(y_j f'_j(\mathbf{x}))| \\
 &\leq \frac{1}{c} \sum_{j=1}^c \mu |f_j(\mathbf{x}) - f'_j(\mathbf{x})| \\
 &\leq \frac{1}{c} \mu c \max_{j \in [c]} |f_j(\mathbf{x}) - f'_j(\mathbf{x})| \\
 &= \mu \|\mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\|_\infty.
 \end{aligned}$$

Then, with the elementary inequality

$$|\max\{a_1, \dots, a_c\} - \max\{b_1, \dots, b_c\}| \leq \max\{|a_1 - b_1|, \dots, |a_c - b_c|\},$$

we proof that the surrogate Subset Loss is  $\mu$ -Lipschitz continuous with respect to the  $\ell_\infty$  norm.

$$\begin{aligned}
 & |L_S(\mathbf{f}(\mathbf{x}), \mathbf{y}) - L_S(\mathbf{f}'(\mathbf{x}), \mathbf{y})| \\
 &= \left| \max_{j \in [c]} \ell(y_j f_j(\mathbf{x})) - \max_{j \in [c]} \ell(y_j f'_j(\mathbf{x})) \right| \\
 &\leq \max_{j \in [c]} |\ell(y_j f_j(\mathbf{x})) - \ell(y_j f'_j(\mathbf{x}))| \\
 &\leq \mu \max_{j \in [c]} |f_j(\mathbf{x}) - f'_j(\mathbf{x})| \\
 &= \mu \|\mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\|_\infty.
 \end{aligned}$$

Finally, we proof that the surrogate Ranking Loss is  $2\mu$ -Lipschitz continuous with respect to the  $\ell_\infty$  norm.

$$\begin{aligned}
 & |L_R(\mathbf{f}(\mathbf{x}), \mathbf{y}) - L_R(\mathbf{f}'(\mathbf{x}), \mathbf{y})| \\
 &= \frac{1}{|Y^+| |Y^-|} \left| \sum_{p \in Y^+} \sum_{q \in Y^-} (\ell(f_p(\mathbf{x}) - f_q(\mathbf{x})) - \ell(f'_p(\mathbf{x}) - f'_q(\mathbf{x}))) \right| \\
 &\leq \max_{p \in Y^+, q \in Y^-} |\ell(f_p(\mathbf{x}) - f_q(\mathbf{x})) - \ell(f'_p(\mathbf{x}) - f'_q(\mathbf{x}))| \\
 &\leq \mu \max_{p \in Y^+, q \in Y^-} |(f_p(\mathbf{x}) - f_q(\mathbf{x})) - (f'_p(\mathbf{x}) - f'_q(\mathbf{x}))| \\
 &\leq \mu \left( \max_{p \in Y^+} |f_p(\mathbf{x}) - f'_p(\mathbf{x})| + \max_{q \in Y^-} |f_q(\mathbf{x}) - f'_q(\mathbf{x})| \right) \\
 &\leq 2\mu \max_{j \in [c]} |f_j(\mathbf{x}) - f'_j(\mathbf{x})| \\
 &= 2\mu \|\mathbf{f}(\mathbf{x}) - \mathbf{f}'(\mathbf{x})\|_\infty.
 \end{aligned}$$

### D.2. Proof of Lemma 3.11

We first introduce the following vector-contraction inequality in (Foster & Rakhlin, 2019):

**Lemma D.1** (Theorem 1 in (Foster & Rakhlin, 2019)). *Let  $\mathcal{F} \subseteq \{\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^c\}$ , and let  $\phi : \mathbb{R}^c \rightarrow \mathbb{R}$  be  $L$ -lipschitz with respect to the  $\ell_\infty$  norm. For any  $1 > \eta > 0$ , there exists a constant  $C > 0$  such that if  $|\phi(\mathbf{f}(x))| \vee \|\mathbf{f}(x)\|_\infty \leq \beta$ , then*

$$\hat{\mathfrak{R}}_D(\phi \circ \mathcal{F}) \leq CL\sqrt{c} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \log^{\frac{3+\eta}{2}} \left( \frac{\beta n}{\max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j)} \right),$$

where  $\hat{\mathfrak{R}}_D(\phi \circ \mathcal{F}) = \mathbb{E}_\epsilon \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i \phi(\mathbf{f}(\mathbf{x}_i)) \right]$ ,  $\tilde{\mathfrak{R}}_n(\mathcal{F}_j) = \sup_{D \in \mathcal{X}^n} \hat{\mathfrak{R}}_D(\mathcal{F}_j)$ .

Next we will refine the proof of Lemma D.1. First, we introduce the result in Theorem 1 in (Foster & Rakhlin, 2019):  $\log \mathcal{N}_\infty(\epsilon, \phi \circ \mathcal{F}, D) \leq \max_i EcLr_i \log^{1+\eta} \frac{e^{2+\eta n}}{r_i \epsilon_0}$ , where  $r_i = 8n \left( \frac{\tilde{\mathfrak{R}}_n(\mathcal{F}_j)}{a\epsilon_0} \right)^2$ ,  $E > 0$ ,  $0 < a < 1$ ,  $\epsilon_0 = \frac{\epsilon}{L}$ .

Then, by direct calculation and proper scaling, we have

$$\begin{aligned} & \log \mathcal{N}_2(\epsilon, \phi \circ \mathcal{F}, D) \\ & \leq \max_j Ec r_i \log^{1+\eta} \frac{e^{2+\eta n}}{r_i \epsilon_0} \\ & \leq \max_j Ec 8n \left( \frac{\tilde{\mathfrak{R}}_n(\mathcal{F}_j)}{a\epsilon_0} \right)^2 \log^{1+\eta} \frac{4}{(\tilde{\mathfrak{R}}_n(\mathcal{F}_j))^2} \quad (\text{Use inequality } 2\tilde{\mathfrak{R}}_n(\mathcal{F}_j) < \epsilon_0 < 1) \\ & \leq \max_j Ec 32n \left( \frac{\tilde{\mathfrak{R}}_n(\mathcal{F}_j)}{a\epsilon_0} \right)^2 \log^{1+\eta} \frac{2}{\tilde{\mathfrak{R}}_n(\mathcal{F}_j)} \\ & \leq \max_j Ec 32n \left( \frac{\tilde{\mathfrak{R}}_n(\mathcal{F}_j)}{a\epsilon_0} \right)^2 \log^{1+\eta} \frac{2\sqrt{2n}}{B} \end{aligned} \tag{12}$$

(Use the similar technique to the proof of Lemma 3.6, the lower bound of  $\tilde{\mathfrak{R}}_n(\mathcal{F}_j) \geq \frac{B}{\sqrt{2n}}$ )

According to Lemma C.5 and combined with (12), we have

$$\begin{aligned} & \hat{\mathfrak{R}}_D(\phi \circ \mathcal{F}) \\ & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^M \sqrt{\log \mathcal{N}_2(\epsilon, \phi \circ \mathcal{F}, D)} d\epsilon \right) \\ & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^M \sqrt{\max_j Ec 32n \left( \frac{\tilde{\mathfrak{R}}_n(\mathcal{F}_j)}{a\epsilon_0} \right)^2 \log^{1+\eta} \frac{2\sqrt{2n}}{B}} d\epsilon \right) \quad (\text{Use inequality (12)}) \\ & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{48\sqrt{2Ec}L}{a} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \int_\alpha^M \epsilon^{-1} d\epsilon \right) \\ & \leq \frac{48\sqrt{2Ec}L}{a} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) + \frac{48\sqrt{2Ec}L}{a} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \frac{aM}{12\sqrt{2Ec}L\tilde{\mathfrak{R}}_n(\mathcal{F}_j)} \\ & \quad (\text{Choose } \alpha = \frac{12\sqrt{2Ec}L}{a} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j)) \\ & \leq \frac{48\sqrt{2Ec}L}{a} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) + \frac{48\sqrt{2Ec}L}{a} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \frac{M\sqrt{n}}{\sqrt{Ec}LB} \quad (\text{Use } 0 < a < 1 \text{ and } \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \geq \frac{B}{\sqrt{2n}}) \end{aligned}$$

Finally, combined with our problem setting, we have

$$\hat{\mathfrak{R}}_D(\mathcal{F}) \leq \frac{48\sqrt{2E}\mu}{a} \sqrt{c} \max_j \tilde{\mathfrak{R}}_n(\mathcal{F}_j) \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right).$$

### D.3. Proof of Theorem 3.12

Using the similar technique to the proof of Theorem 3.4,  $\tilde{\mathfrak{R}}_n(\mathcal{F}_j)$  can be upper bounded by  $\frac{B}{\sqrt{n}}$ , we then have

$$\hat{\mathfrak{R}}_D(\mathcal{F}) \leq \frac{48\sqrt{2E}\mu B}{a} \sqrt{\frac{c}{n}} \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right).$$

Combining with (7), then

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{96\sqrt{2E}\mu B}{a} \sqrt{\frac{c}{n}} \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right) + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

#### D.4. Proof of Corollary 3.14

We first introduce the following lemmas:

**Lemma D.2** (Theorem 5 in (Lei et al., 2019)). *Let  $\mathcal{F}$  be a vector-valued function class of the multi-label learning defined by (1). Let Assumptions 3.1 and 3.9 hold. Given a dataset  $D$  of size  $n$ . Suppose that the loss function is  $L$ -Lipschitz continuous w.r.t.  $\ell_\infty$  norm. Then the  $\hat{\mathfrak{R}}_D(\mathcal{L})$  can be bounded by*

$$\hat{\mathfrak{R}}_D(\mathcal{L}) \leq 16L\sqrt{c \log 2} \tilde{\mathfrak{R}}_{nc}(\tilde{\mathcal{F}}) \left( 1 + \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right) \right).$$

Since  $\log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right) \geq 1$  and  $16\sqrt{\log 2} < 8.8$ , we have

$$\hat{\mathfrak{R}}_D(\mathcal{L}) \leq 32L\sqrt{c \log 2} \tilde{\mathfrak{R}}_{nc}(\tilde{\mathcal{F}}) \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right) \leq 18L\sqrt{c} \tilde{\mathfrak{R}}_{nc}(\tilde{\mathcal{F}}) \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right).$$

Using the similar technique to the proof of Theorem 3.4,  $\tilde{\mathfrak{R}}_{nc}(\tilde{\mathcal{F}})$  can be upper bounded by  $\frac{\Lambda A}{\sqrt{nc}}$ . Then combining with (7) and our problem setting, we have

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{36\mu\Lambda A \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right)}{\sqrt{n}} + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

**Lemma D.3** (Corollary 10 in (Lei et al., 2019)). *Let  $\mathcal{F}$  be a vector-valued function class defined by (1), where  $\alpha(\mathbf{w}) := \|\mathbf{w}\|_{S_p}$ ,  $\beta(\phi(\mathbf{x})) := \max_{i \in [n]} \|\mathbf{x}_i\|_2$ , i.e.,  $\phi$  be the identity map, and  $p \geq 1$ . Assume that the loss function is  $L$ -Lipschitz continuous w.r.t.  $\ell_\infty$  norm. Then, if  $1 \leq p \leq 2$ , for any  $0 < \delta < 1$  with probability of  $1 - \delta$ , we have*

$$\hat{\mathfrak{R}}_D(\mathcal{L}) \leq \frac{16L\sqrt{\log 2} \Lambda \max_{i \in [n]} \|\mathbf{x}_i\|_2}{\sqrt{n}} \left( 1 + \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right) \right) \leq \frac{18L\Lambda \max_{i \in [n]} \|\mathbf{x}_i\|_2 \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right)}{\sqrt{n}}$$

If  $p > 2$ ,

$$\begin{aligned} \hat{\mathfrak{R}}_D(\mathcal{L}) &\leq \frac{16L\sqrt{\log 2} \Lambda \max_{i \in [n]} \|\mathbf{x}_i\|_2 \min\{c, d\}^{\frac{1}{2} - \frac{1}{p}}}{\sqrt{n}} \left( 1 + \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right) \right) \\ &\leq \frac{18L\Lambda \max_{i \in [n]} \|\mathbf{x}_i\|_2 \min\{c, d\}^{\frac{1}{2} - \frac{1}{p}} \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right)}{\sqrt{n}}, \end{aligned}$$

where we scale them by  $\log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right) \geq 1$  and  $16\sqrt{\log 2} < 8.8$ .

This lemma holds for the class of linear functions and is logarithmically dependent on the number of labels. We adjust the constraint on nonlinear mappings such that our bound holds for the general class of nonlinear functions. Then combining with (7) and our problem setting, when  $1 \leq p \leq 2$ , we have:

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{36\mu\Lambda A \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right)}{\sqrt{n}} + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}},$$

when  $p > 2$ , we have:

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{36\mu\Lambda A \min\{c, d\}^{\frac{1}{2} - \frac{1}{p}} \log_2^{\frac{3}{2}} \left( \sqrt{2n^{\frac{3}{2}}c} \right)}{\sqrt{n}} + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

### D.5. Proof of Lemma 3.16

**Proof Sketch:** First, the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded by the empirical  $\ell_\infty$  norm covering number with the refined Dudley's entropy integral inequality. Second, according to the Lipschitz continuity w.r.t the  $\ell_\infty$  norm, the empirical  $\ell_\infty$  norm covering number of  $\mathcal{F}$  can be bounded by the empirical  $\ell_\infty$  norm covering number of  $\mathcal{P}(\mathcal{F})$ . Third, the empirical  $\ell_\infty$  norm covering number of  $\mathcal{P}(\mathcal{F})$  can be bounded by the fat-shattering dimension, and the fat-shattering dimension can be bounded by the worst-case Rademacher complexity of  $\mathcal{P}(\mathcal{F})$ . Hence, the problem is transferred to the estimation of the worst-case Rademacher complexity. Finally, we estimate the lower bound of the worst-case Rademacher complexity of  $\mathcal{P}(\mathcal{F})$ , and then combined with the above steps, the Rademacher complexity of the loss function space associated with  $\mathcal{F}$  can be bounded.

We first introduce the following lemmas:

**Lemma D.4** (Lemma A.2 in (Srebro et al., 2010)). *For any function class  $\mathcal{F}$ , any  $S$  with a finite sample of size  $n$  and any  $\epsilon > 2\hat{\mathfrak{R}}_S(\mathcal{F})$ , we have that*

$$\text{fat}_\epsilon(\mathcal{F}) \leq \frac{16n\hat{\mathfrak{R}}_S^2(\mathcal{F})}{\epsilon^2}.$$

**Lemma D.5** (Theorem 4.4 in (Rudelson & Vershynin, 2006)). *For any function class  $\mathcal{F}$ , any  $S$  with a finite sample of size  $n$ , and any  $\eta \in (0, 1)$  there exist constants  $0 < a < 1$  and  $E > 0$  such that for all  $\epsilon \in (0, 1)$ ,*

$$\log \mathcal{N}_\infty(\epsilon, \mathcal{F}, S) \leq Ev \log(en/v\epsilon) \log^\eta(en/v),$$

where  $v = \text{fat}_{a\epsilon}(\mathcal{F})$ .

**Step 1:** We first derive the relationship between the empirical  $\ell_\infty$  norm covering number  $\mathcal{N}_\infty(\epsilon, \mathcal{L}, D)$  and the empirical  $\ell_\infty$  norm covering number  $\mathcal{N}_\infty(\epsilon, \mathcal{P}(\mathcal{F}), [c] \times D)$ .

For the dataset  $D = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$  with  $n$  i.i.d. examples:

$$\begin{aligned} & \max_i |L(\mathbf{f}(\mathbf{x}_i), \mathbf{y}_i) - L(\mathbf{f}'(\mathbf{x}_i), \mathbf{y}_i)| \\ & \leq \mu \max_i \|\mathbf{f}(\mathbf{x}_i) - \mathbf{f}'(\mathbf{x}_i)\|_\infty \quad (\text{Use Assumption 3.9}) \\ & \leq \mu \max_i \max_j |f_j(\mathbf{x}_i) - f'_j(\mathbf{x}_i)| \\ & \leq \mu \max_i \max_j |p_j(\mathbf{f}(\mathbf{x}_i)) - p_j(\mathbf{f}'(\mathbf{x}_i))|. \quad (\text{The definition of the projection function class } \mathcal{P}(\mathcal{F})) \end{aligned}$$

Then, according to the definition of the empirical  $\ell_\infty$  covering number, we have that an empirical  $\ell_\infty$  cover of  $\mathcal{P}(\mathcal{F})$  at radius  $\epsilon/\mu$  is also an empirical  $\ell_\infty$  cover of the loss function space associated with  $\mathcal{F}$  at radius  $\epsilon$ , and we can conclude that:

$$\mathcal{N}_\infty(\epsilon, \mathcal{L}, D) \leq \mathcal{N}_\infty\left(\frac{\epsilon}{\mu}, \mathcal{P}(\mathcal{F}), [c] \times D\right). \quad (13)$$

**Step 2:** We show that the empirical  $\ell_\infty$  norm covering number of  $\mathcal{P}(\mathcal{F})$  can be bounded by the fat-shattering dimension, and the fat-shattering dimension can be bounded by the worst-case Rademacher complexity of  $\mathcal{P}(\mathcal{F})$ .

According to Lemma D.4, for any  $\epsilon > 2\hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F}))$ , we have

$$\text{fat}_\epsilon(\mathcal{P}(\mathcal{F})) \leq \frac{16nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{\epsilon^2}.$$

Then, combining with Lemma D.5, for any  $\eta \in (0, 1)$  there exist constants  $0 < a < 1$  and  $E > 0$  such that for all  $\epsilon \in (0, 1)$ , we have

$$\log \mathcal{N}_\infty(\epsilon, \mathcal{P}(\mathcal{F}), [c] \times D) \leq Ev \log(enc/v\epsilon) \log^\eta(enc/v),$$

where  $v = \text{fat}_{a\epsilon}(\mathcal{P}(\mathcal{F})) \leq \frac{16nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2}$ , and we set  $d = \frac{16nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2}$ . Hence,

$$\begin{aligned}
 & \log \mathcal{N}_\infty(\epsilon, \mathcal{P}(\mathcal{F}), [c] \times D) \\
 & \leq Ev \log \frac{enc}{v\epsilon} \log^\eta \frac{enc}{v} \\
 & \leq Ed \log \frac{e^{2+\eta}nc}{d\epsilon} \log^\eta \frac{e^{2+\eta}nc}{d} \\
 & \leq Ed \log^{1+\eta} \frac{e^{2+\eta}nc}{d\epsilon} \\
 & \leq E \frac{16nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2} \log^{1+\eta} \left( \frac{1.2}{\hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F}))} \right)^2, \quad (\text{Use inequality } 2\hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F})) < \epsilon < 1) \quad (14)
 \end{aligned}$$

where we use that for any  $s, t > 0$ , the function  $x \mapsto x \log \frac{s}{x} \log^\eta \frac{t}{x}$  is non-decreasing as long as  $s > t > e^{1+\eta}x$  and that  $v \leq nc$  in the second inequality, the third and fourth inequalities are obtained by direct calculation and proper scaling.

Then, we have

$$\begin{aligned}
 & \log \mathcal{N}_\infty(\epsilon, \mathcal{P}(\mathcal{F}), [c] \times D) \\
 & \leq E \frac{16nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2} \log^{1+\eta} \left( \frac{1.2}{\hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F}))} \right)^2 \\
 & \leq E \frac{64nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2} \log^{1+\eta} \left( \frac{1.2}{\hat{\mathfrak{R}}_{[c] \times D}(\mathcal{P}(\mathcal{F}))} \right) \\
 & \leq E \frac{64nc\hat{\mathfrak{R}}_{[c] \times D}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2} \log^{1+\eta} \frac{8\sqrt{nc}}{5B} \quad (\text{Use inequality (10)}) \\
 & \leq E \frac{64nc\tilde{\mathfrak{R}}_{nc}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2} \log^{1+\eta} \frac{8\sqrt{nc}}{5B} \quad (\text{The definition of the worst-case Rademacher complexity}) \quad (15)
 \end{aligned}$$

**Step 3:** According to Lemma C.5 and combined with the above steps, we have

$$\begin{aligned}
 & \hat{\mathfrak{R}}_D(\mathcal{L}) \\
 & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^M \sqrt{\log \mathcal{N}_\infty(\epsilon, \mathcal{L}, D)} d\epsilon \right) \\
 & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^M \sqrt{\log \mathcal{N}_\infty\left(\frac{\epsilon}{\mu}, \mathcal{P}(\mathcal{F}), [c] \times D\right)} d\epsilon \right) \quad (\text{Use inequality (14)}) \\
 & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{12}{\sqrt{n}} \int_\alpha^M \sqrt{\frac{64E\mu^2nc\tilde{\mathfrak{R}}_{nc}^2(\mathcal{P}(\mathcal{F}))}{a^2\epsilon^2} \log^{1+\eta} \frac{8\sqrt{nc}}{5B}} d\epsilon \right) \quad (\text{Use inequality (15)}) \\
 & \leq \inf_{\alpha > 0} \left( 4\alpha + \frac{96\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \log^{\frac{1+\eta}{2}} \left( \frac{8\sqrt{nc}}{5B} \right) \int_\alpha^M \epsilon^{-1} d\epsilon \right) \\
 & \leq \frac{96\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) + \frac{96\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \log^{\frac{1+\eta}{2}} \left( \frac{8\sqrt{nc}}{5B} \right) \cdot \log \frac{aM}{24\sqrt{E}c\mu\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))} \\
 & \quad (\text{Choose } \alpha = \frac{24\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F}))) \\
 & \leq \frac{96\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) + \frac{96\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \log^{\frac{1+\eta}{2}} \left( \frac{8\sqrt{nc}}{5B} \right) \cdot \log \frac{M\sqrt{n}}{\sqrt{E}\mu B} \quad (\text{Use } 0 < a < 1 \text{ and inequality (10)}) \\
 & = \frac{96\sqrt{E}c\mu}{a} \tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{8\sqrt{nc}}{5B} \right) \cdot \log \frac{M\sqrt{n}}{\sqrt{E}\mu B} \right).
 \end{aligned}$$

### D.6. Proof of Theorem 3.17

According to the inequality (11), we have  $\tilde{\mathfrak{R}}_{nc}(\mathcal{P}(\mathcal{F})) \leq \frac{B}{\sqrt{nc}}$ , then

$$\hat{\mathfrak{R}}_D(\mathcal{F}) \leq \frac{96\sqrt{E}\mu B}{a} \frac{1}{\sqrt{n}} \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{8\sqrt{nc}}{5B} \right) \cdot \log \frac{M\sqrt{n}}{\sqrt{E}\mu B} \right).$$

Combining with (7), then

$$R(\mathbf{f}) \leq \hat{R}_D(\mathbf{f}) + \frac{192\sqrt{E}\mu B}{a} \frac{1}{\sqrt{n}} \times \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{8}{5B} \sqrt{nc} \right) \cdot \log \left( \frac{M\sqrt{n}}{\sqrt{E}\mu B} \right) \right) + 3M \sqrt{\frac{\log \frac{2}{\delta}}{2n}}.$$

## E. Generalization Analysis for Macro-Averaged AUC

### E.1. Proof of Theorem 4.2

We first proof the following lemma:

**Lemma E.1.** *Let  $q_\tau : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  be real-valued functions indexed by  $\tau \in T$  where  $T$  is some set. If  $\mathbf{x}_1, \dots, \mathbf{x}_s$  and  $\mathbf{x}'_1, \dots, \mathbf{x}'_t$  are i.i.d.,  $r = \min\{s, t\}$ , then for any convex non-decreasing function  $\psi$ ,*

$$\mathbb{E} \psi \left( \sup_{\tau \in T} \frac{1}{st} \sum_{i=1}^s \sum_{j=1}^t q_\tau(\mathbf{x}_i, \mathbf{x}'_j) \right) \leq \mathbb{E} \psi \left( \sup_{\tau \in T} \frac{1}{r} \sum_{i=1}^r q_\tau(\mathbf{x}_i, \mathbf{x}'_i) \right).$$

*Proof.* The proof of this lemma is inspired by (Cl  men  on et al., 2008).

$$\begin{aligned} & \mathbb{E} \psi \left( \sup_{\tau \in T} \frac{1}{st} \sum_{i=1}^s \sum_{j=1}^t q_\tau(\mathbf{x}_i, \mathbf{x}'_j) \right) \\ &= \mathbb{E} \psi \left( \sup_{\tau \in T} \frac{1}{s!} \sum_{\pi_{\mathbf{x}}} \frac{1}{t!} \sum_{\pi_{\mathbf{x}'}} \frac{1}{r} \sum_{i=1}^r q_\tau(\mathbf{x}_{\pi(i)}, \mathbf{x}'_{\pi(i)}) \right) \\ &\leq \mathbb{E} \psi \left( \frac{1}{s!} \sum_{\pi_{\mathbf{x}}} \frac{1}{t!} \sum_{\pi_{\mathbf{x}'}} \sup_{\tau \in T} \frac{1}{r} \sum_{i=1}^r q_\tau(\mathbf{x}_{\pi(i)}, \mathbf{x}'_{\pi(i)}) \right) \quad (\psi \text{ is nondecreasing}) \\ &\leq \frac{1}{s!} \sum_{\pi_{\mathbf{x}}} \frac{1}{t!} \sum_{\pi_{\mathbf{x}'}} \mathbb{E} \psi \left( \sup_{\tau \in T} \frac{1}{r} \sum_{i=1}^r q_\tau(\mathbf{x}_{\pi(i)}, \mathbf{x}'_{\pi(i)}) \right) \quad (\text{Jensen's inequality}) \\ &= \mathbb{E} \psi \left( \sup_{\tau \in T} \frac{1}{r} \sum_{i=1}^r q_\tau(\mathbf{x}_i, \mathbf{x}'_i) \right). \end{aligned}$$

□

With this lemma, we first prove that  $\mathbb{E} \left[ R(\hat{\mathbf{f}}^*) \right] - R^* \leq \frac{384\sqrt{2E}\mu B}{a} \sqrt{\frac{c}{r_0}} \left( 1 + \log^{\frac{1+\eta}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right).$

$$\begin{aligned} & \mathbb{E} \left[ R(\hat{\mathbf{f}}^*) \right] - R^* = \mathbb{E} \left[ R(\hat{\mathbf{f}}^*) - \hat{R}_D(\hat{\mathbf{f}}^*) + \hat{R}_D(\hat{\mathbf{f}}^*) - R^* \right] \\ &= \mathbb{E} \left[ R(\hat{\mathbf{f}}^*) - \hat{R}_D(\hat{\mathbf{f}}^*) \right] + \mathbb{E} \left[ \hat{R}_D(\hat{\mathbf{f}}^*) - R^* \right] \\ &\leq \mathbb{E} \sup_{\mathbf{f} \in \mathcal{F}} \left| R(\mathbf{f}) - \hat{R}_D(\mathbf{f}) \right| + \mathbb{E} \sup_{\mathbf{f} \in \mathcal{F}} \left| \hat{R}_D(\mathbf{f}) - R(\mathbf{f}) \right| = 2 \mathbb{E} \sup_{\mathbf{f} \in \mathcal{F}} \left| R(\mathbf{f}) - \hat{R}_D(\mathbf{f}) \right| \\ &\leq 2 \mathbb{E} \sup_{\mathbf{f} \in \mathcal{F}} \left| R(\mathbf{f}) - \frac{1}{c} \sum_{j=1}^c \frac{1}{r_j} \sum_{i=1}^{r_j} \ell(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)) \right|. \quad (\text{Use Lemma E.1}) \end{aligned}$$

Let  $\bar{D} = \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_n\}$  be an independent copy of  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , then by the standard symmetrization technique and the Jensen's inequality similar to Preliminaries in Section B, the above inequality can be bounded by:

$$\begin{aligned}
 & 2\mathbb{E}_{D, \bar{D}} \sup_{\mathbf{f} \in \mathcal{F}} \left| \frac{1}{c} \sum_{j=1}^c \frac{1}{r_j} \sum_{i=1}^{r_j} \ell(f_j(\bar{\mathbf{x}}_i) - f_j(\bar{\mathbf{x}}'_i)) - \frac{1}{c} \sum_{j=1}^c \frac{1}{r_j} \sum_{i=1}^{r_j} \ell(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)) \right| \\
 &= 2\mathbb{E}_{D, \bar{D}, \epsilon} \sup_{\mathbf{f} \in \mathcal{F}} \left| \frac{1}{c} \sum_{j=1}^c \frac{1}{r_j} \sum_{i=1}^{r_j} \epsilon_{ij} [\ell(f_j(\bar{\mathbf{x}}_i) - f_j(\bar{\mathbf{x}}'_i)) - \ell(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i))] \right| \\
 &= 4\mathbb{E}_{D, \epsilon} \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{c} \sum_{j=1}^c \frac{1}{r_j} \sum_{i=1}^{r_j} \epsilon_{ij} \ell(f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)) \\
 &= 4\mathfrak{R}(\mathcal{F}) \\
 &\leq \frac{192\sqrt{2E}\mu}{a} \sqrt{c} \max_j \tilde{\mathfrak{R}}(\mathcal{F}_j) \times \left( 1 + \log^{\frac{1+n}{2}} \left( \frac{2\sqrt{2r_j}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{r_j}{c}} \right) \right) \quad (\text{Use Lemma 3.11}) \\
 &= \frac{192\sqrt{2E}\mu}{a} \sqrt{c} \max_j \left( 1 + \log^{\frac{1+n}{2}} \left( \frac{2\sqrt{2r_j}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{r_j}{c}} \right) \right) \mathbb{E}_D \sup_{D \in \mathcal{X}^n} \mathbb{E}_\epsilon \left[ \sup_{\mathbf{f}_j \in \mathcal{F}_j} \frac{1}{r_j} \sum_{i=1}^{r_j} \epsilon_i (f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i)) \right] \\
 &\leq \frac{192\sqrt{2E}\mu}{a} \sqrt{c} \max_j \left( 1 + \log^{\frac{1+n}{2}} \left( \frac{2\sqrt{2r_j}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{r_j}{c}} \right) \right) \mathbb{E}_D \sup_{D \in \mathcal{X}^n} \sup_{\mathbf{f}_j \in \mathcal{F}_j} \frac{1}{r_j} \left( \sum_{i=1}^{r_j} (f_j(\mathbf{x}_i) - f_j(\mathbf{x}'_i))^2 \right)^{\frac{1}{2}} \\
 &\quad (\text{Use Lemma C.2}) \\
 &\leq \frac{384\sqrt{2E}\mu B}{a} \max_j \sqrt{\frac{c}{r_j}} \left( 1 + \log^{\frac{1+n}{2}} \left( \frac{2\sqrt{2r_j}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{r_j}{c}} \right) \right) \\
 &\leq \frac{384\sqrt{2E}\mu B}{a} \sqrt{\frac{c}{r_0}} \left( 1 + \log^{\frac{1+n}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right).
 \end{aligned}$$

Similarly, we can derive that

$$\begin{aligned}
 & R(\hat{\mathbf{f}}^*) - R^* \\
 &= R(\hat{\mathbf{f}}^*) - \hat{R}_D(\hat{\mathbf{f}}^*) + \hat{R}_D(\hat{\mathbf{f}}^*) - R^* \\
 &\leq \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \hat{R}_D(\mathbf{f})] + \sup_{\mathbf{f} \in \mathcal{F}} [\hat{R}_D(\mathbf{f}) - R(\mathbf{f})] \\
 &= 2 \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \hat{R}_D(\mathbf{f})].
 \end{aligned}$$

Let  $D' = \{\mathbf{x}_1, \dots, \mathbf{x}'_k, \dots, \mathbf{x}_n\}$  be the training dataset with only one sample different from  $D$ , where the  $k$ -th sample is replaced by  $\mathbf{x}'_k$ . Then

$$\begin{aligned}
 & \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \hat{R}_{D'}(\mathbf{f})] - \sup_{\mathbf{f} \in \mathcal{F}} [R(\mathbf{f}) - \hat{R}_D(\mathbf{f})] \\
 &\leq \sup_{\mathbf{f} \in \mathcal{F}} [\hat{R}_D(\mathbf{f}) - \hat{R}_{D'}(\mathbf{f})] \\
 &\leq \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{c} \sum_{j=1}^c \frac{1}{|X_j^+| |X_j^-|} \sum_{\mathbf{x}'_i \in X_j^-} (\ell(f_j(\mathbf{x}_k) - f_j(\mathbf{x}'_i)) - \ell(f_j(\mathbf{x}_k) - f_j(\mathbf{x}_i))) \\
 &\leq \frac{M}{r_0}.
 \end{aligned}$$

According to McDiarmid's inequality, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$  over the training dataset  $D$ , the following holds:

$$R(\hat{\mathbf{f}}^*) - R^* \leq \frac{384\sqrt{2E}\mu B}{a} \sqrt{\frac{c}{r_0}} \left( 1 + \log^{\frac{1+n}{2}} \left( \frac{2\sqrt{2n}}{B} \right) \cdot \log \left( \frac{M}{\sqrt{E}\mu B} \sqrt{\frac{n}{c}} \right) \right) + 2M \sqrt{\frac{\log \frac{1}{\delta}}{r_0}}.$$

Note that when we analyze class-imbalance, we focus on whether  $s_j$  and  $t_j$  are seriously imbalanced for a fixed number of  $n$  examples. Hence, the order of the bound is  $\tilde{O}(\sqrt{\frac{c}{r_0}})$ .