# Learning with Multiclass AUC: Theory and Algorithms

Zhiyong Yang, Qianqian Xu*, *Senior Member, IEEE,* Shilong Bao,
Xiaochun Cao, *Senior Member, IEEE,* and Qingming Huang*, *Fellow, IEEE*

✦

**Abstract**—The Area under the ROC curve (AUC) is a well-known ranking metric for problems such as imbalanced learning and recommender systems. The vast majority of existing AUC-optimization-based machine learning methods only focus on binary-class cases, while leaving the multiclass cases unconsidered. In this paper, we start an early trial to consider the problem of learning multiclass scoring functions via optimizing multiclass AUC metrics. Our foundation is based on the M metric, which is a well-known multiclass extension of AUC. We first pay a revisit to this metric, showing that it could eliminate the imbalance issue from the minority class pairs. Motivated by this, we propose an empirical surrogate risk minimization framework to approximately optimize the M metric. Theoretically, we show that: (i) optimizing most of the popular differentiable surrogate losses suffices to reach the Bayes optimal scoring function asymptotically; (ii) the training framework enjoys an imbalance-aware generalization error bound, which pays more attention to the bottleneck samples of minority classes compared with the traditional $O(\sqrt{1/N})$ result. Practically, to deal with the low scalability of the computational operations, we propose acceleration methods for three popular surrogate loss functions, including the exponential loss, squared loss, and hinge loss, to speed up loss and gradient evaluations. Finally, experimental results on 11 real-world datasets demonstrate the effectiveness of our proposed framework.

**Index Terms**—AUC Optimization, Machine Learning.

* corresponding authors

Zhiyong Yang with the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 101408, China (email: yangzhiyong@iie.ac.cn).

Qianqian Xu is with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China, (email: xuqianqian@ict.ac.cn).

- Shilong Bao is with the State Key Laboratory of Information Security, Institute of Information Engineering, Chinese Academy of Sciences, Beijing 100093, China, and also with the School of Cyber Security, University of Chinese Academy of Sciences, Beijing 100049, China (email: baoshilong@iie.ac.cn).

- Xiaochun Cao is with the State Key Laboratory of Information Security, Institute of Information Engineering, Chinese Academy of Sciences, Beijing 100093, China, also with School of Cyber Science and Technology, Sun Yat-sen University, Shenzhen, 518100, China (email: caoxiaochun@iie.ac.cn).

- Q. Huang is with the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 101408, China, also with the Key Laboratory of Big Data Mining and Knowledge Management (BDKM), University of Chinese Academy of Sciences, Beijing 101408, China, also with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China, and also with Peng Cheng Laboratory, Shenzhen 518055, China (e-mail: qmhuang@ucas.ac.cn).

## 1 INTRODUCTION

AUC (Area Under the ROC Curve), which measures the probability that a positive instance has a higher score than a negative instance, is a well-known performance metric for a scoring function's ranking quality. AUC often comes up as a more appropriate performance metric than accuracy in various applications due to its appealing properties, *e.g.*, insensitivity toward label distributions and costs [20], [30]. On one hand, for class-imbalanced tasks such as disease prediction [32], [79] and rare event detection [43], [48], [73], the label distribution is often highly skewed in the sense that the proportion of the majority classes instances significantly dominates the others. As for a typical instance, in the credit card fraud detection dataset released in Kaggle [1], the fraudulent transactions only account for $0.172\%$ of the total records. Accuracy is often not a good choice in this case since it might ignore the performance from the minority classes which are often more crucial than the majority ones. By contrast, the value of AUC does not rest on the label distribution, making it a natural metric under the class-imbalance scenario. On the other hand, AUC has also been adopted as a standard metric for applications such as ads click-through rate prediction and recommender systems [10], [17], [19], [61] where pursuing a correct ranking between positive and negative instances is much more critical than label prediction.

Over the past two decades, the importance of AUC has raised an increasing favor in the machine learning community to explore direct AUC optimization methods, *e.g.*, [3], [8], [23], [35], [36], [55], [58]. However, to our knowledge, the vast majority of related studies merely focus on the binary class scenario. Since real-world pattern recognition problems often involve more than two classes, it is natural to pursue its generalization in the world of multiclass problems. To this end, we present a very early study to the problem of AUC guided machine learning framework under the multiclass setting. Specifically, we propose a universal empirical surrogate risk minimization framework with theoretical guarantees.

First of all, we provide a review of a multiclass generalization of AUC known as the M metric [30]. Specifically, we show that M metric is an appropriate extension of binary AUC in the sense that it could efficiently avoid the imbalanced issue across class ranking pairs. Motivated by this result, we propose to minimize the 0-1 mis-ranking

loss induced by the M metric denoted as MAUC$^\downarrow$. However, directly minimizing the objective is intractable. This intractability has three sources: **(a)** the 0-1 mis-ranking loss MAUC$^\downarrow$ is a discrete and non-differentiable function, making it impossible to perform efficient optimizations; **(b)** the data distribution is not a known *priori*, making the calculation of expectation unavailable; **(c)** the complexity to estimate the loss and gradient function could be as high as $O(N^2 N_C^2)$ in the worst case, where $N$ is the number of samples and $N_C$ is the number of classes.

Targeting at **(a)**, in Sec.4, we investigate how to construct differentiable surrogate risks as replacements for the 0-1 risk. To do this, we derive the set of Bayes optimal score functions under the MAUC$^\downarrow$ criterion. We then provide a general sufficient condition for the fisher consistency. The condition suggests that a large set of popular surrogate loss functions are fisher consistent under certain assumptions in the sense that optimizing the corresponding surrogate expected risk also leads to the Bayes optimal score functions.

Based on the consistency result, in Sec.5, we construct an empirical surrogate risk minimization framework against **(b)**, where we propose an unbiased empirical estimation of the surrogate risks over a training dataset as the objective function to avoid using population-level expectation directly. Moreover, we provide a systematic analysis of its generalization ability. The major challenge here is that the empirical risk function could not be decomposed as a finite sum of independent instance-wise loss terms, making the traditional symmetrization technique not available. To fix this problem, we provide a novel form of Rademacher complexity for MAUC$^\downarrow$. Furthermore, we provide generalization upper bounds for a wide range of model classes, including shallow and deep models. The results consistently enjoy an imbalance-aware property, which pays more attention to the bottleneck samples of minority classes than the traditional result.

In Sec.6, **(c)** is attacked by novel acceleration algorithms to speed up loss and gradient evaluations which are the fundamental calculations for gradient-based optimization methods. The mini-batch/full-batch loss and gradient evaluations could be done, with the proposed algorithms, in a complexity comparable with that for the ordinary instance-wise loss functions.

Finally, in Sec.7, we perform extensive experiments on 11 real-world datasets to validate our proposed framework.

To sum up, our contribution is three-fold:

- **New Framework**: We provide a theoretical framework for empirical risk minimization under the guidance of the M metric.
- **Theoretical Guarantees**: From a theoretical perspective, our framework is soundly supported by consistency analysis and generalization analysis.
- **Fast Algorithms**: From a practical perspective, we propose efficient algorithms for three surrogate losses to accelerate loss and gradient evaluations. The experiments show that the *acceleration ratio could reach* **1,0000+** confronting medium size datasets.

## 2 RELATED WORK

**Binary Class AUC Optimization.** As a motivating early study, [15] points out that maximizing AUC should not be replaced with minimizing the error rate, which shows the necessity to study direct AUC optimization methods. After that, a series of algorithms have been designed for optimizing AUC. At the early stage, the majority of studies focus on a full-batch off-line setting. [3], [8] optimize AUC based on a logistic surrogate loss function and ordinary gradient descent method. RankBoost [22] provides an efficient ensemble-based AUC learning method based on a ranking extension of the AdaBoost algorithm. [36], [77] constructs $SVM^{struct}$-based frameworks that optimize a direct upper bound of the $0-1$ loss version AUC metric instead of its surrogates. Later on, to accommodate big data analysis, researchers start to explore online extensions of AUC optimization methods. [78] provides an early trial for this direction based on the reservoir sampling technique. [23] provides a completely one-pass AUC optimization method for streaming data based on the squared surrogate loss. [76] reformulates the squared-loss-based stochastic AUC maximization problem as a stochastic saddle point problem. The new saddle point problem's objective function only involves summations of instance-wise loss terms, which significantly reduces the burden from the pairwise formulation. [59], [60] further accelerate this framework with tighter convergence rates. On top of the reformulation framework, [75] also provides an acceleration framework for general loss functions where the loss functions are approximated by the Bernstein polynomials. [18] proposes a novel large-scale nonlinear AUC maximization method based on the triply stochastic gradient descent algorithm. [28] proposes a scalable and efficient adaptive doubly stochastic gradient algorithm for generalized regularized pairwise learning problems. Beyond optimization methods, a substantial amount of researches also provide theoretical support for this learning framework from different dimensions, including generalization analysis [2], [12], [64], [70], [71] and consistency analysis [1], [25]. Last but not least, there are also some studies focusing on optimizing the partial area under the ROC curve [4], [56], [57]. This paper takes a further step by providing an early study of the theory and algorithms for AUC-guided machine learning under the more complicated multiclass scenario.

**Multiclass AUC Metrics**. There exist two general ideas for how to define a multiclass AUC metric. The first idea is on top of the belief that multiclass counterparts of the ROC curve should be represented as higher-dimensional surfaces. As a result, AUC is generalized naturally to the volume Under some specifically designed ROC Surfaces (VUS) [21], [54]. However, this idea is restricted by its high complexity to calculate the volume of such high-dimensional spaces. According to [38], calculating VUS for $N$ samples and $N_C$ classes requires $O(N \log N + N^{\lfloor N_C/2 \rfloor})$ time complexity and $O(N^{\lfloor N_C/2 \rfloor})$ space complexity. Another idea then comes out to do it in a much simpler manner, which suggests that one can simply take an average of pairwise binary AUCs [30], [33], [62], [74]. This is based on the intuition that if every pair of classes is well-separated from each other in distribution, one can get reasonable performances. Getting rid of calculations on high dimensional spaces renders this

formulation a low complexity. Due to its simplicity, the M metric proposed in the representative work [30] has been adopted by a series of popular machine learning software such as `scikit-learn` in `Python` and `pROC` in `R`. Most recently, [72] also considers online extensions of multi-class AUC metric to deal with the concept drift issue for streaming data. Unlike this line of research, our main focus in this paper is how to learn valuable models from a proper multiclass extension of the AUC metric.

**AUC Optimization for Multiclass AUC Optimization**. There are few studies that focus on AUC optimization algorithms for multiclass problems, which fall in the following two directions. The first direction of studies focuses on the multipartite ranking problem, which is a natural extension of the bipartite ranking problem where the order is presented with more than two discrete degrees. Hitherto, a substantial amount of efforts have been made to explore the AUC optimization method/theory for the multipartite ranking problem [11], [13], [24], [63], [69]. Moreover, a recent work [66] proposes a novel nonlinear semi-supervised multipartite ranking problem for large-scale datasets. It is noteworthy that multipartite ranking approaches could solve multiclass problems only if the classes are ordinal values for the same semantic concept. For example, the age estimation task could be regarded as a multiclass problem where the class labels are the ages for a given person; the movie rating prediction could be regarded as a multiclass problem where the class labels are the ratings for the same given movie. The major difference here is that we focus on the generic multiclass problems where different labels present different semantic concepts and do not have clear ordinal relations. Thus the studies along this line are not available to our setting in general. Most recently, [44] discusses the possibility of exploring AUC optimization for general multiclass settings. The major differences are as follows. First, [44] only focuses on the square loss function, while our study presents a general framework for multiclass AUC optimization. Second, [44] focuses more on the optimization properties, where the original problem is reformulated as a minimax problem. By contrast, our study focuses on the learning properties for multiclass AUC optimization, such as its generalization ability and the consistency of different loss functions. Moreover, we also present a series of acceleration methods that do not require any reformulation of the optimization problem.

# 3 LEARNING WITH MULTICLASS AUC METRICS: AN OVERVIEW

## 3.1 Preliminary

**Basic Notations**. In this paper, **we will constantly adopt two sets of events**: for pair-wise AUC metrics, $\mathcal{E}^{(ij)}$ denotes the event $y_1 = i, y_2 = j \vee y_1 = j, y_2 = i$, while $\mathcal{E}^{(i)}$ denotes the event $y_1 = i, y_2 \neq i \vee y_1 \neq i, y_2 = i$. Given an event $\mathcal{A}$, $I[\mathcal{A}]$ is the indicator function associated with this event, which equals 1 if $\mathcal{A}$ holds and equals to 0 otherwise. Given a finite dataset $\mathcal{S}$, we denote $N_C$ as the number of classes and $N$ as the total number of sample points in the dataset.

**Settings**. For an $N_C$ class problem, we assume that our samples are drawn from a product space $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. $\mathcal{X}$ is the input feature where $\mathcal{X} \subset \mathbb{R}^d$ and $d$ is the input dimensionality.

$\mathcal{Y}$ is the label space $[N_C]$. Given a label $y_m = i$, we will use the one-hot vector $\boldsymbol{y}_m = [y_m^{(1)}, \cdots, y_m^{(N_C)}]$ to represent it, with $y_m^{(k)} = 0, k \neq i$, and $y_m^{(k)} = 1$ if $k = i$. In this paper, we will adopt the one *vs.* all decomposition [53, Chap.9.4] of the multiclass problem. Under this context, an $N_C$-class ($N_C > 2$) scoring function refers to a set of $N_C$ functions $f = (f^{(1)}, \cdots, f^{(N_C)})$, where $f^{(i)} : \mathcal{X} \to \mathbb{R}$ serves as a *continuous* score function supporting $y = i$.

**Binary Class AUC**. Let $\mathcal{D}_{\mathcal{Z}}$ denote the joint data distribution and $f$ denote a score function estimating the possibility that an underlying instance belongs to the positive class. In the context of binary class problems, AUC is known to have a clear statistical meaning: it is equivalent to the Wilcoxon Statistics [31], namely the possibility that correct ranking takes place between a random paired samples with distinct labels: $\mathrm{AUC}(f) = \mathbb{P}[\Delta(y)\Delta(f) > 0|\mathcal{E}^{(0,1)}] + \frac{1}{2}\mathbb{P}[\Delta(f) = 0|\mathcal{E}^{(0,1)}] = \mathbb{E}_{\boldsymbol{z}_1 \sim \mathcal{D}_{\mathcal{Z}}, \boldsymbol{z}_2 \sim \mathcal{D}_{\mathcal{Z}}} [I[\Delta(y)\Delta(f) > 0] + \frac{1}{2}I[\Delta(f) = 0]|\mathcal{E}^{(0,1)}]$, where $\Delta(y) = y_1 - y_2$, $\Delta(f) = f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2)$. Note that we adopt the convention [1], [12], [25] to score ties with 0.5. See basic notations in the Sec.3.1 for $\mathcal{E}^{(0,1)}$. Hereafter, our calculations on AUC follows this definition.

## 3.2 Motivation

First of all, we start our study by finding a proper multiclass metric from the existing literature. As presented in Sec.2, the main idea to construct multiclass AUC is to express it as an average of binary class AUCs. Just like the way how multiclass classification is transformed into a set of binary classification problems, one can derive multiclass AUC metrics from the following two regimes.

**One *vs.* All Regime (ova)**. Given an ova score function $f = (f^{(1)}, \cdots, f^{(N_C)})$, [62] suggests that one can construct a pairwise AUC score for each $f^{(i)}$. Here, the positive instances are drawn from the $i$-th class, and the negative instances are drawn from the $j$-th distribution conditioned on $j \neq i$. The overall AUC score $\mathrm{AUC}^{\mathrm{ova}}$ is then an average of all $N_C$ pairwise scores, which is defined as follows. Note that here we adopt an equally-weighted average to avoid the imbalance issue.

$$\mathrm{AUC}^{\mathrm{ova}}(f) = \frac{1}{N_C} \sum_{i=1}^{N_C} \mathrm{AUC}_{i|\neg i}(f^{(i)}),$$

where $\mathrm{AUC}_{i|\neg i}(f^{(i)}) = \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2} [I[\Delta(y_{1,2}^{(i)})\Delta(f^{(i)}) > 0|\mathcal{E}^{(i)}] + \frac{1}{2}I[\Delta(f^{(i)}) = 0]|\mathcal{E}^{(i)}]$, $\Delta(y_{1,2}^{(i)}) = y_1^{(i)} - y_2^{(i)}$, $\Delta(f^{(i)}) = f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)$. See basic notations in the Sec.3.1 for the definition of $\mathcal{E}^{(i)}$.

**One *vs.* One Regime (ovo, M metric)**. Alternatively, according to [30], we can formulate a multiclass metric as an average of binary AUC scores for every class pair $(i, j)$, which is defined as:

$$\mathrm{AUC}^{\mathrm{ovo}}(f) = \frac{\sum_{i=1}^{N_C} \sum_{j \neq i} \mathrm{AUC}_{i|j}(f^{(i)})}{N_C(N_C - 1)},$$

where $\mathrm{AUC}_{i|j}(f^{(i)}) = \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2} [I[\Delta(y_{1,2}^{(i)})\Delta(f^{(i)}) > 0|\mathcal{E}^{(ij)}] + \frac{1}{2}I[\Delta(f^{(i)}) = 0]|\mathcal{E}^{(ij)}]$, note that $\mathrm{AUC}_{i|j} \neq \mathrm{AUC}_{j|i}$, since they employ different score functions. See basic notations in Sec.3.1 for the definition of $\mathcal{E}^{(ij)}$.

The following theorem reveals that $\text{AUC}^{\text{ovo}}$ is more insensitive toward the imbalanced distribution of the class pairs than $\text{AUC}^{\text{ova}}$.

**Theorem 1 (Comparison Properties).** *Given the label distribution as* $\mathbb{P}[y = i] = p_i > 0$ *and a multiclass scoring function* $f$. *The following properties hold:*

*(a) We have that :*

$$\text{AUC}^{\text{ova}}(f) = 1/N_C \sum_{i=1}^{N_C} \sum_{j \neq i} (p_j/1-p_i) \cdot \text{AUC}_{i|j}(f^{(i)}). \quad (1)$$

*(b)*

$$\text{AUC}^{\text{ova}}(f) = \text{AUC}^{\text{ovo}}(f), \ \text{when} \ \ p_i = 1/N_C,$$
$$i = 1, 2, \cdots, N_C.$$

*(c) We have* $\text{AUC}^{\text{ova}}(f) = 1$ *if and only if* $\text{AUC}^{\text{ovo}}(f) = 1$.

Thm.1-(a) states that $\text{AUC}^{\text{ova}}$ weights different $\text{AUC}_{i|j}$ with $p_j/1-p_i$, which will overlook the performance of the minority class pairs. On the contrary, $\text{AUC}^{\text{ovo}}$ assigns equal weights for all pair-wise AUCs, which naturally avoids this issue. Thm.1-(b) suggests that $\text{AUC}^{\text{ova}}$ and $\text{AUC}^{\text{ovo}}$ tend to be equivalent when the label distribution is nearly balanced. Thm.1-(c) further shows that $\text{AUC}^{\text{ova}}$ and $\text{AUC}^{\text{ovo}}$ agree with each other when $f$ maximizes the performance. Practically, the following example shows that how the imbalance issue of class pairs affects model selection under different criterions.

**Example 1.** *Consider a three-class classification dataset with a label distribution* $p_1 = 0.5, p_2 = 0.45, p_3 = 0.05$, *we assume that there are two scoring functions* $f_a, f_b$, *where* $\text{AUC}_{i|j}(f_a)$ *are all 1 except that* $\text{AUC}_{1|3}(f_a) = 0.5$, $\text{AUC}_{i|j}(f_b)$ *are all 1 except that* $\text{AUC}_{1|2}(f_b) = 0.8$. *Consequently, we have* $\text{AUC}^{\text{ova}}(f_a) = 98.3$ *and* $\text{AUC}^{\text{ova}}(f_b) = 93.3$, *while* $\text{AUC}^{\text{ovo}}(f_a) = 91.6$ *and* $\text{AUC}^{\text{ovo}}(f_b) = 96.6$.

Obviously, $f_a$ in the example should be a bad scoring function since $f_a^{(1)}$ can not tell apart class-1 and class-3, while $f_b$ is a much better choice since it has good performances in terms of all the pairwise AUCs. We see that $\text{AUC}^{\text{ovo}}$ supports choosing $f_b$ against $f_a$, which is consistent with our expectations. However, $\text{AUC}^{\text{ova}}$ chooses $f_a$ against $f_b$ with a significant margin. This is because that the minority class 3 brings an extremely low weight on $\text{AUC}_{1|3}$ in $\text{AUC}^{\text{ova}}$. This tiny weight makes the fatal disadvantage of $\text{AUC}_{1|3}(f_a)$ totally ignored.

From the theoretical and practical analysis provided above, we can draw the conclusion that $\text{AUC}^{\text{ovo}}$ is a better choice than $\text{AUC}^{\text{ova}}$.

### 3.3 Objective and the Roadmap

**Objective**. Our goal in this paper is then to construct learning algorithms that maximize $\text{AUC}^{\text{ovo}}$. To fit in the standard machine learning paradigm, we follow the widely-adopted convention [1], [12], [25], [56], [57] to cast the maximization problem into an expected-risk-based minimization problem $f \in \text{argmin}_f R(f)$, where:

$$R(f) = \text{MAUC}^{\downarrow} = \sum_{i=1}^{N_C} \sum_{j \neq i} \frac{\mathbb{E}_{z_1, z_2}[\ell_{0-1}^{i,j,1,2}|\mathcal{E}^{(ij)}]}{N_C(N_C - 1)},$$

$$\ell_{0-1}^{i,j,1,2} = \ell_{0-1}(f^{(i)}, \boldsymbol{x}_1, \boldsymbol{x}_2, y_1^{(i)}, y_2^{(i)}),$$
$$= \boldsymbol{I}\left[\Delta(y_{1,2}^{(i)}) \cdot \Delta(f^{(i)}) < 0\right] + \frac{1}{2}\boldsymbol{I}\left[\Delta(f^{(i)}) = 0\right]$$

is the 0-1 mis-ranking loss.

**Roadmap**. The major challenge in this work is that directly minimizing of $R(f)$ is almost impossible in the sense that: **(a)** the 0-1 mis-ranking loss $\text{MAUC}^{\downarrow}$ is not differentiable; **(b)** the calculation of population-level expectation is unavailable; **(c)** the complexity to estimate the ranking loss is $O(N^2 N_C^2)$ in the worst case. In Sec.4-Sec.6, we will present solutions to address **(a)**-**(c)**, respectively.

## 4 BAYES OPTIMAL CLASSIFIER, AND CONSISTENCY ANALYSIS FOR SURROGATE RISK MINIMIZATION

Since the 0-1 mis-ranking loss is a discrete and non-differentiable function, directly solving the optimization problem is almost intractable. In this section, we will construct surrogate risks $R_\ell(f)$ with surrogate losses $\ell$ as differentiable proxies for the 0-1 mis-ranking loss. We start with finding the Bayes optimal scoring function that we should approximate. Then we provide the surrogate risk minimization framework and investigate if it could approximate the true minimization problem well.

**Bayes Optimal Scoring Functions**. First, let us derive what are the best scoring functions that we need to approximate. With a goal to minimize $\text{MAUC}^{\downarrow}$, $f$ reaches the best performance when it realizes the minimum of expected risk under the 0-1 mis-ranking loss. In this paper, since we are dealing with classification problems, we restrict the choice of each $f^{(i)}$ to measurable functions with range $[0, 1]$, *i.e.*,

$$f \in \mathcal{F}_\sigma^{N_C} = \underbrace{\mathcal{F}_\sigma \times \mathcal{F}_\sigma \cdots \times \mathcal{F}_\sigma}_{N_C},$$

where

$$\mathcal{F}_\sigma = \{ \ g : g \text{ is a measurable function with a range } [0, 1] \ \}$$

Given the restriction, the scoring function should be expressed as:

$$f \in \underset{f \in \mathcal{F}_\sigma^{N_C}}{\text{arginf}} \ R(f),$$

Following the convention of machine learning terminologies, we refer to such functions as Bayes optimal scoring functions. The following theorem gives the solution for Bayes optimal scoring functions when the data distribution is a known priori.

**Theorem 2 (Bayes Optimal Scoring Functions).** *Given* $\eta_i(\cdot) = \mathbb{P}[y = i|x]$, $p_i = \mathbb{P}[y = i]$, *we have the following consequences:*

*(a)* $f = \{f^{(i)}\}_{i=1,2,\cdots,N_C} \in \mathcal{F}_\sigma^{N_C}$ *is a Bayes optimal scoring function under the* $\text{MAUC}^{\downarrow}$ *criterion, if:*

$$\Delta(f^{(i)}) \cdot \Delta(\pi) > 0,$$
$$\forall \boldsymbol{x}_1, \boldsymbol{x}_2, \ s.t. \ \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \neq \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1) \quad (2)$$

*where:*

$$\Delta(f^{(i)}) = f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)$$
$$\Delta(\pi) = \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) - \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$$
$$\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) = \sum_{j \neq i} \frac{\eta_i(\boldsymbol{x}_1)\eta_j(\boldsymbol{x}_2)}{2p_i p_j},$$
$$\pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1) = \sum_{j \neq i} \frac{\eta_j(\boldsymbol{x}_1)\eta_i(\boldsymbol{x}_2)}{2p_i p_j}.$$

*(b) Define $\sigma(\cdot)$ as the sigmoid function, $s_i(\boldsymbol{x}) = \eta_i(\boldsymbol{x})/p_i$ and $s_{\setminus i}(\boldsymbol{x}) = \sum_{j \neq i} s_j(\boldsymbol{x})$, then a Bayes optimal scoring function could be given by:*

$$f^{\star(i)}(\boldsymbol{x}) = \begin{cases} \sigma\left(\frac{s_i(\boldsymbol{x})}{s_{\setminus i}(\boldsymbol{x})}\right), & 1 > \eta_i(\boldsymbol{x}) \geq 0, \\ 1, & 1 = \eta_i(\boldsymbol{x}), \end{cases} \quad (3)$$

On top of providing a solution for the optimal scoring function, this theorem also sheds light upon what the optimal scoring functions are really after. Thm.2 shows that the optimal scoring function provides a consistent ranking with $\left(\frac{s_i(\boldsymbol{x})}{s_{\setminus i}(\boldsymbol{x})}\right)$, which could be regarded as a generalized likelihood ratio of $y = i$ vs. $y \neq i$. Here, the posterior distribution $\mathbb{P}(y = i|\boldsymbol{x})$ is weighted by the factor $1/p_i$, which eliminates the dependence on the label distribution since $\mathbb{P}(y = i|\boldsymbol{x})/p_i$ is in proportion to the label-distribution-independent class-conditional distribution $\mathbb{P}(\boldsymbol{x}|y = i)$. This result suggests that the optimal scoring function induced by MAUC$^\downarrow$ is insensitive toward skewed label distribution.

**Surrogate Risk Minimization**. Unfortunately, since the data distribution is not available and the 0-1 loss is not differentiable, the Bayesian scoring functions are intractable even with the solution shown in Thm.2. Practically, we need to replace the 0-1 loss with a convex differentiable loss function $\ell$ to find tractable approximations. Once $\ell$ is fixed, we can naturally minimize the much simpler surrogate risk formulated as $f^\star \in \operatorname{argmin}_f R_\ell(f)$, where

$$R_\ell(f) = \sum_i \frac{R_\ell^{(i)}(f^{(i)})}{N_C(N_C - 1)}$$
$$R_\ell^{(i)}(f^{(i)}) = \sum_{j \neq i} \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2} \left[ \ell(\Delta(y_{1,2}^{(i)})\Delta f^{(i)})|\mathcal{E}^{(ij)} \right].$$

But how could we find out such surrogate losses at all? A potential candidate must be consistent with the 0-1 loss. In other words, one should make sure that Bayes optimal scoring functions be recovered by minimizing the chosen surrogate risk, at least in an asymptotic way. This leads to the following definition of consistency in a limiting sense.

**Definition 1** (MAUC$^\downarrow$ **Consistency**). [1] *$\ell$ is consistent with MAUC$^\downarrow$ if for every function sequence $\{f_t\}_{t=1,2,\dots}$, we have:*

$$R_\ell(f_t) \to \inf_{f \in \mathcal{F}_\sigma^{N_C}} R_\ell(f) \text{ implies } R(f_t) \to \inf_{f \in \mathcal{F}_\sigma^{N_C}} R(f).$$

**Remark 1.** *Note that since the infimum of $f$ is taken over all possible measurable functions with a proper range $[0,1]$, the result is thus irrelevant with the choice of the hypothesis space (linear*

1. This condition should be met for all distributions $\mathcal{D}_{\mathcal{X} \times \mathcal{Y}}$.

*model, Neural Networks). Practically, to reach the theoretical infimum, it is better to consider complicated hypothesis spaces such as deep neural networks.*

Based on this definition, we provide a sufficient condition for MAUC$^\downarrow$ consistency, which is shown in the following theorem. See Appendix B.2 for the proof.

**Theorem 3** (MAUC$^\downarrow$ **Consistency**). *The surrogate loss $\ell$ is consistent with MAUC$^\downarrow$ for all $f \in \mathcal{F}_\sigma^{N_C}$, if it is differentiable, convex, and nonincreasing within $[-1, 1]$ and $\ell'(0) < 0$.*

From the sufficient condition above, we can show that the popular loss functions are all consistent with MAUC$^\downarrow$, which is summarized as the following corollary.

**Corollary 1.** *The following statements hold according to Thm.3:*
1) *Logit loss $\ell_{logit}(x) = \log(1 + \exp(-x))$ is consistent with MAUC$^\downarrow$.*
2) *Exp loss $\ell_{\exp}(x) = \exp(-x)$ is consistent with MAUC$^\downarrow$.*
3) *Square loss $\ell_{sq}(x) = (1 - x)^2$ is consistent with MAUC$^\downarrow$.*
4) *The q-norm hinge loss $\ell_q(x) = (\max(1 - x, 0))^q$ is consistent with MAUC$^\downarrow$, if $q > 1$.*
5) *The generalized hinge loss given by:*

$$\ell_{\epsilon, m}(x) = \begin{cases} m - t, & t \leq 1 - \epsilon \\ (t - 1 - \epsilon)^2/4\epsilon, & 1 - \epsilon \leq t < 1 \\ 0, & otherwise \end{cases}$$

*is consistent with MAUC$^\downarrow$, if $1/2 > \epsilon > 0$.*
6) *The distance-weighted loss given by:*

$$\ell_d(x) = \begin{cases} 1/t, & t > \epsilon \\ 1/\epsilon \cdot (2 - t/\epsilon), & otherwise \end{cases}$$

*is consistent with MAUC$^\downarrow$, if $1 > \epsilon > 0$.*

Note that $\ell_{hinge}(t) = \max(1 - t, 0) = \lim_{\epsilon \to 0} \ell_\epsilon(t)$. This shows that hinge loss is at least a limit point of the set of all consistent losses.

## 5 EMPIRICAL RISK MINIMIZATION AND IMBALANCE AWARE GENERALIZATION ANALYSIS

Thus far, we have defined the suitable replacements for the 0-1 loss and found some popular examples for them. However, the arguments are based on the population version of the risk. Calculating the expectation is hardly available since the data distribution is most likely unknown to us. In this section, we take a step further to explore how to find empirical estimations of the risks given a specific training data set. On top of this, we show theoretical guarantees for such estimations.

### 5.1 Empirical Surrogate Risk Minimization

Given a finite training data $\mathcal{S}$ collected from the data distribution, we start with an unbiased estimation of the expected surrogate risk $R_\ell(f)$ based on $\mathcal{S} = \{\boldsymbol{x}_i, y_i\}_{i=1}^N$. To represent the label frequencies in $\mathcal{S}$, we denote $\mathcal{N}_i = \{\boldsymbol{x}_k : y_k = i, \boldsymbol{z}_k \in \mathcal{S}\}$ as the set of samples having a label $i$. We define $n_i = |\mathcal{N}_i|$ as the number of instances belonging to the $i$-th class in $\mathcal{S}$.

**Proposition 1** (**Unbiased Estimation**). *Define* $\hat{R}_{\ell,\mathcal{S}}(f)$ *as:*

$$\hat{R}_{\ell,\mathcal{S}}(f) = \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \ell^{i,j,m,n},$$

*where $\ell^{i,j,m,n}$ is a shorthand for $\ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n))$. Then $\hat{R}_{\ell,\mathcal{S}}(f)$ is an unbiased estimation of $R_\ell(f)$, in the sense that: $R_\ell(f) = \mathbb{E}_{\mathcal{S}}(\hat{R}_{\ell,\mathcal{S}}(f))$.*

According to the Empirical Risk Minimization (ERM) paradigm, we can turn to minimize $\hat{R}_{\ell,\mathcal{S}}(f)$ based on a finite training dataset $\mathcal{S}$ facing the unknown data distribution. Practically, the function $f$ here is often parameterized by a parameter $\theta$. In this way, we can express $f$ as $f_\theta$. For example, linear models could be defined as $f_\theta(\boldsymbol{x}) = \theta^\top \boldsymbol{x}$. Moreover, to prevent the over-fitting issue, we restrict our choice of $f$ on a specific hypothesis class $\mathcal{H}$. Again, $\mathcal{H}$ is essentially a subset of parameterized functions $f_\theta$ where $\theta$ is restricted in a set $\Theta$. For example, $\mathcal{H}$ could be defined as all the linear models with $||\theta|| \leq \gamma$. With this hypothesis set, we can construct the following optimization problem to learn $f$ within $\mathcal{H}$:

$$(\boldsymbol{OP_1}) \quad \min_{f \in \mathcal{H}} \hat{R}_{\ell,\mathcal{S}}(f). \tag{4}$$

From the parameter perspective, $(\boldsymbol{OP_1})$ could be reformulated equivalently by minimizing over $\theta \in \Theta$. Moreover, the constraint that $\theta \in \Theta$ could be reformulated as a regularization term $Reg_\Theta(\theta)$. Consequently, $(\boldsymbol{OP_1})$ also enjoys an alternative form that could be used for optimization:

$$(\boldsymbol{OP_2}) \quad \min_\theta \hat{R}_{\ell,\mathcal{S}}(f_\theta) + \alpha \cdot Reg_\Theta(\theta). \tag{5}$$

This allows the standard machine learning technologies to come into play. In this section, we will adopt the notations for mathematical convenience: $f$ and $\mathcal{H}$. Instead of explicitly defining $\theta$, we provide the parameterization in the definition of different hypothesis classes.

## 5.2 MAUC$^\downarrow$ **Rademacher Complexity and Its Properties.**

Our next step is then to investigate how well such approximations will perform. From one point, choosing a consistent surrogate loss $\ell$ ensures that minimizing the surrogate risk suffices to find the Bayes optimal scoring function. From another point, based on a small empirical risk, we also need to guarantee that *the training performance generalizes well to the unseen samples such that the expectation $R_\ell(f)$ is small*. To ensure the second point, given the assumption that $f$ is chosen from a hypothesis class $\mathcal{H}$, we will provide a Rademacher-Complexity-based worst-case analysis for the generalization ability in the remainder of this section. The key here is to ensure $R_\ell(f) \leq \hat{R}_\ell(f) + \delta$ with high probability, where $\delta$ goes to zero when $N$ goes to infinity. This way, minimizing $\hat{R}_\ell(f)$ suffices to minimize $R_\ell(f)$.

Unlike traditional machine learning problems that only involve independent instance-wise losses, AUC's pair-wise formulation makes the generalization analysis difficult to be carried out. Specifically, the terms in MAUC formulation exert certain degrees of interdependency. This way, the standard symmetrization scheme [5], [53] is not available for our problem. For instance, the terms $\ell(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2))$ and $\ell(f^{(i)}(\boldsymbol{x}_1') - f^{(i)}(\boldsymbol{x}_2'))$ are interdependent as long as $\boldsymbol{x}_1 = \boldsymbol{x}_1'$

or $\boldsymbol{x}_2 = \boldsymbol{x}_2'$. To address this issue, we provide an extended form of Rademacher complexity for MAUC$^\downarrow$ losses, which is defined as follows.

**Definition 2** (MAUC$^\downarrow$ **Rademacher Complexity**). *The Empirical MAUC$^\downarrow$ Rademacher Complexity over a dataset $\mathcal{S} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^m$, and a hypothesis space $\mathcal{H}$ is defined as:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H}) = \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{H}} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} T^{i,j,m,n} \right],$$

*where*

$$T^{i,j,m,n} = \frac{\sigma_m^{(i)} + \sigma_n^{(j)}}{2} \cdot \frac{\ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n))}{n_i n_j},$$

*for $i = 1, 2, \cdots, N_C$, $\sigma_1^{(i)}, \cdots, \sigma_{n_i}^{(i)}$ are i.i.d Rademacher random variables. The population version of the MAUC$^\downarrow$ Rademacher Complexity is defined as $\mathfrak{R}_{\mathsf{MAUC}^\downarrow}(\ell \circ \mathcal{H}) = \mathbb{E}_{\mathcal{S}} \left[ \hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H}) \right]$.*

The most important property of the MAUC$^\downarrow$ Rademacher complexity is that its magnitude is directly related to the generalization upper bounds, which is shown in the following theorem. The proofs are shown in Appendix.E in the supplementary materials.

**Theorem 4** (**Abstract Generalization Bounds**). *Given dataset $\mathcal{S} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^m$, where the instances are sampled independently, for all multiclass scoring functions $f \in \mathcal{H}$, if $Range\{\ell\} \subseteq [0, B]$, then for any $\delta \in (0, 1)$, the following inequality holds with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\mathcal{S}}(f) + C_1 \cdot \frac{\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H})}{N_C(N_C - 1)} +$$
$$C_2 \cdot \frac{B}{N_C} \cdot \xi(\boldsymbol{Y}) \cdot \sqrt{\frac{\log(2/\delta)}{N}},$$

*where $C_1, C_2$ are universal constants, $\xi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} 1/\rho_i}, \rho_i = \frac{n_i}{N}$.*

According to Thm.4, we can obtain a generalization bound as soon as we can obtain a proper upper bound on the empirical Rademacher complexity $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H})$. In the remainder of this subsection, we will develop general techniques to derive the upper bound of $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H})$ based on the notion of covering number and chaining. With the help of the next few theorems, we can convert the upper bound of $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H})$ to upper bounds of Rademacher complexities for much simpler model classes. These techniques are foundations for the practical results developed in the next subsection. Specifically, Thm.5-(a) is used to prove Thm.6; Thm.6 is employed in the second half of the proof of Thm.8; and Thm.5-(b) is employed in the proof of Thm.9.

We are now ready to present the corresponding results. With the sub-Gaussian property proved in Lem.9 in Appendix E, we can derive chaining upper bounds for the MAUC$^\downarrow$ Rademacher complexity. The foundation here is the notion of covering number, which is elaborated in Def.3 and Def.4.

**Definition 3** ($\epsilon$-covering). *[39] Let $(\mathcal{H}, d)$ be a (pseudo)metric space, and $\Theta \in \mathcal{H}$. $\{h_1, \cdots, h_K\}$ is said to be an $\epsilon$-covering of $\Theta$ if $\Theta \in \bigcup_{i=1}^K \mathcal{B}(h_i, \epsilon)$, i.e., $\forall \theta \in \Theta, \exists i \ s.t. \ d(\theta, h_i) \leq \epsilon$.*

**Definition 4** (Covering Number). *[39] Based on the notations in Def.3, the covering number of $\Theta$ with radius $\epsilon$ is defined as:*

$$\mathfrak{C}(\epsilon, \Theta, d) = \min\{n : \exists \epsilon - covering \ over \ \Theta \ with \ size \ n\}.$$

Based on the definitions above, we can reach the following results. The proof is shown in Appendix E in the supplementary materials. Note that in the following arguments, the covering number is defined on the metric $d_{\infty,\mathcal{S}}$. Specifically, given two vector valued functions $\boldsymbol{s} = (\boldsymbol{s}^{(1)}, \cdots, \boldsymbol{s}^{(N_C)})$, $\tilde{\boldsymbol{s}} = (\tilde{\boldsymbol{s}}^{(1)}, \cdots, \tilde{\boldsymbol{s}}^{(N_C)})$, and the training data $\mathcal{S}$, $d_{\infty,\mathcal{S}}(\boldsymbol{s}, \tilde{\boldsymbol{s}})$ is defined as:

$$d_{\infty,\mathcal{S}}(\boldsymbol{s}, \tilde{\boldsymbol{s}}) = \max_{\boldsymbol{x}_i \in \mathcal{S}, j \in [N_C]} |\boldsymbol{s}^{(j)}(\boldsymbol{z}) - \tilde{\boldsymbol{s}}^{(j)}(\boldsymbol{z})|. \quad (6)$$

**Theorem 5** (Chaining Bounds for MAUC$^\downarrow$ Rademacher Complexity). *Suppose that the score function $s^{(i)}$ maps $\mathcal{X}$ onto a bounded interval $[-R_s, R_s]$, the following properties hold for $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H})$:*

*(a) For a decreasing precision sequence $\{\epsilon_k\}_{k=1}^K$, with $\epsilon_{k+1} = \frac{1}{2}\epsilon_k$, $k = 1, 2, \cdots, K - 1$ and $\epsilon_0 \geq R_s$, we have:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H}) \leq N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_K$$
$$+ 6 \cdot \sum_{k=1}^K \epsilon_k \phi_\ell (N_C - 1) \cdot \xi(\boldsymbol{Y}) \sqrt{\frac{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}}$$

*(b) There exists a universal constant $C$, such that:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H}) \leq C\phi_\ell \inf_{R_s \geq \alpha \geq 0} \left( N_C(N_C - 1)\alpha \right.$$
$$\left. + (N_C - 1) \cdot \xi(\boldsymbol{Y}) \cdot \int_\alpha^{R_s} \sqrt{\frac{\log(\mathfrak{C}(\epsilon, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}} d\epsilon \right)$$

According to Thm.5, we have the following theorem as an extension of a new minorization technique which appears in a recent work [65].

**Theorem 6** (Transformation Upper Bound). *Given the Hypothesis class*

$$\mathsf{soft} \circ \mathcal{F} = \left\{ \boldsymbol{g}(\boldsymbol{x}) = \mathsf{soft}(\boldsymbol{s}(\boldsymbol{x})) : \ \boldsymbol{s} \in \mathcal{F} \right\},$$

*where $\mathsf{soft}(\cdot)$ is the softmax function. Suppose that $\boldsymbol{s}(x) \in [-R_s, R_s]^{N_C}$ and $\ell$ is $\phi_\ell$-Lipschitz continuous, the following inequality holds:*

$$\frac{\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{F})}{N_C(N_C - 1)} \leq \phi_\ell \left( 2^9 \cdot \frac{1}{N_C} \cdot \sqrt{N_C} \cdot \right.$$
$$\xi(\boldsymbol{Y}) \cdot \log^{3/2}(e \cdot R_s \cdot N \cdot N_C) \cdot \hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F})$$
$$\left. + \sqrt{\frac{1}{N}} \right),$$

*where the Rademacher complexity $\hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F})$ is defined as:*

$$\hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F})$$
$$= \mathbb{E}_\sigma \left[ \sup_{f = (f^{(1)}, \cdots, f^{(N_C)}) \in \mathcal{F}} \frac{1}{N \cdot N_C} \sum_{j=1}^{N_C} \sum_{i=1}^N \sigma_j^{(i)} \cdot f^{(j)}(\boldsymbol{x}_i) \right]$$

*where $\{\sigma_j^{(i)}\}_{(i,j)}$ is a sequence of independent Rademacher random variables.*

The result in Thm.6 directly relates the complicated pairwise Rademacher complexity $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{F})$ to the instance-wise Rademacher complexity $\hat{\mathfrak{R}}_{N \cdot N_C}$. This makes the derivation of the generalization upper bound much easier. Once a bound on the ordinary Rademacher complexity over the functional class $\mathcal{F}$ is available, we can directly plug it into this theorem and find a resulting bound over the MAUC$^\downarrow$ complexity.

### 5.3 Generalization Bounds for Deep and Shallow Model Families

Next, we derive the generalization bounds for three hypothesis classes: (1) the $\ell_p$ penalized linear models, (2) deep neural networks with fully-connected layers, and (3) the deep convolutional neural networks.

**Assumption 1** (Common Assumptions). *In this subsection, we require some common assumptions listed as follows:*
- *The sample points in the training dataset $\mathcal{S} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^m$ are sampled independently*
- *$Range\{\ell\} \subseteq [0, B]$*
- *The loss function $\ell$ is $\phi_\ell$-Lipschitz continuous*
- *The input features are sampled from $\mathcal{X} \subset \mathbb{R}^d$, and for all $\boldsymbol{x} \in \mathcal{X}$, we have $||\boldsymbol{x}||_2^2 \leq R_\mathcal{X}$.*

#### 5.3.1 Generalization Bound for $\ell_p$ Penalized Linear Models

We begin with the generalization bound for $\ell_p$ norm penalized linear models. The proof is shown in Appendix F in the supplementary materials.

**Theorem 7** (**Practical Generalization Bounds for Linear Models**). *Define the $\ell_p$ norm penalized linear model as:*

$$\mathcal{H}_{p,\gamma}^{Lin} = \left\{ f = (f^{(1)}, \cdots, f^{(n_C)}) : f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^{(i)}\boldsymbol{x}, \right.$$
$$\left. ||\boldsymbol{W}^{(i)}||_p \leq \gamma \right\},$$

*with $0 < p < \infty$, $\frac{1}{p} + \frac{1}{\bar{p}} = 1$. Based on assumption 1, for all $f \in \mathcal{H}_{p,\gamma}^{Lin}$, we have the following inequality holds with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\ell,\mathcal{S}}(f) + \mathcal{I}_{Lin}\left( \chi(\boldsymbol{Y}), \xi(\boldsymbol{Y}), \delta \right) \cdot \sqrt{\frac{1}{N}},$$

*where*

$$\mathcal{I}_{Lin}\left( \chi(\boldsymbol{Y}), \xi(\boldsymbol{Y}), \delta \right) = \frac{4R_\mathcal{X}\phi_\ell\gamma}{N_C - 1} \cdot \sqrt{\frac{2(\bar{p} - 1)}{N_C}} \cdot \chi(\boldsymbol{Y})$$
$$+ \frac{5B}{N_C} \cdot \sqrt{2\log\left(\frac{2}{\delta}\right)} \cdot \xi(\boldsymbol{Y}).$$

#### 5.3.2 Generalization Bound for Deep Fully-Connected Neural Networks

Next, we take a further step to explore the generalization ability of a specific type of deep neural networks where only fully-connected layers and activation functions exist. The detailed here is as follows.

**Settings**. We denote an $L$-layer deep fully-connected neural network with $N_C$-way output as :

$$f(\boldsymbol{x}) = \boldsymbol{W} f_{\boldsymbol{\omega}, L}(\boldsymbol{x}) = \boldsymbol{W} s(\boldsymbol{\omega}_{L-2} \cdots s(\boldsymbol{\omega}_1 \boldsymbol{x})),$$

The notations are as follows: $s(\cdot)$ is the activation function; $n_{h_j}$ is the number of hidden neurons for the $j$-th layer; $\boldsymbol{\omega}_j \in \mathbb{R}^{n_{h_{j+1}} \times n_{h_j}}, j = 1, 2, \cdots, L - 2$ are the weights for the first $L - 1$ layers; $\boldsymbol{W} \in \mathbb{R}^{n_{h_{L-1}} \times N_C}$ is the weight for the output layer. Moreover, the output from the $i$-th layer of the network is defined as :

$$f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^{(i)^\top} f_{\boldsymbol{\omega},L}(\boldsymbol{x}),$$

where $\boldsymbol{W}^{(i)^\top}$ is the $i$-th row of $\boldsymbol{W}$. In the next theorem, we focus on a specific hypothesis class for such networks where the product of weight norms

$$\Pi_{\boldsymbol{W},\boldsymbol{\omega}} = ||\boldsymbol{W}||_F \cdot \prod_{j=1}^{L-2} ||\boldsymbol{\omega}_j||_F$$

are no more than $\gamma$. We denote such a hypothesis class as $\mathcal{H}_{\gamma,R_s,n_h}^{DNN}$:

$$\mathcal{H}_{\gamma,R_s,n_h}^{DNN} = \left\{ f : f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^{(i)^\top} f_{\boldsymbol{\omega},L}(\boldsymbol{x}), ||f^{(i)}||_\infty \leq R_s, \right.$$
$$\left. i = 1, \cdots, N_C, \; \Pi_{\boldsymbol{W},\boldsymbol{\omega}} \leq \gamma \right\},$$

To obtain the final output, we perform a softmax operation over $f(\boldsymbol{x})$. Thus, the valid model under this setting could be chosen from the following hypothesis class:

$$\text{soft} \circ \mathcal{H}_{\gamma,R_s,n_h}^{DNN} = \left\{ g : g^{(i)} = \frac{\exp(f^{(i)}(\boldsymbol{x}))}{\sum_{j=1}^{N_C} \exp(f^{(j)}(\boldsymbol{x}))}, \right.$$
$$\left. f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN} \right\}.$$

We have the following bound for this type of models. The result here is a merge of two independent results we proposed in Appendix E and Appendix F in the supplementary material, which is based on the Talagrand contraction properties and the chaining technology shown in Thm.5 and Thm.6.

**Theorem 8 (Practical Generalization Bounds for Deep Models).** *On top of assumption 1, if we further assume that $s(\cdot)$ is a 1-Lipshitz and positive homogeneous activation function, then for all $f \in \text{soft} \circ \mathcal{H}_{\gamma,R_s,n_h}^{DNN}$, we have the following inequality holds with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_\mathcal{S}(f) + \min\left( \mathcal{I}_{DNN,1}, \; \mathcal{I}_{DNN,2} \right) \cdot \sqrt{\frac{1}{N}}$$

*where $\chi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} \sum_{j \neq i} 1/\rho_i \rho_j}$, $\xi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} 1/\rho_i}$, $\rho_i = \frac{n_i}{N}$,*

$$\mathcal{I}_{DNN,1} = C_1 \frac{\sqrt{2}}{2} \phi_\ell \cdot \frac{\chi(\boldsymbol{Y})}{N_C - 1} + \left( \frac{\sqrt{2} C_1 R_\mathcal{X} \phi_\ell \gamma}{2} \cdot C_3 \right.$$
$$\left. + \frac{C_2 B}{N_C} \cdot \sqrt{2 \log\left(\frac{2}{\delta}\right)} \right) \cdot \xi(\boldsymbol{Y}),$$

$$\mathcal{I}_{DNN,2} = C_1 \phi_\ell \left( \frac{2^9}{N_C} \cdot \xi(\boldsymbol{Y}) \cdot \log^{3/2}(K \cdot N \cdot N_C) \right.$$
$$\gamma \cdot R_\mathcal{X} \cdot (\sqrt{2 \log(2) L} + 1) + 1 \Big)$$
$$+ C_2 \frac{B \cdot \sqrt{\log(2/\delta)} \cdot \xi(\boldsymbol{Y})}{N_C},$$

$C_1, C_2$ *are universal constants as thm.4, $K = e \cdot R_s$, $C_3 = \frac{\sqrt{L \log 2} + \sqrt{N_C}}{\sqrt{N_C - 1}}$.*

### 5.3.3 Generalization Bound for Deep Convolutional Neural Networks

Now we use the result in Thm.5 to derive a generalization bound for a class of deep neural networks where fully-connected layers and convolutional layers coexist. In a nutshell, the result is essentially an application of Thm.5 to a recent idea appeared in [49].

**Settings.** Now we are ready to introduce the setting of the deep neural networks employed in the forthcoming theoretical analysis, which is adopted from [49]. We focus on the deep neural networks with $N_{conn}$ fully-connected layers and $N_{conv}$ convolutional layers. The $i$-th convolutional layer has a kernel $\boldsymbol{K}^{(i)} \in \mathbb{R}^{k_i \times k_i \times c_{i-1} \times c_i}$. Recall that convolution is a linear operator. For a given kernel $\boldsymbol{K}$, we denote its associated matrix as $op(\boldsymbol{K})$, such that $\boldsymbol{K}(\boldsymbol{x}) = op(\boldsymbol{K})\boldsymbol{x}$. Moreover, we assume that, each time, the convolution layer is followed by a componentwise non-linear activation function and an optional pooling operation. We assume that the activation functions and the pooling operations are all 1-Lipschitz. For the $i$-th fully-connected layer, we denote its weight as $\boldsymbol{V}^{(i)}$. Above all, the complete parameter set of a given deep neural network could be represented as $\boldsymbol{P} = \{\boldsymbol{K}^{(1)}, \cdots, \boldsymbol{K}^{(N_{conv})}, \boldsymbol{V}^{(1)}, \cdots, \boldsymbol{V}^{(N_{conn})}\}$. Again, we also assume that the loss function is $\phi_\ell$-Lipschitz and $Range\{\ell\} \subseteq [0, B]$. Finally, given two deep neural networks with parameters $\boldsymbol{P}$ and $\tilde{\boldsymbol{P}}$, we adopt a metric $d_{NN}(\cdot, \cdot)$ to measure their distance:

$$d_{NN}(\boldsymbol{P}, \tilde{\boldsymbol{P}}) = \sum_{i=1}^{N_{conv}} ||op(\boldsymbol{K}^{(i)}) - op(\tilde{\boldsymbol{K}}^{(i)})||_2$$
$$+ \sum_{i=1}^{N_{conn}} ||\boldsymbol{V}^{(i)} - \tilde{\boldsymbol{V}}^{(i)}||_2.$$

**Constraints Over the Parameters**. First, we define $\mathcal{P}_\nu^{(0)}$ as the class for initialization of the parameters:

$$\mathcal{P}_\nu^{(0)} = \left\{ \boldsymbol{P} : \left( \max_{i \in \{1, \cdots, N_{conv}\}} ||op(\boldsymbol{K}^{(i)})||_2 \right) \right.$$
$$\left. \leq 1 + \nu, \; \left( \max_{j \in \{1, \cdots, N_{conn}\}} ||\boldsymbol{V}^{(j)}||_2 \right) \leq 1 + \nu \right\}.$$

Now we further assume that the learned parameters should be chosen from a class denoted by $\mathcal{P}_{\beta,\nu}$, where the distance between the learned parameter and the fixed initialization residing in $\mathcal{P}_\nu^{(0)}$ is no bigger than $\beta$:

$$\mathcal{P}_{\beta,\nu} = \left\{ \boldsymbol{P} : d_{NN}(\boldsymbol{P}, \tilde{\boldsymbol{P}}_0) \leq \beta, \; \tilde{\boldsymbol{P}}_0 \in \mathcal{P}_\nu^{(0)} \right\}.$$

We have the following result for this class of deep neural networks, which is an application of Thm.5. The proof is shown in Appendix F in the supplementary materials.

**Theorem 9.** *Denote the hypothesis class,*

$$\text{soft} \circ \mathcal{F}_{\beta,\nu} = \left\{ g(\boldsymbol{x}) = \text{soft}(s_{\boldsymbol{P}}(\boldsymbol{x})) : \; s_{\boldsymbol{P}} \in \mathcal{F}_{\beta,\nu} \right\},$$

$$\mathcal{F}_{\beta,\nu} = \{s_{\boldsymbol{P}} : \mathbb{R}^{N N_L - 1} \to \mathbb{R}^{N_C} \mid \boldsymbol{P} \in \mathcal{P}_{\beta,\nu},$$
$$Range(\boldsymbol{s}_P) \subseteq [-R_s, R_s]^{N_C} \}.$$

*On top of Assumption 1, if we further assume that*

$$R_s > 1/\min\left\{ \sqrt{N}, \right.$$
$$\left. \frac{\xi(\boldsymbol{Y})}{N_C} \cdot \sqrt{N_{par}\left(\nu N_L + \beta + \log(3R_{\mathcal{X}} \cdot \beta \cdot N)\right)} \right\},$$

*then for all multiclass scoring functions $f \in \mathsf{soft} \circ \mathcal{F}_{\beta,\nu}$, $\forall \delta \in (0, 1)$, the following inequalities hold with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\mathcal{S}}(f) + C_1 \alpha_1 \alpha_2$$
$$+ C_2 \cdot \frac{B\xi(\boldsymbol{Y})}{N_C} \cdot \sqrt{\frac{\log(2/\delta)}{N}},$$

*where*

$$\alpha_1 = \tilde{C} \cdot \frac{\phi_\ell \cdot R_s \cdot \xi(\boldsymbol{Y})}{N_C},$$
$$\alpha_2 = \sqrt{\frac{N_{par}\left(\nu N_L + \beta + \log\left(3\beta R_{\mathcal{X}} N\right)\right)}{N}},$$

*$\tilde{C}, C_1, C_2$ are universal constants, $N_L = N_{conv} + N_{conn}$, $N_{par}$ is total number of parameters in the neural network.*

### 5.4 Summary

Note that there are two imbalance-aware factors that appear in the upper bounds above. The first factor $\xi(\boldsymbol{Y})$ captures the imbalance of class label distribution in $\mathcal{S}$. With an analogous spirit, $\chi(\boldsymbol{Y})$ captures the degree of imbalance of the class pair distribution in $\mathcal{S}$. To improve generalization, the training data $\mathcal{S}$ must simultaneously have a large sample size $N$ and a small degree of imbalance captured by $(\xi(\boldsymbol{Y}), \chi(\boldsymbol{Y}))$. In other words, our bounds suggest that blindly increasing the sample size of the training dataset will not improve generalization. To really improve generalization, one should increase data points of the minority classes which are the sources for performance bottleneck. Consequently, compared with the ordinary $O(\sqrt{1/N})$ result, our bound is much more aware of the imbalance issue hidden behind the training data.

## 6 EFFICIENT COMPUTATIONS IN OPTIMIZATION

Till now, we have conducted a series of theoretical analyses for our framework. In this section, we turn our focus to the practical issues we must face during optimization. As shown in the previous sections, the complicated formulations of MAUC$^\downarrow$ surrogate losses bring a great burden to the fundamental operations of the downstream optimization algorithm. Specifically, given $T_\ell$ as the time complexity for loss evaluation of a single sample pair and given $T_{grad}$ as that for the gradient evaluation, we can find that even a single full(mini)-batch loss and gradient evaluation takes $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j T_\ell)$ and $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j T_{grad})$, which scales almost quadratically to the sample (batch) size. In general, such complexities are unaffordable facing medium and large-scale datasets. However, we find that for some of the popular surrogate losses, the pairwise computation could be largely simplified. Consequently, we propose acceleration

algorithms for loss and gradient evaluations for three well-known losses: exponential loss, squared loss, and hinge loss. The time complexity with/without acceleration is shown in Tab.1, which shows the effectiveness of our proposed algorithms.

Generally, the empirical surrogate risk functions $\hat{R}_\ell$ have the following abstract form:

$$\hat{R}_\ell = \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \ell^{i,j},$$

where:

$$\ell^{i,j} = \ell\left(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)\right),$$

$$f^{(i)}(\boldsymbol{x}) = g_i(\boldsymbol{W} h_{\boldsymbol{\theta}}(\boldsymbol{x})),$$

$$\boldsymbol{W} = [\boldsymbol{w}^{(1)}, \cdots, \boldsymbol{w}^{(N_C)}]^\top.$$

We adopt this general form since it covers a lot of popular models. For example, when $g_i(\cdot)$ is defined as the activation function of the last layer of a neural network (say a softmax function), $\boldsymbol{w}^{(i)}$ are the weights of the last layer, and $h_{\boldsymbol{\theta}}(\cdot)$ is the neural net where the last layer is excluded, $f^{(i)}(\boldsymbol{x})$ becomes a deep neural network architecture with the last layer designed as a fully-connected layer. As an another instance, if both $g_i(\boldsymbol{x}) = \boldsymbol{x}$ and $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{x}$, then we reach a simple linear multiclass scoring function. Note that the scalability with respect to sample size only depends on the choice of $\ell$. The choice of $g_i(\cdot)$ and $h_{\boldsymbol{\theta}}(\cdot)$ only affects the instance-wise chain rule. This allows us to provide a general acceleration framework once the surrogate loss is fixed.

*Due to the limited space of this paper, we only provide loss evaluation acceleration methods in the main paper. The readers are referred to Appendix G for more details, where we also present a discussion to deal with the acceleration for general loss functions.*

### 6.1 Exponential Loss

For the exponential loss, we can simplify the computations by a factorization scheme:

$$\hat{R}_{exp} = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \exp\left(\alpha \cdot \left(f^{i,j,m,n}\right)\right),$$
$$= \sum_{i=1}^{N_C} \underbrace{\left(\sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \exp(\alpha \cdot f^{(i)}(\boldsymbol{x}_m))\right)}_{(a_i)} \cdot$$
$$\underbrace{\left(\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \exp(-\alpha \cdot f^{(i)}(\boldsymbol{x}_n))\right)}_{(b_i)},$$

where $f^{i,j,m,n} = f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)$. From the derivation above, the loss evaluation could be done by first calculating $(a_i), (b_i)$ separately and then performing the multiplication, which only takes $O(N N_C T_\ell)$. This is a significant improvement compared with the original $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j N_C T_\ell)$ result.

TABLE 1
Acceleration for three losses , where $\bar{N} = \sum_{i=1}^{N_C} n_i \log n_i + (N - n_i) \log(N - n_i)$

| Algorithms | loss | gradient | requirement |
|---|---|---|---|
| exp + acceleration | $O(N_C \cdot N \cdot T_\ell)$ | $O(N_C \cdot N \cdot T_{grad})$ | $\min_i n_i \gg 2$ |
| squared + acceleration | $O(N_C \cdot N \cdot T_\ell)$ | $O(N_C \cdot N \cdot T_{grad})$ | $e^{\frac{1}{2}(N-n_i)} \gg n_i \gg N - e^{\frac{1}{2}n_i}$ |
| hinge + acceleration | $O(N_C \cdot \bar{N} \cdot T_\ell)$ | $O(N_C \cdot \bar{N} \cdot T_{grad})$ | $\min_i n_i \gg 2$ |
| w/o acceleration | $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j \cdot T_\ell)$ | $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j \cdot T_{grad})$ | \ |

## 6.2 Hinge loss

First we put down the hinge surrogate loss as:

$$\hat{R}_{hinge} = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \left( \alpha - \left( f^{i,j,m,n} \right) \right)_+ .$$

The key of our acceleration is to notice that the terms are non-zero only if $\left( f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n) \right) \leq \alpha$. Moreover, for these non-zero terms $\max(x, 0) = x$. This means that the hinge loss degenerates to an identity function for the activated nonzero terms, which enjoys efficient computation. So the key step in our algorithm is to find out the non-zero terms in an efficient manner.

Given a fixed class $i$ and an instance $\boldsymbol{x}_m \in \mathcal{N}_i$, we denote:

$$\mathcal{A}^{(i)}(\boldsymbol{x}_m) = \{ \boldsymbol{x}_n \notin \mathcal{N}_i, \alpha > f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n) \}.$$

With $\mathcal{A}^{(i)}(\boldsymbol{x}_m)$, one can reformulate $\hat{R}_{hinge}$ as:

$$\hat{R}_{hinge} = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \Bigg[ \left( \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} \right) \cdot$$
$$\left( \alpha - f^{(i)}(\boldsymbol{x}_m) \right) + \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} \cdot f^{(i)}(\boldsymbol{x}_n) \Bigg],$$
$$= \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \left[ \delta^{(i)}(\boldsymbol{x}_m) \cdot (\alpha - f^{(i)}(\boldsymbol{x}_m)) + \Delta^{(i)}(\boldsymbol{x}_m) \right],$$

where

$$\delta^{(i)}(\boldsymbol{x}_m) = \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j},$$
$$\Delta^{(i)}(\boldsymbol{x}_m) = \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} f^{(i)}(\boldsymbol{x}_n). \qquad (7)$$

According to this reformulation, once $\delta^{(i)}(\boldsymbol{x})$ and $\Delta^{(i)}(\boldsymbol{x})$ are all fixed, we can calculate the loss function within $O(N)$. So the efficiency bottleneck comes from the calculations of $\delta^{(i)}(\boldsymbol{x})$ and $\Delta^{(i)}(\boldsymbol{x})$.

Next we propose an efficient way to calculate $\delta^{(i)}(\boldsymbol{x}_m), \Delta^{(i)}(\boldsymbol{x}_m), \mathcal{A}^{(i)}(\boldsymbol{x}_m)$. To do so, given a specific class $i$, we first sort the relevant (resp. irrelevant) instances descendingly according to their scores $f^{(i)}(\boldsymbol{x})$:

$$f^{(i)}(\boldsymbol{x}_0^\downarrow) \geq f^{(i)}(\boldsymbol{x}_1^\downarrow) \cdots, \geq f^{(i)}(\boldsymbol{x}_{n_i-1}^\downarrow), \quad \boldsymbol{x}_i^\downarrow \in \mathcal{N}_i,$$
$$f^{(i)}(\tilde{\boldsymbol{x}}_0^\downarrow) \geq f^{(i)}(\tilde{\boldsymbol{x}}_1^\downarrow) \cdots, \geq f^{(i)}(\tilde{\boldsymbol{x}}_{N-n_i-1}^\downarrow), \quad \tilde{\boldsymbol{x}}_i^\downarrow \notin \mathcal{N}_i.$$

It immediately follows that:

$$\mathcal{A}^{(i)}(\boldsymbol{x}_0^\downarrow) \subseteq \mathcal{A}^{(i)}(\boldsymbol{x}_1^\downarrow) \subseteq \cdots \subseteq \mathcal{A}^{(i)}(\boldsymbol{x}_{n_i}^\downarrow).$$

This allows us to find out all $\mathcal{A}^{(i)}(\boldsymbol{x}_k^\downarrow), k = 1, 2, \cdots, n_i$, within a single pass of the whole dataset with an efficient dynamic programming.

Based on the construction of $\mathcal{A}_k^{(i)} = \mathcal{A}^{(i)}(\boldsymbol{x}_k^\downarrow)$, we provide the following recursive rules to obtain $\delta^{(i)}(\boldsymbol{x})$ and $\Delta^{(i)}(\boldsymbol{x})$:

$$\delta^{(i)}(\boldsymbol{x}_{k+1}^\downarrow) = \delta^{(i)}(\boldsymbol{x}_k^\downarrow) + \sum_{j \neq i} \sum_{\mathcal{N}_j \cap \mathcal{A}_{k+1}^{(i)} \setminus \mathcal{A}_k^{(i)}} \frac{1}{n_i n_j},$$
$$\Delta^{(i)}(\boldsymbol{x}_{k+1}^\downarrow) = \Delta^{(i)}(\boldsymbol{x}_k^\downarrow) +$$
$$\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}_{k+1}^{(i)} \setminus \mathcal{A}_k^{(i)}} \frac{1}{n_i n_j} f^{(i)}(\boldsymbol{x}_n).$$

These recursive rules induce an efficient dynamic programming which only requires a single pass of the whole dataset. Putting all together, we then come to a $O(N_C \cdot \bar{N} \cdot T_l)$ time complexity speed-up for hinge loss evaluation. *An implementation of this idea is detailed in the Appendix G.*

## 6.3 Squared Loss

For each fixed class $i$, we construct an affinity matrix $\boldsymbol{Aff}^{(i)}$ such that

$$\boldsymbol{Aff}_{m,n}^{(i)} = \begin{cases} \dfrac{1}{n_i n_{y_m}}, & n \in \mathcal{N}_i, m \notin \mathcal{N}_i, \\[2mm] \dfrac{1}{n_i n_{y_n}}, & m \in \mathcal{N}_i, n \notin \mathcal{N}_i, \\[2mm] 0, & \text{Otherwise,} \end{cases}$$

which could be written in a matrix form,

$$\boldsymbol{Aff}^{(i)} = \boldsymbol{D}^{(i)}(\boldsymbol{1} - \mathbf{Y}^{(i)})\mathbf{Y}^{(i)\top} + \mathbf{Y}^{(i)}(\boldsymbol{1} - \mathbf{Y}^{(i)})^\top \boldsymbol{D}^{(i)},$$

Then the empirical risk function under the squared surrogate loss could be reformulated as:

$$\hat{R}_{sq} = \sum_{i=1}^{N_C} \boldsymbol{\Delta}_{sq}^{(i)\top} \mathcal{L}^{(i)} \boldsymbol{\Delta}_{sq}^{(i)},$$

where

$$\boldsymbol{\Delta}_{sq}^{(i)} = \mathbf{Y}^{(i)} - f^{(i)}(\boldsymbol{X}),$$
$$f^{(i)}(\boldsymbol{X}) = [f^{(i)}(\boldsymbol{x}_1), \cdots, f^{(i)}(\boldsymbol{x}_N)]^\top, \qquad (8)$$

$\mathcal{L}^{(i)} = diag(\boldsymbol{Aff}^{(i)}\boldsymbol{1}) - \boldsymbol{Aff}^{(i)}$ is the graph Laplacian matrix associated with $\boldsymbol{Aff}^{(i)}$. Based on this reformulation, one can find that the structure of $\mathcal{L}^{(i)}$ provides an efficient factorization scheme to speed-up loss evaluation. To see this, we have:

$$\hat{R}_{sq} = \sum_{i=1}^{N_C} \frac{1}{2} \boldsymbol{\Delta}_{sq}^{(i)\top} (\boldsymbol{\kappa}^{(i)} \odot \boldsymbol{\Delta}_{sq}^{(i)}) - \Delta_1^{(i)} \cdot \Delta_2^{(i)},$$

where

$$\Delta_2^{(i)} = \mathbf{Y}^{(i)\top} \boldsymbol{\Delta}_{sq}^{(i)}$$
$$\Delta_1^{(i)} = \boldsymbol{\Delta}_{sq}^{(i)\top} \boldsymbol{D}^{(i)}(\boldsymbol{1} - \mathbf{Y}^{(i)})$$
$$\boldsymbol{\kappa}^{(i)} = n_i \boldsymbol{D}^{(i)}(\boldsymbol{1} - \mathbf{Y}^{(i)}) + \frac{N_C - 1}{n_i} \mathbf{Y}^{(i)}.$$

Taking all together, we see that the accelerated evaluation requires only $O(N \cdot N_C \cdot T_l)$ time.

# 7 EXPERIMENTS

## 7.1 Datasets

Now we start the empirical analyses on the real-world datasets. First, we describe all the datasets involved in the forthcoming discussions. Generally speaking, the datasets come from three types of sources: (a) LIBSVM website, (b) KEEL website, and (c) others. Note that for all the datasets prefixed with Imb, we use an imbalanced subset sampled from the original dataset to perform our experiments. Here we only present a brief summary. More detailed contents could be found in Appendix H.2.

(a) **LIBSVM Datasets** [2], which includes: **Shuttle**, **Svmguide-2**, **SegmentImb**. [3]

(b) **KEEL Datasets**. [4], which includes: **Balance**, **Dermatology**, **Ecoli**, **New Thyroid**, **Page Blocks**, **Yeast**.

(c) **Other Datasets**:

(1) **CIFAR-100-Imb**. We sampled an imbalanced version of the CIFAR-100 [5] dataset, which originally contains 100 image classes each with 600 instances. The 100 classes are encoded as $1, \cdots 100$.

(2) **User-Imb**. The original dataset is collected from TalkingData, a famous third-party mobile data platform from China, which predicts mobile users' demographic characteristics based on their app usage records. The dataset is collected for the Kaggle Competition named Talking Data Mobile User Demographics. [6] The raw features include logged events, app attributes, and device information. There are 12 target classes 'F23-', 'F24-26','F27-28','F29-32', 'F33-42', 'F43+', 'M22-', 'M23-26', 'M27-28', 'M29-31', 'M32-38', 'M39+', which describe the demographics (gender and age) of users. In our experiments, we sample an imbalanced subset of the original dataset to leverage a class-skewed dataset.

(3) **iNaturalist2017.** iNaturalist Challenge 2017 dataset[7] is a large-scale image classification benchmark with 675,170 images covering 5,089 different species of plants and animals. We split the dataset into the training set, validation set, and test set at a ratio of 0.7:0.15:0.15. Since directly training deep models in such a large-scale dataset is time-consuming, we instead generate 2048-d features with a ResNet-50 model pre-trained on ImageNet for each image and train models with three fully-connected layers for all methods. We utilize Adam optimizer to train the models, with an initial learning rate of $10^{-5}$. To ensure all categories are covered in a mini-batch, the batch size is set to 8196. Other hyperparameters are the same as those in the CIFAR-100-Imb dataset.

## 7.2 Competitors

**Choice of Competitors**. Specifically, our goal is to validate that direct MAUC optimization techniques leverage better MAUC performance than other imbalanced learning methods. Motivated by this, the competitors include: a **Baseline** that does not deal with the imbalance distribution; **Over-sampling and Under-sampling Methods** that tackle imbalance issues by re-sampling; and **Imbalanced Loss Functions** that tackle imbalance issues by reformulating the loss function. Moreover, we implement our framework with the squared loss, exponential loss, and hinge loss with the acceleration methods proposed in Sec.6. Note that, for the sampling-based methods, we use the code from the python Lib `Imbalanced-learn` [40] [8] to implement these competitors. The details are listed as follows:

1) **Standard Baseline: LR**. This is the baseline model where no measures are taken against imbalance issues. For **LR**, we adopt a multiclass version of the logistic regression. For traditional datasets, **LR** is built upon a linear model. For deep learning datasets, **LR** is built upon a common deep backbone.

2) **Over-sampling Methods**. For these competitors, we first plug in an oversampling method to generate a more balanced dataset, then we use **LR** to train a model on the new dataset. Here we adopt the following two oversampling methods as our competitors:

   - **BM**. (BorderlineSMOTE) [29]: This method is a variant of the SMOTE method (Synthetic Minority Oversampling Technique), which restricts the oversampling process to the hard minority samples which lie at the borderline of the decision boundary.

   - **MM**. (MWMOTE) [6]: MWMOTE is another variant of SMOTE. It first identifies the hard-to-learn informative minority class samples and assigns them weights according to their Euclidean distance from the nearest majority class samples.

3) **Under-sampling Methods**. For these competitors, we first adopt an undersampling method to generate a more balanced dataset, then we use **LR** to train a model on the new dataset. Here we adopt the following two methods as our competitors:

   - **IHT**. (InstanceHardnessThreshold) [67]. This method works by filtering out the hard and noisy samples from the majority classes based on the instance hardness measure proposed in [67].

   - **NM**. (NearMiss) [50]: It adopts an under-sampling idea to make majority class samples surround most of the minority class samples.

   - **TL**. (TomekLinks) [68] It adopts the Tomek Links method to remove redundant samples that fail to contribute to the outline of the decision boundary from the majority class.

4) **Imbalanced Loss Functions** (for deep learning): For CIFAR-100-Imb and User-Imb dataset, we also compare our method with some of the recently proposed imbalanced loss functions for deep learning.

   - **Focal Loss** [42]: It tackles the imbalance problem by

---

2. https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/

3. The goal of the original dataset is to predict image segmentation based on 19 hand-crafted features. To obtain an imbalanced version of the dataset, we sample 17, 33, 26, 13, 264, 231, 165 instances respectively from each of the 7 classes.

4. http://www.keel.es/

5. https://www.cs.toronto.edu/~kriz/cifar.html

6. https://www.kaggle.com/c/talkingdata-mobile-user-demographics/data

7. https://www.kaggle.com/c/inaturalist-challenge-at-fgvc-2017

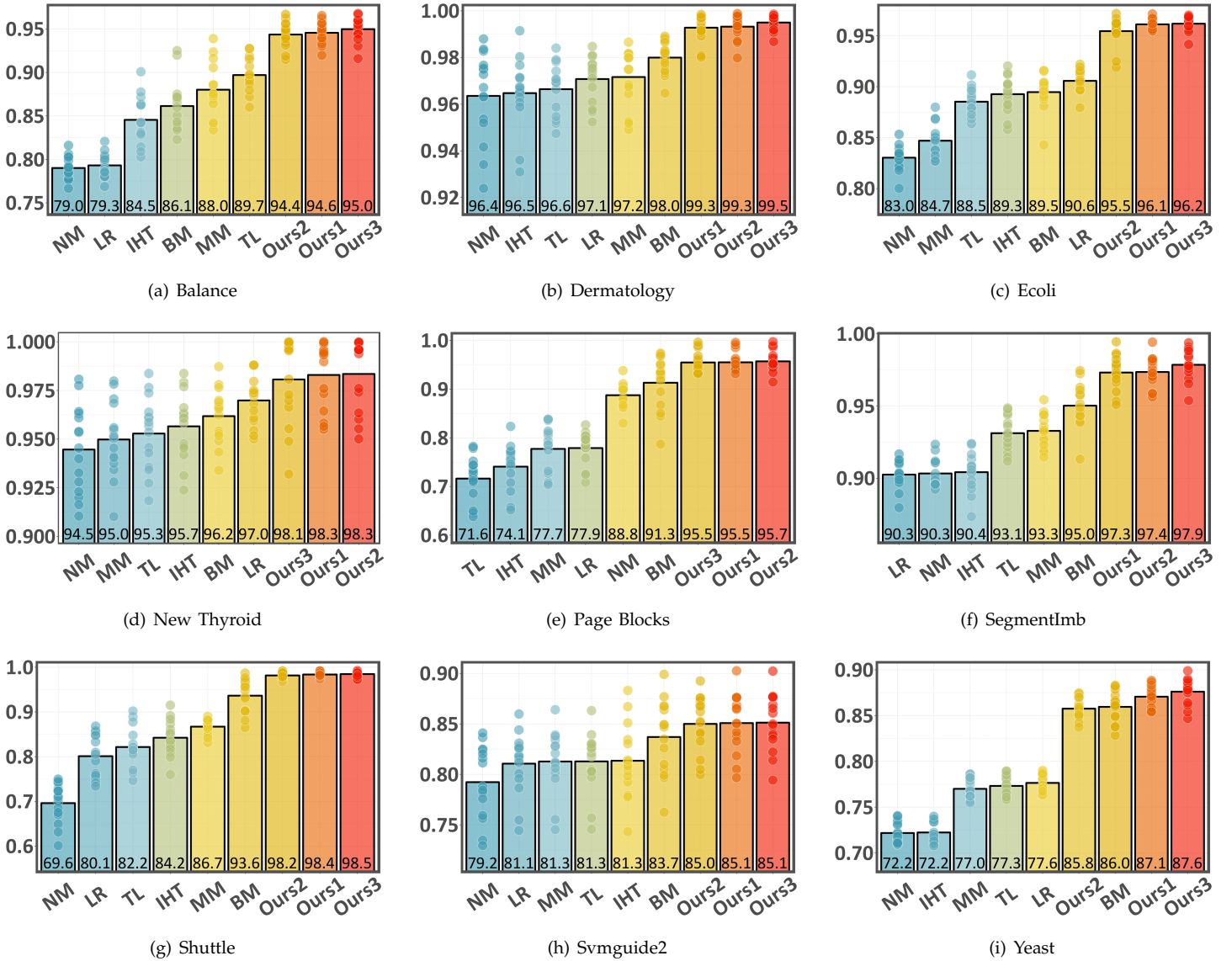8. https://imbalanced-learn.readthedocs.io/en/stable/#

Fig. 1. **Coarse-grained Performance Comparison.** We plot the overall MAUC$^\uparrow$ over the 15 splits against different algorithms. Here each bar presents the performance of a specific algorithm. The bar itself captures the mean performance, and the scatters distributed along a bar are the results over 15 splits.

adding a modulating factor to the cross-entropy loss to highlight the hard and minority samples during the training process.

- **CB-CE**: It refers to the loss function that applies the reweighting scheme proposed in [16] on the cross-entropy loss.
- **CB-Focal**: It refers to the loss function that applies the reweighting scheme proposed in [16] on the Focal loss.
- **LDAM** [9]: It proposes a Label-distribution-aware margin loss based on the minimum margin per class. Here we adopt the smooth relaxation for cross-entropy loss proposed therein.

5) Existing Deep AUC optimization methods:
- **DeepAUC:** [45] is a state-of-the-art stochastic AUC maximization algorithm developed for the deep neural network. It solves the AUC maximization from the saddle point problem. To ensure a fair comparison,

we implement the algorithm that extends to the multi-class AUC problem by PyTorch. In addition, there are two optimization configurations in [45], including Proximal Primal-Dual Stochastic Gradient (PPD-SG) and Proximal Primal-Dual AdaGrad (PPD-AdaGrad). PPD-AdaGrad is employed in our experiment because it usually demonstrates better performance than PPD-SG in most cases.

6) Our Methods:
- **Ours1**: An implementation of our learning framework with the square surrogate loss $\ell_{sq}(\alpha, t) = (\alpha - t)^2$
- **Ours2**: An implementation of our learning framework with the exponential loss $\ell_{\exp}(\alpha, t) = \exp(-\alpha t)$
- **Ours3**: An implementation of our learning framework with the hinge loss $\ell_{hinge}(\alpha, t) = (\alpha - t)_+$

## 7.3 Implementation details

**Infrastructure**. All the experiments are carried out on a ubuntu 16.04.6 server equipped with Intel(R) Xeon(R) CPU E5-2620 v4 cpu and a TITAN RTX GPU. The codes are implemented via `python 3.6.7`, the basic dependencies are: `pytorch` (v-1.1.0), `sklearn` (v-0.21.3), `numpy` (v-1.16.2). For traditional datasets, we implement our proposed algorithms with the help of the `sklearn` and `numpy`. For hinge loss, we use Cython to accelerate the dynamic programming algorithm. For the deep learning datasets, our proposed algorithms are implemented with `pytorch`.

**Evaluation Metric**. Given a trained scoring function $f = (f^{(1)}, \cdots, f^{(N_C)})$, all the forthcoming results are evaluated with the MAUC metric (the larger the better).

$$\mathrm{MAUC}^{\uparrow} = \frac{1}{N_C \cdot (N_C - 1)} \cdot$$
$$\sum_{i=1}^{N_C} \sum_{j \neq i} \frac{\left| \left\{ (\boldsymbol{x}_1, \boldsymbol{x}_2) : \boldsymbol{x}_1 \in \mathcal{N}_i, \ \boldsymbol{x}_2 \in \mathcal{N}_j, f^{(i)}(\boldsymbol{x}_1) > f^{(i)}(\boldsymbol{x}_2) \right\} \right|}{n_i n_j}$$

**Optimization methods** As for the optimization methods, we adopt the nestrov-like acceleration method for non-convex optimization algorithm [41], [46], [47] for training linear model composed with softmax and we adopt adam to train deep neural networks. We adopt SGD method to train the models for User-Imb dataset, while we use the adam [37], [80] to train the models for the CIFAR-100-Imb dataset.

*The readers are referred to Sec.H for more details about how the models are trained in this paper.*

TABLE 2
Performance comparison based on MAUC$^{\uparrow}$ with Deep Learning

| type | method | CIFAR-100 | User | iNaturalist |
|---|---|---|---|---|
| Baseline | CE | 59.80 | 55.26 | 66.84 |
| Sampling | BM | 59.07 | 58.95 | 78.18 |
| | MM | 58.52 | 55.35 | 67.07 |
| | IHT | 60.56 | 55.87 | 67.08 |
| | NM | 60.24 | 54.52 | 66.74 |
| | TL | 60.28 | 52.99 | 66.77 |
| Loss | FOCAL | 60.03 | 55.34 | 66.97 |
| | CBCE | 60.04 | 58.91 | 77.71 |
| | CBFOCAL | 60.42 | 58.75 | 71.12 |
| | LDAM | 60.26 | 56.47 | 73.18 |
| AUC | DeepAUC | 60.27 | 59.93 | 62.48 |
| | Ours1 | 61.90 | 60.64 | 79.67 |
| | Ours2 | **62.08** | **60.94** | **84.87** |
| | Ours3 | 59.64 | 59.52 | 79.16 |

## 7.4 Results and Analysis

**Traditional Datasets**. The average performances of 15 repetitions for the nine traditional datasets are shown in Fig.1, where the scatters show the 15 observations over different dataset splits, and the bar plots show the average performance over 15 repetitions. Consequently, we have the following observations: 1) The best performance of our proposed algorithm consistently surpasses all the competitors significantly on all the datasets. More specifically, the best algorithm among Ours1, Our2 and Our3 outperforms the

best competitors by a margin of 5.3, 1.5, 5.6, 1.3, 4.4, 2.9, 4.9, 1.4, 1.8 for *Balance*, *Dermatology*, *Ecoli*, *New Thyroid*, *Page Blocks*, *SegmentImb*, *Shuttle*, *Svmguide2*, and *Yeast*, respectively. It turns out that our improvements are significant in most cases. 2) It could always be observed that some of the sampling-based competitors fail to outperform the original LR algorithm. One reason for this phenomenon might be that the sampling-based methods are not directly designed for optimizing AUC. Another possible reason is that most of the sampling methods adopt an ova strategy to perform multiclass sampling. This has a similar negative effect, as shown in Thm.1 for AUC$^{\mathrm{ova}}$, where the imbalance issue across class pairs is not taken into consideration. 3) Based on the fact that an imbalanced dataset's performance bottleneck comes from its minority classes, we investigate a closer look at the performance comparison on minority class pairs. Specifically, we visualize the finer-grained comparison for 5 class pairs having the smallest frequency for each data, which is shown in Fig.2. Note that we only present the results that are greater than 0.7 to have a clearer view of the differences across top algorithms. The results show that our improvement over the minority class pairs is even more significant, especially for the Balance, Ecoli, Page Blocks, SegmentImb, Shuttle, and Yeast.

**Deep Learning Datasets**. The performance comparisons are shown in the Tab.2. Consequently, when using either the squared loss (Ours1) or the exponential loss (Ours2), our proposed algorithm consistently outperforms all the competitors for all the datasets. Note that DeepAUC has a reasonable performance on the first two datasets. However, it fails in the last dataset. The possible reason here is that DeepAUC requires an extra set of parameters with the size of which depending on $N_C$. Hence, for a large dataset like iNaturalist, training DeepAUC becomes significantly more difficult than other methods. Next, we show the finer-grained comparisons of the minority class pairs in Tab.3 and Tab.7.4. For CIFAR-100-Imb and User-Imb, we adopt a similar to the traditional datasets where the performance comparisons in terms of AUC$_{i|j}$ on the class pairs $(i, j)$ with bottom-5 pairwise frequency $p_i p_j$. For the CIFAR-100-Imb dataset, the best algorithm among Ours1, Our2, and Our3 significantly outperform their best competitor for all the class pairs except the 4-th one. Moreover, even though the hinge loss (Ours3) does not provide an improvement of MAUC$^{\uparrow}$, it at least mitigates the imbalanced issue by providing good performances on the first four minority pairs. For User-Imb, the corresponding improvements produced by our proposed methods are not that sharp as the previous results. The possible reasons are that: (1) the imbalance degree of User-Imb is relatively moderate; (2) even the minority classes in this dataset have a sufficient amount of samples as high as 400. These two traits automatically alleviate the imbalance issue in this dataset, making the difference between different imbalance-aware methods seem less significant. For the iNaturalist 2017 Dataset, since the tail classes are too rate to differentiate the performance from the different algorithms, we turn to report the result for bottom 0-5%, 6-10%, 11-15%, 16-20%, 21-25% class pairs. Again, we can see that our proposed algorithm outperforms the competitors consistently.

TABLE 3
**Fine-grained Comparison Over the Minority Class Pairs (Deep Learning Datasets).** Again, we provide a finer-grained comparison result on the minority class pairs of the two deep learning datasets. Here 1st, 2nd, 3rd, 4th, 5th are the **bottom 5 class pairs** $(i, j)$ with smallest $p_i p_j$.

| type | method | CIFAR-100-Imb | | | | | User-Imb | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | bottom pairs | 1st | 2nd | 3rd | 4th | 5th | 1st | 2nd | 3rd | 4th | 5th |
| Baseline | CE | 48.82 | 53.28 | 50.74 | 58.44 | 56.17 | 49.43 | 49.51 | 53.72 | 58.06 | 49.73 |
| Sampling | BM | 50.46 | 50.65 | 51.97 | 54.44 | 56.79 | 54.64 | 53.04 | 56.36 | 58.44 | 58.50 |
| | MM | 46.45 | 38.70 | 39.07 | 52.61 | 44.71 | 52.77 | 49.68 | 54.82 | 60.42 | 48.39 |
| | IHT | 46.18 | 51.60 | 53.53 | 53.67 | 58.63 | 56.49 | 54.77 | 58.02 | 61.60 | 52.00 |
| | NM | 49.01 | 55.40 | 54.35 | 60.78 | 61.75 | 49.38 | 51.21 | 53.99 | 56.05 | 52.31 |
| | TL | 48.75 | 52.88 | 50.43 | 52.11 | 55.83 | 47.55 | 47.96 | 50.11 | 52.35 | 50.92 |
| Imbalanced | FOCAL | 51.25 | 52.58 | 48.89 | 58.17 | 60.38 | 49.83 | 49.88 | 54.16 | 58.86 | 49.77 |
| | CBCE | 48.42 | 53.33 | 52.53 | 55.06 | 58.67 | **61.48** | **59.05** | 61.60 | 64.42 | **58.82** |
| | CBFOCAL | 49.74 | 51.18 | 49.45 | 58.22 | 58.29 | 61.45 | 58.98 | 61.39 | 64.09 | 58.75 |
| | LDAM | 48.68 | 50.45 | 52.57 | 57.17 | 59.96 | 50.46 | 52.18 | 55.20 | 60.64 | 49.33 |
| AUC | DeepAUC | 48.55 | 53.03 | 54.68 | **62.06** | 61.88 | 51.56 | 55.30 | 55.27 | 55.45 | 57.34 |
| | Ours1 | 51.18 | 57.48 | 59.39 | 60.33 | 62.04 | 58.74 | 58.20 | 62.41 | 64.49 | 57.58 |
| | Ours2 | 52.57 | 54.68 | 57.31 | 59.33 | **65.38** | 59.02 | 58.54 | **62.90** | **65.26** | 57.55 |
| | Ours3 | **54.34** | **60.70** | **61.18** | 58.17 | 65.13 | 59.52 | 56.59 | 60.09 | 63.56 | 56.37 |

TABLE 4
Peformance Comparison on iNaturalist 2017 Dataset.

| type | method | iNaturalist | | | | |
|---|---|---|---|---|---|---|
| | bottom pairs | 0-5% | 5-10% | 10-15% | 15-20% | 20-25% |
| Baseline | CE | 62.62 | 62.45 | 62.33 | 62.55 | 62.78 |
| Sampling | BM | 78.68 | 78.20 | 77.96 | 78.06 | 78.06 |
| | MM | 62.83 | 62.65 | 62.57 | 62.78 | 63.04 |
| | IHT | 63.08 | 62.82 | 62.66 | 62.80 | 63.07 |
| | NM | 62.66 | 62.44 | 62.30 | 62.45 | 62.71 |
| | TL | 62.50 | 62.32 | 62.23 | 62.46 | 62.75 |
| Imbalanced | FOCAL | 62.97 | 62.69 | 62.53 | 62.73 | 63.01 |
| | CBCE | 83.09 | 81.22 | 80.18 | 79.50 | 79.02 |
| | CBFOCAL | 78.65 | 75.90 | 74.49 | 73.55 | 72.91 |
| | LDAM | 70.94 | 70.48 | 70.34 | 70.49 | 70.84 |
| AUC | DeepAUC | 51.31 | 51.34 | 51.31 | 51.21 | 50.92 |
| | Ours1 | 80.29 | 79.37 | 78.98 | 79.02 | 79.00 |
| | Ours2 | **85.14** | **84.33** | **83.93** | **83.94** | **83.92** |
| | Ours3 | 78.65 | 75.90 | 74.49 | 73.55 | 72.91 |

# 8 CONCLUSION

This paper provides an early study of AUC-guided machine learning for multiclass problems. Specifically, we propose a novel framework to learn scoring functions by optimizing the popular M metric where employs surrogate losses for the 0-1 loss to leverage a differentiable objective function. Utilizing consistency analysis, we show that a series of surrogate loss fucntions are fisher consistent with the 0-1 loss based M metric. Moreover, we provide an empirical surrogate risk minimization framework to minimize the MAUC$^\downarrow$ with guaranteed generalization upper bounds. Practically, we propose acceleration methods for three implementations of our proposed framework. Finally, empirical studies on 11 real-world datasets show the efficacy of our proposed algorithms.

# 9 ACKOWNLEDGMENT

# REFERENCES

[1] S. Agarwal. Surrogate regret bounds for bipartite ranking via strongly proper losses. *Journal of Machine Learning Research*, 15(1):1653–1674, 2014.

[2] S. Agarwal, T. Graepel, R. Herbrich, S. Har-Peled, and D. Roth. Generalization bounds for the area under the roc curve. *Journal of Machine Learning Research*, 6(Apr):393–425, 2005.

[3] A. H. Alan and B. Raskutti. Optimising area under the roc curve using gradient descent. *International Conference on Machine Learning*, pages 49–56, 2004.

[4] Z. Bai, X.-L. Zhang, and J. Chen. Partial auc optimization based deep speaker embeddings with class-center learning for text-independent speaker verification. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6819–6823. IEEE, 2020.

[5] P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

[6] S. Barua, M. M. Islam, X. Yao, and K. Murase. Mwmote–majority weighted minority oversampling technique for imbalanced data set learning. *IEEE Transactions on Knowledge and Data Engineering*, 26(2):405–425, 2012.

[7] S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.

[8] T. Calders and S. Jaroszewicz. Efficient auc optimization for classification. *European Conference on Principles of Data Mining and Knowledge Discovery*, pages 42–53, 2007.

[9] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems*, 2019.

[10] Y. Chen, B. Chen, X. He, C. Gao, Y. Li, J.-G. Lou, and Y. Wang. λopt: Learn to regularize recommender models in finer levels. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 978–986, 2019.

[11] S. Clémençon and M. Achab. Ranking data with continuous labels through oriented recursive partitions. In *Advances in Neural Information Processing Systems*, pages 4600–4608, 2017.

[12] S. Clémençon, G. Lugosi, N. Vayatis, et al. Ranking and empirical minimization of u-statistics. *The Annals of Statistics*, 36(2):844–874, 2008.

[13] S. Clémençon, S. Robbiano, and N. Vayatis. Ranking data with ordinal labels: optimality and pairwise aggregation. *Machine Learning*, 91(1):67–104, 2013.

[14] C. Cortes, V. Kuznetsov, M. Mohri, and S. Yang. Structured prediction theory based on factor graph complexity. *Advances in Neural Information Processing Systems*, pages 2514–2522, 2016.

[15] C. Cortes and M. Mohri. Auc optimization vs. error rate minimization. *Advances in Neural Information Processing Systems*, pages 313–320, 2003.

[16] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie. Class-balanced loss based on effective number of samples. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 9268–9277, 2019.

[17] L. Dai, Y. Yin, C. Qin, T. Xu, X. He, E. Chen, and H. Xiong. Enterprise cooperation and competition analysis with a sign-oriented preference network. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 774–782, 2020.

[18] Z. Dang, X. Li, B. Gu, C. Deng, and H. Huang. Large-scale nonlinear auc maximization via triply stochastic gradients. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[19] J. Ding, F. Feng, X. He, G. Yu, Y. Li, and D. Jin. An improved sampler for bayesian personalized ranking by leveraging view data. In *The Web Conference*, 2018.

[20] T. Fawcett. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8):861–874, 2006.

[21] C. Ferri, J. Hernández-Orallo, and M. A. Salido. Volume under the ROC surface for multi-class problems. *European Conference on Machine Learning*, pages 108–120, 2003.

[22] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer. An efficient boosting algorithm for combining preferences. *Journal of machine learning research*, 4(Nov):933–969, 2003.

[23] W. Gao, L. Wang, R. Jin, S. Zhu, and Z. Zhou. One-pass auc optimization. *International Conference on Machine Learning*, pages 906–914, 2013.

[24] W. Gao and W. Wang. Analysis of k-partite ranking algorithm in area under the receiver operating characteristic curve criterion. *International Journal of Computer Mathematics*, 95(8):1527–1547, 2018.

[25] W. Gao and Z. Zhou. On the consistency of AUC pairwise optimization. *International Joint Conference on Artificial Intelligence*, pages 939–945, 2015.

[26] J. L. Gardy, C. Spencer, K. Wang, M. Ester, G. E. Tusnady, I. Simon, S. Hua, K. deFays, C. Lambert, K. Nakai, et al. Psort-b: Improving protein subcellular localization prediction for gram-negative bacteria. *Nucleic Acids Research*, 31(13):3613–3617, 2003.

[27] N. Golowich, A. Rakhlin, and O. Shamir. Size-independent sample complexity of neural networks. pages 297–299, 2018.

[28] B. Gu, Z. Huo, and H. Huang. Scalable and efficient pairwise learning to achieve statistical accuracy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3697–3704, 2019.

[29] H. Han, W.-Y. Wang, and B.-H. Mao. Borderline-smote: a new over-sampling method in imbalanced data sets learning. *International Conference on Intelligent Computing*, pages 878–887, 2005.

[30] D. J. Hand and R. J. Till. A simple generalisation of the area under the roc curve for multiple class classification problems. *Machine Learning*, 45(2):171–186, 2001.

[31] J. A. Hanley and B. J. McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982.

[32] H. Hao, H. Fu, Y. Xu, J. Yang, F. Li, X. Zhang, J. Liu, and Y. Zhao. Open-narrow-synechiae anterior chamber angle classification in AS-OCT sequences. In *International Conference on Medical Image Computing and Computer Assisted Intervention*, 2020.

[33] P. Honzik, P. Kucera, O. Hyncica, and V. Jirsik. Novel method for evaluation of multi-class area under receiver operating characteristic. *International Conference on Soft Computing, Computing with Words and Perceptions in System Analysis, Decision and Control*, pages 1–4, 2009.

[34] C.-W. Hsu, C.-C. Chang, C.-J. Lin, et al. A practical guide to support vector classification, 2003.

[35] T. Joachims. A support vector method for multivariate performance measures. *International Conference on Machine Learning*, pages 377–384, 2005.

[36] T. Joachims. Training linear svms in linear time. *the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 217–226, 2006.

[37] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning Representations*, 2015.

[38] T. Lane. Position paper: Extensions of roc analysis to multi-class domains. *ICML Workshop on Cost-sensitive Learning*, 2000.

[39] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.

[40] G. Lemaître, F. Nogueira, and C. K. Aridas. Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*, 18(1):559–563, 2017.

[41] H. Li and Z. Lin. Accelerated proximal gradient methods for nonconvex programming. *Advances in Neural Information Processing Systems*, pages 379–387, 2015.

[42] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *IEEE international conference on computer vision*, pages 2980–2988, 2017.

[43] C. Liu, Q. Zhong, X. Ao, L. Sun, W. Lin, J. Feng, Q. He, and J. Tang. Fraud transactions detection via behavior tree with local intention calibration. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2020.

[44] M. Liu, Z. Yuan, Y. Ying, and T. Yang. Stochastic auc maximization with deep neural networks. In *International Conference on Learning Representations (ICLR)*, 2020.

[45] M. Liu, Z. Yuan, Y. Ying, and T. Yang. Stochastic AUC maximization with deep neural networks. In *ICLR*, 2020.

[46] R. Liu, S. Cheng, L. Ma, X. Fan, and Z. Luo. Deep proximal unrolling: Algorithmic framework, convergence analysis and applications. *IEEE Transactions on Image Processing*, 28(10):5013–5026, 2019.

[47] R. Liu, Y. Zhang, S. Cheng, X. Fan, and Z. Luo. A theoretically guaranteed deep optimization framework for robust compressive sensing mri. *AAAI Conference on Artificial Intelligence*, pages 4368–4375, 2019.

[48] W. Liu, W. Luo, D. Lian, and S. Gao. Future frame prediction for anomaly detection - A new baseline. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 6536–6545, 2018.

[49] P. M. Long and H. Sedghi. Generalization bounds for deep convolutional neural networks. In *International Conference on Learning Representations*, 2020.

[50] I. Mani and I. Zhang. knn approach to unbalanced data distributions: a case study involving information extraction. *ICML Workshop on Learning from Imbalanced Datasets*, 126, 2003.

[51] A. Maurer. A vector-contraction inequality for rademacher complexities. *International Conference on Algorithmic Learning Theory*, pages 3–17, 2016.

[52] C. McDiarmid. Concentration. pages 195–248. Springer, 1998.

[53] M. Mohri, A. Rostamizadeh, and A. Talwalkar. Foundations of machine learning. 2018.

[54] D. Mossman. Three-way rocs. *Medical Decision Making*, 19(1):78–89, 1999.

[55] H. Narasimhan and S. Agarwal. A structural SVM based approach for optimizing partial AUC. *International Conference on Machine Learning*, pages 516–524, 2013.

[56] H. Narasimhan and S. Agarwal. A structural SVM based approach for optimizing partial AUC. In *International Conference on Machine Learning*, pages 516–524, 2013.

[57] H. Narasimhan and S. Agarwal. Svm$_{pauc}$$^{tight}$: a new support vector method for optimizing partial AUC based on a tight convex upper bound. In *International Conference on Knowledge Discovery and Data Mining*, pages 167–175, 2013.

[58] H. Narasimhan and S. Agarwal. Support vector algorithms for optimizing the partial area under the ROC curve. *Neural Computation*, 29(7):1919–1963, 2017.

[59] M. Natole, Y. Ying, and S. Lyu. Stochastic proximal algorithms for auc maximization. *International Conference on Machine Learning*, pages 3707–3716, 2018.

[60] M. A. Natole, Y. Ying, and S. Lyu. Stochastic auc optimization algorithms with linear convergence. *Frontiers in Applied Mathematics and Statistics*, 5:30, 2019.

[61] F. Pan, S. Li, X. Ao, P. Tang, and Q. He. Warm up cold-start advertisements: Improving ctr predictions via learning to learn id embeddings. In *International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 695–704, 2019.

[62] F. J. Provost and P. M. Domingos. Tree induction for probability-based ranking. *Machine Learning*, 52(3):199–215, 2003.

[63] S. Rajaram and S. Agarwal. Generalization bounds for k-partite ranking. In *NIPS Workshop on Learning to Rank*, pages 18–23, 2005.

[64] L. Ralaivola, M. Szafranski, and G. Stempfel. Chromatic pac-bayes bounds for non-iid data: Applications to ranking and stationary $\beta$-mixing processes. *Journal of Machine Learning Research*, 11(Jul):1927–1956, 2010.

[65] H. W. Reeve and A. Kaban. Optimistic bounds for multi-output prediction. *arXiv preprint arXiv:2002.09769*, 2020.

[66] W. Shi, B. Gu, X. Li, and H. Huang. Quadruply stochastic gradient method for large scale nonlinear semi-supervised ordinal regression

auc optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5734–5741, 2020.

[67] M. R. Smith, T. Martinez, and C. Giraud-Carrier. An instance level analysis of data complexity. *Machine Learning*, 95(2):225–256, 2014.

[68] I. Tomek. Two modifications of cnn. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-6(11):769–772, 1976.

[69] K. Uematsu and Y. Lee. Statistical optimality in multipartite ranking and ordinal regression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(5):1080–1094, 2014.

[70] N. Usunier, M.-R. Amini, and P. Gallinari. A data-dependent generalisation error bound for the auc. *ICML Workshop on ROC Analysis in Machine Learning*, 2005.

[71] N. Usunier, M. R. Amini, and P. Gallinari. Generalization error bounds for classifiers trained with interdependent data. *Advances in Neural Information Processing Systems*, pages 1369–1376, 2006.

[72] S. Wang and L. L. Minku. Auc estimation and concept drift detection for imbalanced data streams with multiple classes. In *2020 International Joint Conference on Neural Networks*, pages 1–8, 2020.

[73] H. Wu, Z. Hu, J. Jia, Y. Bu, X. He, and T.-S. Chua. Mining unfollow behavior in large-scale online social networks via spatial-temporal interaction. In *AAAI Conference on Artificial Intelligence*, volume 34, pages 254–261, 2020.

[74] B. Yang. The extension of the area under the receiver operating characteristic curve to multi-class problems. *Isecs International Colloquium on Computing, Communication, Control, and Management*, pages 463–466, 2009.

[75] Z. Yang, W. Shen, Y. Ying, and X. Yuan. Stochastic auc optimization with general loss. *Communications on Pure & Applied Analysis*, 19(8), 2020.

[76] Y. Ying, L. Wen, and S. Lyu. Stochastic online auc maximization. *Advances in Neural Information Processing Systems*, pages 451–459, 2016.

[77] X. Zhang, A. Saha, and S. Vishwanathan. Smoothing multivariate performance measures. *Journal of Machine Learning Research*, 13(1):3623–3680, 2012.

[78] P. Zhao, S. C. Hoi, R. Jin, and T. Yang. Online auc maximization. In *International Conference on Machine Learning*, pages 233–240, 2011.

[79] K. Zhou, S. Gao, J. Cheng, Z. Gu, H. Fu, Z. Tu, J. Yang, Y. Zhao, and J. Liu. Sparse-gan: Sparsity-constrained generative adversarial network for anomaly detection in retinal oct image. In *IEEE International Symposium on Biomedical Imaging*, pages 1227–1231, 2020.

[80] F. Zou, L. Shen, Z. Jie, W. Zhang, and W. Liu. A sufficient condition for convergences of adam and rmsprop. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 11127–11135, 2019.

**Qianqian Xu** received the B.S. degree in computer science from China University of Mining and Technology in 2007 and the Ph.D. degree in computer science from University of Chinese Academy of Sciences in 2013. She is currently an Associate Professor with the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China. Her research interests include statistical machine learning, with applications in multimedia and computer vision. She has authored or coauthored 40+ academic papers in prestigious international journals and conferences (including T-PAMI, IJCV, T-IP, NeurIPS, ICML, CVPR, AAAI, etc), among which she has published 6 full papers with the first author's identity in ACM Multimedia. She served as member of professional committee of CAAI, and member of online program committee of VALSE, etc.

**Shilong Bao** received the B.S. degree in College of Computer Science and Technology from Qingdao University in 2019. He is currently pursuing the M.S. degree with University of Chinese Academy of Sciences. His research interest is machine learning and data mining.

**Xiaochun Cao**, Professor of the Institute of Information Engineering, Chinese Academy of Sciences. He received the B.E. and M.E. degrees both in computer science from Beihang University (BUAA), China, and the Ph.D. degree in computer science from the University of Central Florida, USA, with his dissertation nominated for the university level Outstanding Dissertation Award. After graduation, he spent about three years at ObjectVideo Inc. as a Research Scientist. From 2008 to 2012, he was a professor at Tianjin University. He has authored and coauthored over 100 journal and conference papers. In 2004 and 2010, he was the recipients of the Piero Zamperoni best student paper award at the International Conference on Pattern Recognition. He is a fellow of IET and a Senior Member of IEEE. He is an associate editor of IEEE Transactions on Image Processing, IEEE Transactions on Circuits and Systems for Video Technology and IEEE Transactions on Multimedia.

**Zhiyong Yang** received the M.Sc. degree in computer science and technology from University of Science and Technology Beijing (USTB) in 2017, and the Ph.D. degree from University of Chinese Academy of Sciences (UCAS) in 2021. He is currently a postdoctoral research fellow with the University of Chinese Academy of Sciences. His research interests lie in machine learning and learning theory, with special focus on AUC optimization, meta-learning/multi-task learning, and learning theory for recommender systems. He has authored or coauthored several academic papers in top-tier international conferences and journals including T-PAMI/ICML/NeurIPS/CVPR. He served as a reviewer for several top-tier journals and conferences such as T-PAMI, ICML, NeurIPS and ICLR.

**Qingming Huang** is a chair professor in the University of Chinese Academy of Sciences and an adjunct research professor in the Institute of Computing Technology, Chinese Academy of Sciences. He graduated with a Bachelor degree in Computer Science in 1988 and Ph.D. degree in Computer Engineering in 1994, both from Harbin Institute of Technology, China. His research areas include multimedia computing, image processing, computer vision and pattern recognition. He has authored or coauthored more than 400 academic papers in prestigious international journals and top-level international conferences. He is the associate editor of IEEE Trans. on CSVT and Acta Automatica Sinica, and the reviewer of various international journals including IEEE Trans. on PAMI, IEEE Trans. on Image Processing, IEEE Trans. on Multimedia, etc. He is a Fellow of IEEE and has served as general chair, program chair, track chair and TPC member for various conferences, including ACM Multimedia, CVPR, ICCV, ICME, ICMR, PCM, BigMM, PSIVT, etc.

Fig. 2. **Fine-grained Comparison Over the Minority Class Pairs (Traditional Datasets).** The x-axis gives the frequency rank of the class pairs $(i, j)$, *i.e.*, $p_i p_j$, where a larger rank represents a larger frequency. The y-axis represents the $\mathrm{AUC}_{i|j}$ of the corresponding class pairs. Each line in a plot then captures the minority class pair performance of a given algorithm. To have a clearer look at tendency, we carry out two filtering processes before we visualize the plot: (a) For those datasets which has more than 5 class pairs, **we only visualize the bottom-5 pairs** in terms of the frequency to turn our focus to the minority pairs in the dataset. (b) To have a clearer look at the difference of the top competitors, we filter out the pairs with a smaller $\mathrm{AUC}_{i|j}$ than $0.7$.

# Appendices

## CONTENTS

# A COMPARISON PROPERTIES BETWEEN $AUC^{ova}$ AND $AUC^{ovo}$

In this section, we will provide a comparison result for $AUC^{ova}$ and $AUC^{ovo}$, which suggests that employing $AUC^{ova}$ does not consider the imbalance issue across class pairs. Then we will provide a practical example to compare these two metrics.

**Restate of Theorem 1.** *Given the label distribution as $\mathbb{P}[y = i] = p_i > 0$ and any multiclass scoring function $f$. The following properties hold:*

*(a) We have that:*

$$\text{AUC}^{\text{ova}}(f) = 1/N_C \sum_{i=1}^{N_C} \sum_{j \neq i} (p_j/1-p_i) \cdot \text{AUC}_{i|j}(f^{(i)}). \tag{9}$$

*(b)*

$$\text{AUC}^{\text{ova}}(f) = \text{AUC}^{\text{ovo}}(f), \ \text{when} \ \ p_i = 1/N_C, \quad i = 1, 2, \cdots, N_C.$$

*(c) We have $\text{AUC}^{\text{ova}}(f) = 1$ if and only if $\text{AUC}^{\text{ovo}}(f) = 1$.*

**Proof.**
*Proof of (a). :For any event $\mathcal{C}$, we have:*

$$\mathbb{P}\left[\mathcal{C}, \mathcal{E}^{(i)}\right] = \sum_{j \neq i} \mathbb{P}\left[\mathcal{C}|\mathcal{E}^{(ij)}\right] \mathbb{P}\left[\mathcal{E}^{(ij)}\right].$$

$$\begin{aligned}
\text{AUC}_{i|\neg i} &= \mathbb{P}\left[(y_1^{(i)} - y_2^{(i)}) \cdot (f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)) > 0|\mathcal{E}^{(i)}\right] + \frac{1}{2}\mathbb{P}\left[f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) = 0|\mathcal{E}^{(i)}\right] \\
&= \sum_{j \neq i} \frac{\mathbb{P}\left[\mathcal{E}^{(ij)}\right]}{\mathbb{P}\left[\mathcal{E}^{(i)}\right]} \cdot \left(\mathbb{P}\left[(y_1^{(i)} - y_2^{(i)}) \cdot (f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)) > 0|\mathcal{E}^{(ij)}\right]\right. \\
&\quad \left. + \frac{1}{2}\mathbb{P}\left[f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) = 0|\mathcal{E}^{(ij)}\right]\right) \\
&= \sum_{j \neq i} \frac{\mathbb{P}\left[\mathcal{E}^{(ij)}\right]}{\mathbb{P}\left[\mathcal{E}^{(i)}\right]} \text{AUC}_{i|j} \\
&= \sum_{j \neq i} \frac{p_j}{1 - p_i} \text{AUC}_{i|j}
\end{aligned}$$

*It then follows that :*

$$\begin{aligned}
\text{AUC}^{\text{ovo}} &= \frac{1}{N_C} \sum_{i=1}^{N_C} \text{AUC}_{i|\neg i} \\
&= \frac{1}{N_C} \sum_{i=1}^{N_C} \sum_{j \neq i} \frac{p_j}{1 - p_i} \text{AUC}_{i|j}.
\end{aligned}$$

*Proof of (b). When $p_i = \frac{1}{N_C}$, $\forall i = 1, 2, \cdots, N_C$, we have $\frac{p_j}{1-p_i} = \frac{p_i}{1-p_j} = \frac{1}{N_C-1}$, which completes the proof with the conclusion of (a).*

*Proof of (c). Since $\text{AUC}^{\text{ovo}}$ and $\text{AUC}^{\text{ova}}$ are convex combinations of $\text{AUC}_{i|j}$, and that $\text{AUC}_{i|j} \in [0, 1]$, we have*

$$\text{AUC}^{\text{ovo}} = 1 \leftrightarrow \text{AUC}_{i|j} = 1, \ \forall i \neq j \leftrightarrow \text{AUC}^{\text{ova}} = 1.$$

$\square$

# B CONSISTENCY ANALYSIS

In this section, we will provide all the proofs for Consistency Analysis.

## B.1 The Bayes Optimal Scoring Function

First, we show the derivation of the Bayes optimal scoring function set under the MAUC$^\downarrow$ criterion.

**Restate of Theorem 2.** *Given $\eta_i(\cdot) = \mathbb{P}[y = i|x]$, $p_i = \mathbb{P}[y = i]$, we have the following consequences:*

*(a)* $f = \{f^{(i)}\}_{i=1,2,\cdots,N_C} \in \mathcal{F}_\sigma^{N_C}$ *is a Bayes optimal scoring function under the MAUC$^\downarrow$ criterion, if:*

$$\left( f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) \right) \cdot \left( \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) - \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1) \right) > 0, \ \forall \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \neq \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1),$$

*where:*

$$\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) = \sum_{j \neq i} \frac{\eta_i(\boldsymbol{x}_1)\eta_j(\boldsymbol{x}_2)}{2p_i p_j}, \quad \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1) = \sum_{j \neq i} \frac{\eta_j(\boldsymbol{x}_1)\eta_i(\boldsymbol{x}_2)}{2p_i p_j}.$$

*(b) Define $\sigma(\cdot)$ as the sigmoid function, $s_i(\boldsymbol{x}) = \eta_i(\boldsymbol{x})/p_i$ and $s_{\setminus i}(\boldsymbol{x}) = \sum_{j \neq i} s_j(\boldsymbol{x})$, then a Bayesian scoring function could be given by:*

$$f^{\star(i)}(\boldsymbol{x}) = \begin{cases} \sigma\left( \dfrac{s_i(\boldsymbol{x})}{s_{\setminus i}(\boldsymbol{x})} \right), & 0 \leq \eta_i(\boldsymbol{x}) < 1 \\ 1, & \eta_i(\boldsymbol{x}) = 1. \end{cases} \tag{10}$$

**Proof.**

*Proof of (a)*

*First let us summarize the notations we will use in the forthcoming derivation in Tab.5.*

### TABLE 5
### Notations and descriptions.

| Notation | Description |
|---|---|
| $\mathcal{Y}_{ij}$ | is the event $[y_1 = i, y_2 = j]$ |
| $\mathcal{E}_{i,j}$ | is the event $\mathcal{Y}_{ij}$ or $\mathcal{Y}_{ji}$ |
| $y_1^{(i)}$ | $= 1$ if $y_1 = i$, otherwise $y_1^{(i)} = 0$ |
| $y_2^{(i)}$ | $= 1$ if $y_2 = i$, otherwise $y_2^{(i)} = 0$ |
| $\boldsymbol{x}_{1,2}$ | refers to the couple $(\boldsymbol{x}_1, \boldsymbol{x}_2)$ |
| $y_{1,2}$ | refers to the couple $(y_1, y_2)$ |
| $\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2})$ | is event that $(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)) \cdot (y_1^{(i)} - y_2^{(i)}) < 0$ |
| $\mathcal{G}_i((1,0), \boldsymbol{x}_{1,2})$ | is event that $(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)) \cdot (1 - 0) < 0$ |
| $\mathcal{G}_i((0,1), \boldsymbol{x}_{1,2})$ | is event that $(f^{(i)}\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)) \cdot (0 - 1) < 0$ |
| $\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2})$ | is the event that $f^{(i)}(\boldsymbol{x}_1) = f^{(i)}(\boldsymbol{x}_2)$ |
| $\eta_i(\boldsymbol{x})$ | $\mathbb{P}(y = i|\boldsymbol{x})$ |
| $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2)$ | $\sum_{j \neq i}(\frac{\eta_i(\boldsymbol{x}_1)\eta_j(\boldsymbol{x}_2)}{2p_i p_j})$ |
| $\pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$ | $\sum_{j \neq i}(\frac{\eta_j(\boldsymbol{x}_1)\eta_i(\boldsymbol{x}_2)}{2p_i p_j})$ |

*First of all, we prove that the condition holds for:*

$$f \in \operatorname*{arginf}_{f \in \mathcal{M}} R(f)$$

*where*

$$\mathcal{M} = \{f = (f^{(1)}, \cdots, f^{(N_C)}), f^{(i)} \text{ is a measurable function.}\}$$

*First of all, we can remove the dependence on the conditioning of $\mathcal{E}_{i,j}$. We have*

$$
\begin{aligned}
N_C \cdot (N_C - 1) \cdot \text{MAUC}^\downarrow &= \sum_{i=1}^{N_C} \sum_{j \neq i} \left( \mathbb{P}\big[\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2})|\mathcal{E}_{i,j}\big] + \frac{1}{2} \cdot \mathbb{P}\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2})|\mathcal{E}_{i,j}\big] \right) \\
&= \sum_{i=1}^{N_C} \sum_{j \neq i} \frac{\mathbb{P}\big[\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2}), \mathcal{E}_{i,j}\big]}{\mathbb{P}\big[\mathcal{E}_{i,j}\big]} + \frac{1}{2} \frac{\mathbb{P}\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2}), \mathcal{E}_{i,j}\big]}{\mathbb{P}\big[\mathcal{E}_{i,j}\big]} \\
&= \sum_{i=1}^{N_C} \sum_{j \neq i} \frac{\mathbb{P}\big[\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2}), \mathcal{E}_{i,j}\big]}{2p_i p_j} + \frac{1}{2} \frac{\mathbb{P}\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2}), \mathcal{E}_{i,j}\big]}{2p_i p_j}
\end{aligned}
\tag{11}
$$

*Now we only need to expand the joint possibility by expectation, which leads to:*

$$\frac{\mathbb{P}\big[\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2}), \mathcal{E}_{i,j}\big]}{2p_ip_j} = \frac{\mathbb{E}_{\boldsymbol{x}_{1,2}, y_{1,2}}\big[\boldsymbol{I}[\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2})] \cdot \boldsymbol{I}[\mathcal{E}_{i,j}]\big]}{2p_ip_j}$$

$$= \frac{\mathbb{E}_{\boldsymbol{x}_{1,2}}\Big[\mathbb{E}_{y_{1,2}|\boldsymbol{x}_{1,2}}\big[\boldsymbol{I}[\mathcal{G}_i(y_{1,2}, \boldsymbol{x}_{1,2})] \cdot \boldsymbol{I}[\mathcal{E}_{i,j}]\big]\Big]}{2p_ip_j} \tag{12}$$

$$= \frac{\mathbb{E}_{\boldsymbol{x}_{1,2}}\Big[\eta_i(\boldsymbol{x}_1) \cdot \eta_j(\boldsymbol{x}_2) \cdot \boldsymbol{I}[\mathcal{G}_i\,((1,0), \boldsymbol{x}_{1,2})] + \eta_j(\boldsymbol{x}_1) \cdot \eta_i(\boldsymbol{x}_2) \cdot \boldsymbol{I}[\mathcal{G}_i\,((0,1), \boldsymbol{x}_{1,2})]\Big]}{2p_ip_j}.$$

*Similarly, we have:*

$$\frac{\mathbb{P}\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2}), \mathcal{E}_{i,j}\big]}{2p_ip_j} = \frac{\mathbb{E}_{\boldsymbol{x}_{1,2}}\Big[\eta_i(\boldsymbol{x}_1) \cdot \eta_j(\boldsymbol{x}_2) \cdot \boldsymbol{I}\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2})\big] + \eta_j(\boldsymbol{x}_1) \cdot \eta_i(\boldsymbol{x}_2) \cdot \boldsymbol{I}\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2})\big]\Big]}{2p_ip_j} \tag{13}$$

*Above all, by combining Eq.(11)-Eq.(13), we reach that :*

$$N_C \cdot (N_C - 1) \cdot \mathsf{MAUC}^{\downarrow} = \sum_{i=1}^{N_C} \mathbb{E}_{x_{1,2}} \Big[ \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \cdot \boldsymbol{I}\big[\mathcal{G}_i\,((1,0), \boldsymbol{x}_{1,2})\big] + \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)) \cdot \boldsymbol{I}\big[\mathcal{G}_i\,((0,1), \boldsymbol{x}_{1,2})\big]$$

$$+ \frac{1}{2} \cdot \big(\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) + \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)\big) \cdot I\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2})\big]\Big]. \tag{14}$$

*Fix a sub-scoring function $f^{(i)}$ and $\boldsymbol{x}_{1,2}$, we only need to minimize the following quantity $L^{(i)}$ to reach the Bayes optimal solution:*

$$L^{(i)} = \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \cdot \boldsymbol{I}\big[\mathcal{G}_i\,((1,0), \boldsymbol{x}_{1,2})\big] + \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)) \cdot \boldsymbol{I}\big[\mathcal{G}_i\,((0,1), \boldsymbol{x}_{1,2})\big] + \frac{1}{2} \cdot \big(\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) + \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)\big) \cdot I\big[\mathcal{G}_{0,i}(\boldsymbol{x}_{1,2})\big] \tag{15}$$

*Obviously ,*

$$L^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) = \begin{cases} \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1), & f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) > 0, \\[2mm] \dfrac{\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) + \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)}{2}, & f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) = 0, \\[2mm] \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2), & f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) < 0. \end{cases} \tag{16}$$

*To obtain the Bayes optimal scoring function, we must minimize $L(\boldsymbol{x}_1, \boldsymbol{x}_2)$ for any $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{X}$.*
*When $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) = \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$, $L(\boldsymbol{x}_1, \boldsymbol{x}_2)$ stays as a constant, which is irrelevant to the choice of $f^{(i)}$.*
*When $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) < \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$, we must have $f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) < 0$ to minimize $L(\boldsymbol{x}_1, \boldsymbol{x}_2)$.*
*When $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) > \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$, we must have $f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) > 0$ to minimize $L(\boldsymbol{x}_1, \boldsymbol{x}_2)$.*
*Above all, by choosing $f^{(i)}$ such that $\forall \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{X}, \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \neq \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$*

$$\Big(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)\Big) \cdot \Big(\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) - \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)\Big) > 0,$$

*one can reach the minimum of $R^{(i)}(f^{(i)})$ as $\mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2} \min\{\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2), \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)\}$.*
*This ends the proof of the condition for:*

$$f \in \underset{f \in \mathcal{M}}{\arg\inf}\, R(f)$$

*To extend the conclusion to:*

$$f \in \underset{f \in \mathcal{F}_\sigma^{N_C}}{\arg\inf}\, R(f)$$

*One only need to notice that if $f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) > 0$, then for any strictly increasing measurable function $\varphi$, we have: $\varphi(f^{(i)}(\boldsymbol{x}_1)) - \varphi(f^{(i)}(\boldsymbol{x}_2)) > 0$. Obviousely, a similar result holds for $f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) < 0$. This implies if $f \in \mathcal{M}$ is Bayes optimal, then*

$$\varphi \circ f = (\varphi \circ f^{(1)}, \cdots, \varphi \circ f^{(N_C)})$$

*is also Bayes optimal with $\varphi \circ f \in \mathcal{F}_\sigma^{N_C}$. The proof is then completed.*

*Proof of (b)*
*Fixing class $i$, if $\eta_i(\boldsymbol{x}_k) \neq 1$, $k = 1, 2$, then by dividing $s_{\setminus i}(\boldsymbol{x}_1) \cdot s_{\setminus i}(\boldsymbol{x}_2)$ in Eq.(2) in Thm.2, we find that $f^{\star(i)}(\boldsymbol{x})$ could be chosen as*

$\sigma(s_i(\boldsymbol{x})/s_{\setminus i}(\boldsymbol{x}))$.

*Otherwise, if $\eta_i(\boldsymbol{x}_1) = 1, \eta_i(\boldsymbol{x}_2) < 1$ (note that $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \neq \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$) then $s_{\setminus i}(\boldsymbol{x}_1) = 0$. To reach $f^{\star(i)}(\boldsymbol{x}_1) > f^{\star(i)}(\boldsymbol{x}_2)$, we can set $f^{\star(i)}(\boldsymbol{x}_1) = 1, f^{\star(i)}(\boldsymbol{x}_2) = \sigma(s_i(\boldsymbol{x}_2)/s_{\setminus i}(\boldsymbol{x}_2))$.*

*In the last case, we have $\eta_i(\boldsymbol{x}_1) < 1, \eta_i(\boldsymbol{x}_2) = 1$, we can reach that an optimal solution could be $f(\boldsymbol{x}_1) = \sigma(s_i(\boldsymbol{x}_1)/s_{\setminus i}(\boldsymbol{x}_1)), f(\boldsymbol{x}_2) = 1$. Since $\boldsymbol{x}_1, \boldsymbol{x}_2$ is arbitrarily chosen, a Bayesian scoring function could be then set as: $f^{\star(i)}(\boldsymbol{x}) = \sigma(s_i(\boldsymbol{x})/s_{\setminus i}(\boldsymbol{x}))(< 1)$ if $\eta_i(\boldsymbol{x}) < 1$; $f^{\star(i)}(\boldsymbol{x}) = 1$ otherwise.* □

## B.2 Consistency Analysis of Surrogate Losses

Based on Thm.2, we provide the following result for some popular surrogate losses.

**Restate of Theorem 3.** *The surrogate loss $\ell$ is consistent with $\mathsf{MAUC}^{\downarrow}$ if it is differentiable, convex, nonincreasing within $[-1, 1]$ and $\ell'(0) < 0$.*

**Proof.** *First, recall that we have the following definition for the surrogate risk:*

$$R_\ell(f) = \sum_i \frac{R_\ell^{(i)}(f^{(i)})}{N_C(N_C - 1)}$$

$$R_\ell^{(i)}(f^{(i)}) = \sum_{j \neq i} \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2} \left[ \ell(\Delta(y_{1,2}^{(i)})\Delta f^{(i)}) | \mathcal{E}^{(ij)} \right],$$

*where $R_\ell^{(i)}(f^{(i)})$ is the corresponding risk for $f^{(i)}$.*

*Denote:*

$$f^\star = \underset{f \in \mathcal{F}_\sigma}{\mathrm{arginf}}\, R_\ell(f)$$

*with $f^\star = (f^{\star(1)}, f^{\star(2)}, \cdots, f^{\star(N_C)})$.*

*Moreover, it is obvious that:*

$$\inf_{f \in \mathcal{F}_\sigma} R_\ell(f) = \frac{1}{N_C \cdot (N_C - 1)} \cdot \sum_i \inf_{f^{(i)} \in \mathcal{F}_\sigma} R_\ell^{(i)}(f^{(i)}). \tag{17}$$

*In this sense, each binary scoring function $f^{\star(i)}$ could be solved independently for different subproblems, i.e.*

$$f^{\star(i)} = \underset{f^{(i)} \in \mathcal{F}_\sigma}{\mathrm{arginf}}\, R_\ell^{(i)}(f^{(i)})$$

*Similar to the derivation in Thm.2, we could rewrite $R_\ell^{(i)}$ as:*

$$R_\ell^{(i)}(f^{(i)}) = \frac{1}{2} \int_{\mathcal{X}} \int_{\mathcal{X}} \left[ \eta_+^{(i)}(\boldsymbol{x}) \cdot \eta_-^{(i)}(\boldsymbol{x}') \cdot \ell(f^{(i)}(\boldsymbol{x}) - f^{(i)}(x')) + \eta_+^{(i)}(\boldsymbol{x}') \cdot \eta_-^{(i)}(\boldsymbol{x}) \cdot \ell(f^{(i)}(\boldsymbol{x}') - f^{(i)}(\boldsymbol{x})) \right] d\mathbb{P}(\boldsymbol{x}) d\mathbb{P}(\boldsymbol{x}')$$

*where*

$$\eta_+^{(i)}(\boldsymbol{x}) = \frac{\eta_i(\boldsymbol{x})}{p_i}, \ \eta_-^{(i)}(\boldsymbol{x}) = \sum_{j \neq i} \frac{\eta_j(\boldsymbol{x})}{p_j}.$$

*Denote $\mathsf{Bayes}_\sigma^{(i)}$ as the set of $f^{(i)}$, where $f$ is a Bayes optimal socring functions, i.e.*

$$\mathsf{Bayes}_\sigma^{(i)} = \{f^{(i)} : \left( f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2) \right) \cdot (\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) - \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)) > 0, \ \forall(\boldsymbol{x}_1, \boldsymbol{x}_2) \ s.t. \ \pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) \neq \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)\}$$

*Now proof the theorem with the following steps.*

---

**Claim 1.** $f^{\star(i)} \in \mathsf{Bayes}_\sigma^{(i)}$

---

*We prove the claim by contradiction. To do this, we assume that claim does not hold in the sense that $\exists \boldsymbol{x}_1, \boldsymbol{x}_2$ s.t. $f^{\star(i)}(\boldsymbol{x}_1) \leq f^{\star(i)}(\boldsymbol{x}_2)$ but $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) > \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$.*

**Case 1:** $\eta_-^{(i)}(\boldsymbol{x}_1) > 0, \eta_-^{(i)}(\boldsymbol{x}_2) > 0$

*Define an intermediate function $\delta_h(\gamma) = R_\ell^{(i)}(f^{\star(i)} + \gamma h)$, we must have $\delta_h'(0) = 0, \forall h$.*

*Given any pair of instance $(\boldsymbol{x}_1, \boldsymbol{x}_2) \in \mathcal{X}$, with $\boldsymbol{x}_1 \neq \boldsymbol{x}_2$, and for all $\ell$ satisfying the sufficient condition, we prove that $f^{\star(i)}(\boldsymbol{x}_1) > f^{\star(i)}(\boldsymbol{x}_2)$ must imply that $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) > \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$ by contradiction.*

*Picking*

$$h_1(x) = \frac{\boldsymbol{I}\,[x = x_1]}{\eta_-(\boldsymbol{x})},$$

*using $\delta'_{h_1}(0) = 0$, we have:*

$$\int_{X\setminus \boldsymbol{x}_1} \frac{\eta_+^{(i)}(\boldsymbol{x}_1)}{\eta_-^{(i)}(\boldsymbol{x}_1)} \cdot \eta_-^{(i)}(\boldsymbol{x}) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_1) - f^{\star(i)}(\boldsymbol{x})\right) - \eta_+^{(i)}(\boldsymbol{x}) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_1)\right) d\mathbb{P}(\boldsymbol{x}) = 0 \tag{18}$$

*Picking*

$$h_2(x) = \frac{\boldsymbol{I}\,[x = x_2]}{\eta_-(\boldsymbol{x})},$$

*using $\delta'_{h_2}(0) = 0$, we have:*

$$\int_{X\setminus \boldsymbol{x}_2} \frac{\eta_+^{(i)}(\boldsymbol{x}_2)}{\eta_-^{(i)}(\boldsymbol{x}_2)} \cdot \eta_-^{(i)}(\boldsymbol{x}) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_2) - f^{\star(i)}(\boldsymbol{x})\right) - \eta_+^{(i)}(\boldsymbol{x}) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_2)\right) d\mathbb{P}(\boldsymbol{x}) = 0 \tag{19}$$

*Taking (18) - (19), we have:*

$$\underbrace{\int_{X\setminus\{\boldsymbol{x}_1,\boldsymbol{x}_2\}} \eta_+^{(i)}(\boldsymbol{x}) \cdot \left(\ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_2)\right) - \ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_1)\right)\right) d\mathbb{P}(\boldsymbol{x})}_{(a)}$$

$$+ \underbrace{\int_{X\setminus\{\boldsymbol{x}_1,\boldsymbol{x}_2\}} \frac{\eta_-^{(i)}(\boldsymbol{x})}{\eta_-^{(i)}(\boldsymbol{x}_2)\eta_-^{(i)}(\boldsymbol{x}_1)} \left(\pi^{(i)}(\boldsymbol{x}_1,\boldsymbol{x}_2) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_1) - f^{\star(i)}(\boldsymbol{x})\right) - \pi^{(i)}(\boldsymbol{x}_2,\boldsymbol{x}_1) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_2) - f^{\star(i)}(\boldsymbol{x})\right)\right) d\mathbb{P}(\boldsymbol{x})}_{(b)} \tag{20}$$

$$+ \underbrace{\left(\frac{\mathbb{P}(\boldsymbol{x}_1)}{\eta_-^{(i)}(\boldsymbol{x}_1)} + \frac{\mathbb{P}(\boldsymbol{x}_2)}{\eta_-^{(i)}(\boldsymbol{x}_2)}\right) \cdot \left(\pi^{(i)}(\boldsymbol{x}_1,\boldsymbol{x}_2) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_1) - f^{\star(i)}(\boldsymbol{x})\right) - \pi^{(i)}(\boldsymbol{x}_2,\boldsymbol{x}_1) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_2) - f^{\star(i)}(\boldsymbol{x})\right)\right)}_{(c)}$$

$$= 0$$

*We now construct a contradiction with the equation above.*

*First notice that*

$$\ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_2)\right) \le \ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_1)\right)$$

*since by assumption $\ell$ is nonconvex, i.e., $\ell'$ is nondecreasing. This implies $(a) \le 0$.*

*With a similar spirit, we have:*

$$\ell'\left(f^{\star(i)}(\boldsymbol{x}_1) - f^{\star(i)}(\boldsymbol{x})\right) \le \ell'\left(f^{\star(i)}(\boldsymbol{x}_2) - f^{\star(i)}(\boldsymbol{x})\right) \le 0.$$

*Together with $\pi^{(i)}(\boldsymbol{x}_1,\boldsymbol{x}_2) > \pi^{(i)}(\boldsymbol{x}_2,\boldsymbol{x}_1)$, we have $(b) \le 0$.*

*Now, we prove that $(c) < 0$:*

**Case a:** *If $f^{\star(i)}(\boldsymbol{x}_1) = f^{\star(i)}(\boldsymbol{x}_2)$, we have :*

$$(c) = (\pi^{(i)}(\boldsymbol{x}_1,\boldsymbol{x}_2) - \pi^{(i)}(\boldsymbol{x}_2,\boldsymbol{x}_1)) \cdot \ell'(0) < 0$$

*since $\ell(0) < 0$.*

**Case b** *If $f^{\star(i)}(\boldsymbol{x}_1) < f^{\star(i)}(\boldsymbol{x}_2)$, we have:*

$$\ell'\left(f^{\star(i)}(\boldsymbol{x}_1) - f^{\star(i)}(\boldsymbol{x}_2)\right) \le \ell'(0) < 0,$$

$$\ell'\left(f^{\star(i)}(\boldsymbol{x}_2) - f^{\star(i)}(\boldsymbol{x}_1)\right) \le 0,$$

*and $\pi^{(i)}(\boldsymbol{x}_2,\boldsymbol{x}_1) > \pi^{(i)}(\boldsymbol{x}_1,\boldsymbol{x}_2)$, we have:$(c) < 0$*

*Putting all togther, we have that $(a) + (b) + (c) < 0$, which contradicts with eq.(21). Since the choice of $(\boldsymbol{x}_1,\boldsymbol{x}_2)$ is arbitrary, this shows that*

$$f^{\star(i)} \in \mathsf{Bayes}_\sigma^{(i)}, \ if \ \eta_-^{(i)}(\boldsymbol{x}_1) > 0, \eta_-^{(i)}(\boldsymbol{x}_2) > 0$$

**Case 2:** $\eta_-^{(i)}(\boldsymbol{x}_1) = 0, \eta_-^{(i)}(\boldsymbol{x}_2) > 0$

*Since :*

$$\sum_i \mathbb{P}\,[\boldsymbol{x}_1, y = i] = \sum_i \eta_i(\boldsymbol{x}_1)\mathbb{P}(\boldsymbol{x}_1) = \mathbb{P}(\boldsymbol{x}_1)$$

*We have :*

$$\sum_i \eta_i(\boldsymbol{x}_1) = 1.$$

*This further implies that $\eta_i(\boldsymbol{x}_1)$ reaches its maximum 1, which means $\eta_+^{(i)}(\boldsymbol{x}_1) > \eta_+^{(i)}(\boldsymbol{x}_2)$.*
*Similar to the proof of Case 1, we can set $h_1(\boldsymbol{x}) = \boldsymbol{I}\left[\boldsymbol{x} = \boldsymbol{x}_1\right]$, $h_2(\boldsymbol{x}) = \boldsymbol{I}\left[\boldsymbol{x} = \boldsymbol{x}_2\right]$ and obtain the following equation.*

$$\underbrace{\int_{X\setminus\{\boldsymbol{x}_1,\boldsymbol{x}_2\}} \eta_-^{(i)}(\boldsymbol{x}) \cdot \left(\eta_+^{(i)}(\boldsymbol{x}_1) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_1)\right) - \eta_+^{(i)}(\boldsymbol{x}_2) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}) - f^{\star(i)}(\boldsymbol{x}_2)\right)\right) d\mathbb{P}(\boldsymbol{x})}_{(a)}$$

$$+ \underbrace{\int_{X\setminus\{\boldsymbol{x}_1,\boldsymbol{x}_2\}} \eta_+^{(i)}(\boldsymbol{x}) \cdot \eta_-^{(i)}(\boldsymbol{x}_2) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_2) - f^{\star(i)}(\boldsymbol{x})\right) d\mathbb{P}(\boldsymbol{x})}_{(b)} \tag{21}$$

$$+ \underbrace{(\mathbb{P}(\boldsymbol{x}_1) + \mathbb{P}(\boldsymbol{x}_2)) \cdot \eta_+^{(i)}(\boldsymbol{x}_1) \cdot \eta_-^{(i)}(\boldsymbol{x}_2) \cdot \ell'\left(f^{\star(i)}(\boldsymbol{x}_1) - f^{\star(i)}(\boldsymbol{x}_2)\right)}_{(c)}$$

$$= 0$$

*Similarly, one can show that $(a) + (b) + (c) < 0$, which contradicts with the optimal condition.*

$$f^{\star(i)} \in \mathsf{Bayes}_\sigma^{(i)}, \ if \ \eta_-^{(i)}(\boldsymbol{x}_1) = 0, \eta_-^{(i)}(\boldsymbol{x}_2) > 0$$

**Case 3:** $\eta_-^{(i)}(\boldsymbol{x}_1) > 0, \eta_-^{(i)}(\boldsymbol{x}_2) = 0$. *In this sense, we have $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) < \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$, which contradicts with the assumption that $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) > \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$.* **Thus this is impossible to observe.**
**Case 4:** $\eta_-^{(i)}(\boldsymbol{x}_1) = 0, \eta_-^{(i)}(\boldsymbol{x}_2) = 0$. *We have $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) = \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$, which contradicts with the assumption that $\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2) > \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)$.* **Thus this is impossible to observe.**
*Then the proof of Claim 1 is completed.*

---

**Claim 2.**

$$\inf_{f^{(i)}\notin\mathsf{Bayes}_\sigma} R_\ell^{(i)}(f^{(i)}) > \inf_{f^{(i)}\in\mathcal{F}_\sigma} R_\ell^{(i)}(f^{(i)}) \tag{22}$$

---

*Claim 2 follows that $\operatorname{arginf}_{f^{(i)}\in\mathcal{F}_\sigma} R_\ell^{(i)}(f^{(i)}) = \mathsf{Bayes}_\sigma^{(i)}$.*

---

**Claim 3.** *For any sequence $\{f_t\}_{t\in\mathbb{N}_+}$, such that $f_i \in \mathcal{F}_\sigma$,*

$$R_\ell^{(i)}(f_t) \to R_\ell^{(i)}(f^{\star(i)}) \ implies \ R^{(i)}(f_t) \to \inf_{f\in\mathcal{F}_\sigma} R^{(i)}(f)$$

---

*Since $\operatorname{arginf}_{f\in\mathcal{F}_\sigma} R^{(i)}(f) = \mathsf{Bayes}_\sigma^{(i)}$, we only need to prove that $\lim_{t\to\infty} f_t \in \mathsf{Bayes}_\sigma^{(i)}$. Define $\delta^{(i)}$ as:*

$$\delta^{(i)} = \inf_{f^{(i)}\notin\mathsf{Bayes}_\sigma^{(i)}} R_\ell^{(i)}(f^{(i)}) - \inf_{f^{(i)}\in\mathcal{F}_\sigma} R_\ell^{(i)}(f^{(i)}) > 0$$

*Suppose that $\lim f_t \notin \mathcal{F}_\sigma$, then for large enough $T$, we have:*

$$R_\ell^{(i)}(f_T) - R_\ell^{(i)}(f^{\star(i)}) > \delta^{(i)},$$

*which contradicts with the fact that $R_\ell^{(i)}(f_t) \to R_\ell^{(i)}(f^{\star(i)})$. This shows that Eq.(22) holds.*
*Since the argument above is based on an arbitrarily chosen class $i$, Eq.(22) holds for all $i$.*

---

**Claim 4.** *Given a sequence $\{f\}_t$ with $f_t = (f_t^{(1)}, \cdots, f_t^{(N_C)})$, we have:*

$$R_\ell(f_t) \to R_\ell(f) \ implies \ R(f_t) \to \inf_{f\in\mathcal{F}_\sigma} R(f)$$

---

*This directly follows that:*

$$R(f_t) = \frac{\sum_{i=1}^{N_C} R^{(i)}(f_t^{(i)})}{N_C \cdot (N_C - 1)}, \ R_\ell(f_t) = \frac{\sum_{i=1}^{N_C} R_\ell^{(i)}(f_t^{(i)})}{N_C \cdot (N_C - 1)}.$$

*and Claim 4.*

*The proof is then completed according to Def.1.* □

## C  UNBIASED ESTIMATION OF $R_\ell$

In this section, we provide a derivation of an unbiased estimation of $R_\ell$.

**Restate of Proposition 1.** $\hat{R}_{\ell,\mathcal{S}}(f)$ as:

$$\hat{R}_{\ell,\mathcal{S}}(f) = \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)).$$

$\hat{R}_{\ell,\mathcal{S}}(f)$ is an unbiased estimation of $R_\ell(f)$, in the sense that: $R_\ell(f) = \underset{\mathcal{S}}{\mathbb{E}}(\hat{R}_{\ell,\mathcal{S}}(f))$, where $\mathcal{N}_i$ denotes the set of all samples in $\mathcal{S}$ with $y = i$.

**Proof.**

$$\underset{\mathcal{S}}{\mathbb{E}}(\hat{R}_{\ell,\mathcal{S}}(f)) = \underset{\boldsymbol{Y}}{\mathbb{E}}\left[\underset{\boldsymbol{X}|\boldsymbol{Y}}{\mathbb{E}}(\hat{R}_{\ell,\mathcal{S}}(f))\right]$$

*Moreover, we have*

$$\underset{\boldsymbol{X}|\boldsymbol{Y}}{\mathbb{E}}(\hat{R}_{\ell,\mathcal{S}}(f))$$

$$= \frac{1}{N_C(N_C - 1)} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_1 \in \mathcal{N}_i} \sum_{\boldsymbol{x}_2 \in \mathcal{N}_j} \frac{1}{n_i n_j} \underset{\boldsymbol{x}_1, \boldsymbol{x}_2}{\mathbb{E}}\left[\ell(f^{(i)}, \boldsymbol{x}_1, \boldsymbol{x}_2)|y_1 = i, y_2 = j\right]$$

$$= \frac{1}{N_C(N_C - 1)} \sum_{i=1}^{N_C} \sum_{j \neq i} \underset{\boldsymbol{x}_1, \boldsymbol{x}_2}{\mathbb{E}}\left[\ell(f^{(i)}, \boldsymbol{x}_1, \boldsymbol{x}_2)|y_1 = i, y_2 = j\right].$$

*Since $\mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2}\left[\ell(f^{(i)}, \boldsymbol{x}_1, \boldsymbol{x}_2)|y_1 = i, y_2 = j\right]$ does not depend on the distribution of $\boldsymbol{Y}$, we have:*

$$\underset{\boldsymbol{Y}}{\mathbb{E}}\left[\underset{\boldsymbol{X}|\boldsymbol{Y}}{\mathbb{E}}(\hat{R}_{\ell,\mathcal{S}}(f))\right] = \underset{\boldsymbol{X}|\boldsymbol{Y}}{\mathbb{E}}(\hat{R}_{\ell,\mathcal{S}}(f)).$$

*Now it only remains to prove that*

$$\underset{\boldsymbol{x}_1, \boldsymbol{x}_2}{\mathbb{E}}\left[\ell\left(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)\right)|y_1 = i, y_2 = j\right] = \underset{\boldsymbol{z}_1, \boldsymbol{z}_2}{\mathbb{E}}\left[\ell\left(\Delta(y^{(i)})\Delta(f^{(i)})\right)|\mathcal{E}^{(ij)}\right].$$

*To see this, we have:*

$$\underset{\boldsymbol{z}_1, \boldsymbol{z}_2}{\mathbb{E}}\left[\ell\left(\Delta(y^{(i)})\Delta(f^{(i)})\right)|\mathcal{E}^{(ij)}\right]$$

$$= \underset{\boldsymbol{z}_1, \boldsymbol{z}_2}{\mathbb{E}}\left[\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2)\ell\left(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)\right) + \pi^{(i)}(\boldsymbol{x}_2, \boldsymbol{x}_1)\ell\left(f^{(i)}(\boldsymbol{x}_2) - f^{(i)}(\boldsymbol{x}_1)\right)\right]$$

$$= 2\underset{\boldsymbol{x}_1, \boldsymbol{x}_2}{\mathbb{E}}\left[\pi^{(i)}(\boldsymbol{x}_1, \boldsymbol{x}_2)\ell\left(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)\right)\right]$$

$$= \underset{\boldsymbol{x}_1, \boldsymbol{x}_2}{\mathbb{E}}\left[\frac{\eta^{(i)}(\boldsymbol{x}_1)\eta^{(j)}(\boldsymbol{x}_2)}{p_i p_j}\ell\left(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)\right)\right]$$

$$= \underset{\boldsymbol{x}_1, \boldsymbol{x}_2}{\mathbb{E}}\left[\ell\left(f^{(i)}(\boldsymbol{x}_1) - f^{(i)}(\boldsymbol{x}_2)\right)|y_1 = i, y_2 = j\right].$$

*The proof is thus completed.* □

## D  PRELIMINARY FOR GENERALIZATION ANALYSES

### D.1  Concentration Inequalities

In this subsection, we provide a brief introduction to the concentration bounds that are employed throughout the next two sections.

### D.1.1 Bounded Difference Property

**Definition 5** (Bounded Difference Property). *Given independent random variables $X_1, \cdots, X_n$, with $X_i \in \mathbb{X}$, A function $f(X_1, X_2, \cdots, X_n)$ is said to has the bounded difference property if there exist non-negative constants $c_1, c_2, \cdots, c_n$, such that:*

$$\sup_{x_1, x_2, \cdots, x_n, x_i'} |f(x_1, \cdots, x_n) - f(x_1, \cdots, x_{i-1}, x_i', \cdots, x_n)| \leq c_i, \ \forall 1 \leq i \leq n. \tag{23}$$

For all functions satisfying the Bounded Difference Property, we have the following Bounded Difference Inequality over the generated momentum functions.

**Lemma 1** (Bounded Difference Inequality). *[7, Prop.6.1, Thm.6.2] Assume that $Z = f(X_1, \cdots, X_n)$, with $X_i$s being independent, satisfies the bounded difference property with constants $c_1, c_2, \cdots, c_n$. Denote*

$$v = \frac{1}{4} \sum_{i=1}^{n} c_i^2 \tag{24}$$

*then we have:*

$$\log \mathbb{E}\left[\exp\left(\lambda(Z - \mathbb{E}[Z])\right)\right] \leq \frac{\lambda^2 v^2}{2}, \tag{25}$$

*holds for all $\lambda > 0$.*

### D.1.2 Maximal Inequality

**Lemma 2** (Maximal Inequality). *[7, Sec.2.5] Let $Z_1, \cdots, Z_n$ be real-valued random variables where a $v > 0$ exists, such that for every $i = 1, 2, \cdots, n$, we have $\log\left(\mathbb{E}\left[\exp\left(\lambda Z_i\right)\right]\right) \leq \frac{\lambda v^2}{2}$, then we have:*

$$\mathbb{E}\left[\max_{i=1,2,\cdots,n} Z_i\right] \leq \sqrt{2v \log n}.$$

### D.1.3 Mcdiarmid Inequality

First let us review the fundamental concentration inequality that is adopted in our proof.

**Lemma 3** (Mcdiarmid inequality). *[52] Let $X_1, \cdots, X_m$ be independent random variables all taking values in the set $\mathcal{X}$. Let $f : \mathcal{X} \to \mathbb{R}$ be a function of $X_1, \cdots, X_m$ that satisfies:*

$$\sup_{\boldsymbol{x}, \boldsymbol{x}'} |f(x_1, \cdots, x_i, \cdots, x_m) - f(x_1, \cdots, x_i' \cdots, x_,)| \leq c_i,$$

*with $\boldsymbol{x} \neq \boldsymbol{x}'$. Then for all $\epsilon > 0$,*

$$\mathbb{P}[\mathbb{E}(f) - f \geq \epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^{m} c_i^2}\right).$$

## D.2 Properties of Rademacher Averages

The following lemmas are instrumental for our derivation of MAUC$^\downarrow$ Rademacher Complexity upper bounds in Lem.10 and Lem.11.

First the Khinchin-Kahane inequality bridges the $\ell_p$ norm average to the $\ell_q$ norm average, given that $1 < p < q < \infty$.

**Lemma 4** (Khinchin-Kahane inequality). *[39] $\forall n \in \mathbb{N}$, let $\sigma_1, \cdots, \sigma_n$ be a family of i.i.d. Rademacher variables. Then for any $1 < p < q < \infty$, we have:*

$$\left(\mathbb{E}_{\sigma}\left[|\sum_{i=1}^{n} \sigma_i f_i|^q\right]\right)^{1/q} \leq \left(\frac{q-1}{p-1}\right)^{1/2} \cdot \left(\mathbb{E}_{\sigma}\left[|\sum_{i=1}^{n} \sigma_i f_i|^p\right]\right)^{1/p}.$$

Dealing with the Rademacher complexities for vector-valued function, we need the vector version of the well-known Talagrand contraction lemma, see [14], [51] for a proof.

**Lemma 5** (vector version Talagrand contraction lemma). *Let $\mathcal{X}$ be any set, $n \in \mathcal{N}$, $(\boldsymbol{x}_1; \cdots; \boldsymbol{x}_N) \in X_n$, let $\mathcal{F}$ be a class of $k$-component-functions $f = (f^{(1)}, \cdots, f^{(k)}) : X \to \ell_2^k$ and let $h : \ell_2^k \to \mathbb{R}$ have Lipschitz norm $\phi$. Then,*

$$\mathbb{E}_{\sigma}\left[\sup_{f \in \mathcal{F}} \sum_{i=1}^{N} \sigma_i h(f(\boldsymbol{x}_i))\right] \leq \sqrt{2}\phi \, \mathbb{E}_{\sigma}\left[\sup_{f \in \mathcal{F}} \sum_{i=1}^{N} \sum_{k=1}^{K} \sigma_{i,k}(f^{(k)}(\boldsymbol{x}_i))\right],$$

*where $\sigma_1 \cdots, \sigma_N$, and $\sigma_{1,1}, \cdots, \sigma_{1,K}, \cdots, \sigma_{N,K}$ are two sequences of independent Rademacher random variables.*

The following Lemma is a simple extension of Lem.1 in [27] to double indexed Rademacher random sequence, which is the key to the derivation of Lem.11.

**Lemma 6.** *Let $s$ be a 1-Lipschitz, positive-homogeneous activation function which is applied element-wise. Then for any class of vector-valued functions $\mathcal{F}$, and any convex and monotonically increasing function $g : \mathbb{R} \to [0, \infty)$, we have:*

$$\mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}, \|W\|_F \leq R_W} g \left( \left\| \sum_{i=1}^{N} \sum_{c=1}^{N_C} \sigma_{i,c} \cdot s(\boldsymbol{W} f(\boldsymbol{x}_i)) \right\| \right) \right] \leq 2 \, \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} g \left( R_W \cdot \left\| \sum_{i=1}^{N} \sum_{c=1}^{N_C} \sigma_{i,c} \cdot f(\boldsymbol{x}_i) \right\| \right) \right]$$

*where $\{\sigma_{i,c}\}$ is a sequence of double indexed Rademacher random variables.*

### D.3 Lipschitz Properties of Softmax

In this section, we provide an $\ell_2$ Lipschitz Constant for softmax component function.

**Lemma 7** (Lipschitz Constant for softmax components). *Given $\mathcal{X} \in \mathbb{R}^K$, the function $\mathsf{soft}_i$, $i = 1, 2, \cdots, K$ is a mapping $\boldsymbol{x} = (x_1, \cdots, x_k) \in \mathcal{X} \to [0, 1]$, defined as:*

$$\mathsf{soft}_i(\boldsymbol{x}) = \frac{\exp(x_i)}{\sum_{j=1}^{K} \exp(x_j)},$$

*then $\mathsf{soft}_i(\cdot)$ is $\frac{\sqrt{2}}{2}$-Lipschitz continuous with respect to vector $\ell_2$ norm.*

**Proof.** *We only need to prove that*

$$\sup_{\boldsymbol{x} \in \mathcal{X}} \left\| \nabla_{\boldsymbol{x}} \mathsf{soft}_i(\boldsymbol{x}) \right\|_2 \leq \frac{\sqrt{2}}{2}.$$

*For any $\boldsymbol{x} \in \mathcal{X}$, we have:*

$$\frac{\partial \mathsf{soft}_i(\boldsymbol{x})}{\partial x_j} = \mathsf{soft}_i(\boldsymbol{x}) \cdot (I[i = j] - \mathsf{soft}_j(\boldsymbol{x})), \; i, j = 1, \cdots, K.$$

*Then we have:*

$$\left\| \nabla_{\boldsymbol{x}} \mathsf{soft}_i(\boldsymbol{x}) \right\|_2 = \left( \mathsf{soft}_i^2(\boldsymbol{x}) \cdot \sum_{j \neq i} \mathsf{soft}_j^2(\boldsymbol{x}) + \mathsf{soft}_i(\boldsymbol{x})^2 \cdot (1 - \mathsf{soft}_j(\boldsymbol{x}))^2 \right)^{1/2}.$$

*Since $\mathsf{soft}_i^2(\boldsymbol{x}) \leq \mathsf{soft}_i(\boldsymbol{x})$, $(1 - \mathsf{soft}_j(\boldsymbol{x}))^2 \leq (1 - \mathsf{soft}_j(\boldsymbol{x}))$, and $\sum_{j \neq i} \mathsf{soft}_j^2(\boldsymbol{x}) \leq (1 - \mathsf{soft}_i(\boldsymbol{x}))$, we have:*

$$\left\| \nabla_{\boldsymbol{x}} \mathsf{soft}_i(\boldsymbol{x}) \right\|_2 \leq \left( 2 \cdot \mathsf{soft}_i(\boldsymbol{x}) \cdot (1 - \mathsf{soft}_i(\boldsymbol{x})) \right)^{1/2} \leq \frac{\sqrt{2}}{2},$$

*which ends the proof.*

## E MAUC$^{\downarrow}$ RADEMACHER COMPLEXITY AND ITS PROPERTIES

### E.1 MAUC$^{\downarrow}$ Symmetrization

In this section, we provide the derivation for MAUC$^{\downarrow}$ Symmetrization, which is the key technology for Thm.4. The main idea is that exchanging instances, but not the terms, does not change the value of

$$\mathbb{E}_{\mathcal{S}, \mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}}(f) - \hat{R}_{\mathcal{S}'}(f) \right) \right].$$

**Lemma 8** (MAUC$^{\downarrow}$ symmetrization). *The following inequality holds:*

$$\mathbb{E}_{\mathcal{S}, \mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}}(f) - \hat{R}_{\mathcal{S}'}(f) \right) \right] \leq \frac{4 \mathfrak{R}_{\mathsf{MAUC}^{\downarrow}}(\ell \circ \mathcal{H})}{N_C(N_C - 1)},$$

*with the labels $\boldsymbol{Y}$ fixed for $\mathcal{S}$ and $\mathcal{S}'$.*

**Proof.** *Denote $\ell(f^{(i)}, \boldsymbol{x}_m, \boldsymbol{x}'_n) = \ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}'_n))$ and define*

$$\begin{aligned} T^{i,j,m,n} = &\frac{\sigma_m^{(i)} + \sigma_n^{(j)}}{2} \ell(f^{(i)}, \boldsymbol{x}'_m, \boldsymbol{x}'_n) + \frac{\sigma_m^{(i)} - \sigma_n^{(j)}}{2} \ell(f^{(i)}, \boldsymbol{x}'_m, \boldsymbol{x}_n) \\ &- \frac{\sigma_m^{(i)} - \sigma_n^{(j)}}{2} \ell(f^{(i)}, \boldsymbol{x}_m, \boldsymbol{x}'_n) - \frac{\sigma_m^{(i)} + \sigma_n^{(j)}}{2} \ell(f^{(i)}, \boldsymbol{x}_m, \boldsymbol{x}_n), \end{aligned} \qquad (26)$$

we first show that

$$\mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'}\left[\sup_{f\in\mathcal{H}}\left(\hat{R}_{\mathcal{S}'}(f)-\hat{R}_{\mathcal{S}}(f)\right)\right]$$

$$=\frac{1}{N_C\cdot(N_C-1)}\mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'}\mathop{\mathbb{E}}_{\sigma}\left[\sup_{f\in\mathcal{H}}\sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j}\frac{1}{n_i n_j}\cdot T^{i,j,m,n}\right].$$ (27)

Given $\mathcal{S}=\{(\boldsymbol{x}_i,y_i)\}_{i=1}^m$, $\mathcal{S}'=\{(\boldsymbol{x}_i',y_i)\}_{i=1}^m$, where we sort the instances such that $y_i\leq y_j$ for $i\leq j$ for the sake of convenience, since the instances are drawn independently, we employ the fact that:

$$\mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'}\left[\sup_{f\in\mathcal{H}}\left(\hat{R}_{\mathcal{S}'}(f)-\hat{R}_{\mathcal{S}}(f)\right)\right]=\mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'}\left[\sup_{f\in\mathcal{H}}\left(\hat{R}_{\tilde{\mathcal{S}}'}(f)-\hat{R}_{\tilde{\mathcal{S}}}(f)\right)\right]$$ (28)

where $\tilde{\mathcal{S}}$ and $\tilde{\mathcal{S}}'$ are permuted datasets obtained from $\mathcal{S}$ and $\mathcal{S}'$ respectively after exchanging $0\leq x\leq N$ samples of $(\boldsymbol{x}_i,y_i)$ and $(\boldsymbol{x}_i',y_i')$ sharing the same index.

Next, we show that for any sequence of i.i.d Rademacher random variables, $\sigma=(\sigma_1^{(1)},\cdots,\sigma_{n_1}^{(1)},\cdots,\sigma_1^{(2)},\cdots,\sigma_{n_C}^{(N_C)})$ and any $\mathcal{S}$, and $\mathcal{S}'$, there exists a pair of permuted datasets $\tilde{\mathcal{S}}^\sigma$ and $\tilde{\mathcal{S}}'^\sigma$ such that

$$\sup_{f\in\mathcal{H}}\left[\sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j}\frac{1}{n_i n_j}\cdot T^{i,j,m,n}\right]=\sup_{f\in\mathcal{H}}\left[R_{\ell,\tilde{\mathcal{S}}^{\sigma'}}(f)-R_{\tilde{\mathcal{S}}^\sigma}(f).\right]$$ (29)

Specifically, we show this by induction.

**Base Case.** We assume that $\mathcal{S}=\{(\boldsymbol{x}_1,1),(\boldsymbol{x}_2,2)\}$ $\mathcal{S}'=\{(\boldsymbol{x}_1',1),(\boldsymbol{x}_2',2)\}$, $\sigma=(\sigma_1^{(1)},\sigma_1^{(2)})$. Here we define

$$T^1=T^{1,2,1,1},T^2=T^{2,1,1,1}.$$

We have the following cases:

(a) $\sigma_1^{(1)}=1,\sigma_1^{(2)}=1$. We have:

$$T^1=\ell(f^{(1)},\boldsymbol{x}_1',\boldsymbol{x}_2')-\ell(f^{(1)},\boldsymbol{x}_1,\boldsymbol{x}_2),$$
$$T^2=\ell(f^{(2)},\boldsymbol{x}_2',\boldsymbol{x}_1')-\ell(f^{(2)},\boldsymbol{x}_2,\boldsymbol{x}_1),$$ (30)

This shows that

$$\sup_{f\in\mathcal{H}}\left[T^1+T^2\right]=\sup_{f\in\mathcal{H}}\left[\hat{R}_{\mathcal{S}'}(f)-\hat{R}_{\mathcal{S}}(f)\right].$$ (31)

This suggests that $\tilde{\mathcal{S}}_\sigma=\mathcal{S}$, $\tilde{\mathcal{S}}'_\sigma=\mathcal{S}'$.

(b) $\sigma_1^{(1)}=1,\sigma_1^{(2)}=-1$.

$$T^1=\ell(f^{(1)},\boldsymbol{x}_1',\boldsymbol{x}_2)-\ell(f^{(1)},\boldsymbol{x}_1,\boldsymbol{x}_2'),$$
$$T^2=\ell(f^{(2)},\boldsymbol{x}_2,\boldsymbol{x}_1')-\ell(f^{(2)},\boldsymbol{x}_2',\boldsymbol{x}_1),$$ (32)

This shows that

$$\sup_{f\in\mathcal{H}}\left[T^1+T^2\right]=\sup_{f\in\mathcal{H}}\left[R_{\ell,\tilde{\mathcal{S}}^{\sigma'}}(f)-R_{\tilde{\mathcal{S}}^\sigma}(f)\right].$$ (33)

where $\tilde{\mathcal{S}}_\sigma$ and $\tilde{\mathcal{S}}'_\sigma$ are obtained by exchanging $(\boldsymbol{x}_2,2)$ and $(\boldsymbol{x}_2',2)$

(c) $\sigma_1^{(1)}=-1,\sigma_1^{(2)}=1$. One can show that the corresponding $\tilde{\mathcal{S}}_\sigma,\tilde{\mathcal{S}}'_\sigma$ are obtained by exchanging $(\boldsymbol{x}_1,1)$ and $(\boldsymbol{x}_1',1)$.

(d) $\sigma_1^{(1)}=-1,\sigma_1^{(2)}=-1$. One can show that the corresponding $\tilde{\mathcal{S}}_\sigma,\tilde{\mathcal{S}}'_\sigma$ are obtained by completely exchanging $\mathcal{S}$ and $\mathcal{S}'$.

The arguments above complete the proof for base case.

**Recursion.** Given

$$\mathcal{S}^-=\{(\boldsymbol{x}_i,y_i)\}_{i=1}^k,\mathcal{S}^{-'}=\{(\boldsymbol{x}_i',y_i)\}_{i=1}^k,\forall\sigma^-\in\{-1,1\}^k,$$

we assume that $\mathcal{S}^-,\mathcal{S}^{-'}$ and $\sigma^-$ satisfy our conclusion, which is realized by $\tilde{\mathcal{S}}_{\sigma^-}^-$ and $\tilde{\mathcal{S}}_{\sigma^-}^{-'}$. We now show that given a new pair $(\boldsymbol{x}_n,i)$ , $(\boldsymbol{x}_n',i)$ and $\sigma_n^{(i)}\in\{-1,1\}$ such that

$$\mathcal{S}=\mathcal{S}^-\cup\{(\boldsymbol{x}_n,i)\},\mathcal{S}'=\mathcal{S}^{-'}\cup\{(\boldsymbol{x}_n',i)\},\sigma=(\sigma^-,\sigma_{new}^{(i)})$$

still satisfies our conclusion. Obviously only $T^{i,*,n,*}$ and $T^{*,i,*,n}$ are involved with the new sample, where $*$ simply refers to all the possible choices.

For the case that $\sigma_n^{(i)}=-1$, we have the following cases: First, we examine the terms $T^{i,*,n,*}$ by taking any specific $T^{i,j,n,m}$.

(a) $\sigma_m^{(j)}=1$, we have:

$$T^{i,j,n,m}=\ell(f^{(i)},\boldsymbol{x}_n,\boldsymbol{x}_m')-\ell(f^{(i)},\boldsymbol{x}_n',\boldsymbol{x}_m).$$

(b) $\sigma_m^{(j)} = -1$, we have:

$$T^{i,j,n,m} = \ell(f^{(i)}, \boldsymbol{x}_n, \boldsymbol{x}_m) - \ell(f^{(i)}, \boldsymbol{x}'_m, \boldsymbol{x}'_n).$$

Moreover, for any $T^{j,i,n,m}$, we have:

(a) $\sigma_m^{(j)} = 1$, we have:

$$T^{j,i,m,n} = \ell(f^{(j)}, \boldsymbol{x}'_m, \boldsymbol{x}_n) - \ell(f^{(j)}, \boldsymbol{x}_m, \boldsymbol{x}'_n).$$

(b) $\sigma_m^{(j)} = -1$, we have:

$$T^{i,j,n,m} = \ell(f^{(j)}, \boldsymbol{x}_m, \boldsymbol{x}_n) - \ell(f^{(j)}, \boldsymbol{x}'_m, \boldsymbol{x}'_n).$$

This shows that

$$\tilde{\mathcal{S}}_\sigma = \tilde{\mathcal{S}}_{\sigma^-}^- \cup \{(\boldsymbol{x}'_n, i)\}, \ \tilde{\mathcal{S}}'_\sigma = \tilde{\mathcal{S}}_{\sigma^-}^{-\prime} \cup \{(\boldsymbol{x}_n, i)\}.$$

Similarly if $\sigma_n^{(i)} = 1$, we have

$$\tilde{\mathcal{S}}_\sigma = \tilde{\mathcal{S}}_{\sigma^-}^- \cup \{(\boldsymbol{x}_n, i)\}, \ \tilde{\mathcal{S}}'_\sigma = \tilde{\mathcal{S}}_{\sigma^-}^{-\prime} \cup \{(\boldsymbol{x}'_n, i)\}.$$

In either case, we know that $\tilde{\mathcal{S}}_\sigma$ and $\tilde{\mathcal{S}}'_\sigma$ are obtained by exchanging the corresponding instances of $\mathcal{S}$ and $\mathcal{S}'$. By the arguments in the base case and the recursion, we complete the proof with:

$$\sup_{f \in \mathcal{H}} \left[ \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot T^{i,j,m,n} \right] = \sup_{f \in \mathcal{H}} \left[ R_{\ell, \tilde{\mathcal{S}}_{\sigma'}}(f) - R_{\tilde{\mathcal{S}}_\sigma}(f) \right]. \tag{34}$$

It immediately suggests that

$$
\begin{aligned}
& \mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \mathop{\mathbb{E}}_{\sigma} \left[ \sup_{f \in \mathcal{H}} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot T^{i,j,m,n} \right] \\
& = \mathop{\mathbb{E}}_{\sigma} \mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( R_{\ell, \tilde{\mathcal{S}}_{\sigma'}}(f) - R_{\tilde{\mathcal{S}}_\sigma}(f) \right) \right] \\
& = \frac{1}{2^N} \cdot \sum_{\sigma} \mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( R_{\ell, \tilde{\mathcal{S}}_{\sigma'}}(f) - R_{\tilde{\mathcal{S}}_\sigma}(f) \right) \right] \\
& = \frac{1}{2^N} \cdot \sum_{\sigma} \mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}'}(f) - \hat{R}_{\mathcal{S}}(f) \right) \right] \\
& = \frac{2^N}{2^N} \cdot \mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}'}(f) - \hat{R}_{\mathcal{S}}(f) \right) \right] \\
& = \mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \left[ \sup_{f \in \mathcal{H}} \left( \hat{R}_{\mathcal{S}'}(f) - \hat{R}_{\mathcal{S}}(f) \right) \right].
\end{aligned}
\tag{35}
$$

The proof then follows that

$$\mathop{\mathbb{E}}_{\mathcal{S},\mathcal{S}'} \mathop{\mathbb{E}}_{\sigma} \left[ \sup_{f \in \mathcal{H}} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot T^{i,j,m,n} \right] \leq 4\mathfrak{R}_{\mathsf{MAUC}^\downarrow}(\ell \circ \mathcal{H}). \tag{36}$$

## E.2 A General Result

In this section, we will provide a proof for the general result for $\mathsf{MAUC}^\downarrow$ generalization bound, which is based on Lem.3, and Lem.8.

**Restate of Theorem 4.** *Given dataset $\mathcal{S} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^m$, where the instances are sampled independently, for all multiclass scoring functions $f \in \mathcal{H}$, if $\mathrm{dom}\,\ell = [0, B]$, $\forall \delta \in (0, 1)$, the following inequalities hold with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\mathcal{S}}(f) + \frac{4\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H})}{N_C(N_C - 1)} + \frac{10B}{N_C} \cdot \xi(\boldsymbol{Y}) \cdot \sqrt{\frac{\log(2/\delta)}{2N}}$$

*where $\xi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{n_C} 1/\rho_i}$.*

**Proof.** *Given $\mathcal{S}'$ as another dataset being independent with $\mathcal{S}$, since $\hat{R}_\mathcal{S}(f) = \mathbb{E}_{\mathcal{S}'}\hat{R}_{\mathcal{S}'}(f)$, using the Jensen's inequality, we have:*

$$\mathbb{E}_\mathcal{S}\left[\sup_{f\in\mathcal{H}}\left(\mathbb{E}_\mathcal{S}\hat{R}_\mathcal{S}(f) - \hat{R}_\mathcal{S}(f)\right)\right] = \mathbb{E}_\mathcal{S}\sup_{f\in\mathcal{H}}\mathbb{E}_{\mathcal{S}'}\left[\hat{R}_\mathcal{S}(f) - \hat{R}_{\mathcal{S}'}(f)\right]$$

$$\leq \mathbb{E}_{\mathcal{S},\mathcal{S}'}\left[\sup_{f\in\mathcal{H}}\left(\hat{R}_\mathcal{S}(f) - \hat{R}_{\mathcal{S}'}(f)\right)\right]$$

*From Lem.8, we have:*

$$\mathbb{E}_{\mathcal{S},\mathcal{S}'}\left[\sup_{f\in\mathcal{H}}\left(\hat{R}_\mathcal{S}(f) - \hat{R}_{\mathcal{S}'}(f)\right)\right] \leq 4\frac{1}{N_C\cdot(N_C-1)}\mathfrak{R}_{\mathsf{MAUC}^\downarrow}(\ell\circ\mathcal{H}). \tag{37}$$

*Given $\mathcal{S}$, define $\mathcal{S}_m = (\mathcal{S}\setminus\{(\boldsymbol{x}_m, y_m)\})\cup\{(\boldsymbol{x}'_m, y_m)\}$. We now proceed to derive the bounded difference property for the loss function. Suppose that $y_m = i$, we have:*

$$\begin{aligned}
d_i &= |\sup_{f\in\mathcal{H}}(\mathbb{E}_\mathcal{S}\hat{R}_\mathcal{S}(f) - \hat{R}_\mathcal{S}(f)) - \sup_{f\in\mathcal{H}}(\mathbb{E}_{\mathcal{S}_i}\hat{R}_{\mathcal{S}_i}(f) - \hat{R}_{\mathcal{S}_i}(f))|\\
&\leq \sup_{f\in\mathcal{H}}|\hat{R}_\mathcal{S}(f) - \hat{R}_{\mathcal{S}_i}(f)|\\
&= 1/N_C(N_C-1)\sup_{f\in\mathcal{H}}\left[\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j}\frac{1}{n_in_j}|\ell(f^{(i)}, \boldsymbol{x}_m, \boldsymbol{x}_n) - \ell(f^{(i)}, \boldsymbol{x}'_m, \boldsymbol{x}_n)|\right.\\
&\quad\left. + \sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j}\frac{1}{n_in_j}|\ell(f^{(j)}, \boldsymbol{x}_n, \boldsymbol{x}_m) - \ell(f^{(j)}, \boldsymbol{x}_n, \boldsymbol{x}'_m)|\right]\\
&\leq \frac{2B}{N_Cn_i}
\end{aligned}$$

*Thus we have $\sum_{i=1}^N d_i^2 = 4B^2/N_C^2\cdot\sum_{i=1}^{n_C} 1/n_i$. According to Eq.(37) and Lem.3, we know that fixing the labels $\boldsymbol{Y}$, with probability at least $1-\delta/2$, the following inequality holds:*

$$\mathbb{E}_\mathcal{S}\hat{R}_\mathcal{S}(f) \leq \hat{R}_\mathcal{S}(f) + \frac{4}{N_C(N_C-1)}\mathfrak{R}_{\mathsf{MAUC}^\downarrow}(\ell\circ\mathcal{H}) + \frac{2B}{N_C}\xi(\boldsymbol{Y})\cdot\sqrt{\frac{\log(2/\delta)}{2N}} \tag{38}$$

*Using an analogous argument, we can prove that, fixing the labels $\boldsymbol{Y}$, with probability at least $1-\delta/2$, the following inequality holds:*

$$\frac{\mathfrak{R}_{\mathsf{MAUC}^\downarrow}(\ell\circ\mathcal{H})}{N_C(N_C-1)} \leq \frac{\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell\circ\mathcal{H})}{N_C(N_C-1)} + \frac{2B}{N_C}\xi(\boldsymbol{Y})\cdot\sqrt{\frac{\log(2/\delta)}{2N}}. \tag{39}$$

*Combining Eq.(38) and Eq.(39), fixing $\boldsymbol{Y}$, we have with probability at least $1-\delta$ that:*

$$\mathbb{E}_\mathcal{S}\hat{R}_\mathcal{S}(f) \leq \hat{R}_\mathcal{S}(f) + \frac{4}{N_C(N_C-1)}\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell\circ\mathcal{H}) + \frac{10B}{N_C}\xi(\boldsymbol{Y})\cdot\sqrt{\frac{\log(2/\delta)}{2N}} \tag{40}$$

*Following the arguments in Thm.8 in [2], we have Eq.(40) holds with probability (over $\boldsymbol{X}$ and $\boldsymbol{Y}$) at least $1-\delta$.*

### E.3 Sub-Gaussian Property of the MAUC$^\downarrow$ Rademacher Complexity

**Definition 6** (Sub-Gaussian Stochastic Process). *A stochastic process $\theta\mapsto X_\theta$ with indexing set $\mathcal{T}$ is sub-Gaussian to a pseudo-metric $d$ on $\mathcal{T}$, if for all $\theta, \theta'\in\mathcal{T}$ and all $\lambda\in\mathbb{R}$,*

$$\mathbb{E}\left[\exp\left(\lambda\cdot(X_\theta - X_{\theta'})\right)\right] \leq \exp\left(\frac{\lambda^2 d(\theta, \theta')^2}{2}\right). \tag{41}$$

Define $T_f(\sigma) = \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} T^{i,j,m,n}$, and recall that

$$T^{i,j,m,n} = \frac{\sigma_m^{(i)} + \sigma_n^{(j)}}{2}\cdot\frac{\ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n))}{n_in_j},$$

with the function $f$ chosen from an index set (in this case a function class) $\mathcal{F}$. We see that $T_f(\sigma)_{f\in\mathcal{F}}$ is essentially a stochastic process of the Rademacher random variables $\sigma$, which is revealed by the following lemma:

**Lemma 9.** *Given an input feature set $\mathcal{D}_\mathcal{X} = \{\boldsymbol{x}_1, \cdots, \boldsymbol{x}_N\}$, $\mathcal{Z}_\mathcal{D}$ is defined as*

$$\{(\boldsymbol{x}, y) : \boldsymbol{x}\in\mathcal{D}_\mathcal{X}, y\in\{1, 2, \cdots, N_C\}\},$$

which is the cartesian product between $\mathcal{D}_{\mathcal{X}}$ and the label space. Moreover, $\forall z = (x, i) \in \mathcal{Z}_{\mathcal{D}}$, $s(z) = s^{(i)}(x)$. With the notations above, if $\ell$ is $\phi_\ell$-Lipschitz continuous, we have that, $\{C_G \cdot T_f(\sigma)\}_{f \in \mathcal{H}}$ with $f = (f^{(1)}, \cdots, f^{N_C})$ and $f^{(i)} = \mathsf{soft}^{(i)} \circ s$ is a sub-Gaussian stochastic process in the sense that :

$$\mathbb{E}_\sigma \left[ \exp \left( \lambda \cdot C_G \cdot \left( T_f(\sigma) - T_{\tilde{f}}(\sigma) \right) \right) \right] \leq \exp \left( \frac{\lambda^2 d_\infty(s, \tilde{s})^2}{2} \right). \tag{42}$$

where

$$C_G = \frac{1}{\phi_\ell \cdot (N_C - 1) \cdot \xi(Y) \cdot \sqrt{\frac{1}{N}}}, \quad d_{\infty, \mathcal{S}}(s, \tilde{s}) = \max_{z \in \mathcal{Z}_{\mathcal{D}}} |s(z) - \tilde{s}(z)|. \tag{43}$$

**Proof.** Denote $\sigma$ as a set of Rademacher random variables in the following form:

$$\sigma = (\sigma_1^{(1)}, \cdots, \sigma_{n_1}^{(1)}, \cdots, \sigma_1^{(i)}, \cdots, \sigma_k^{(i)}, \cdots, \sigma_{n_1}^{(i)}, \cdots, \sigma_{n_{N_C}}^{(N_C)})$$

and denote $\sigma_{\backslash(i,k)}$ as another set of Rademacher random variables with all entries of which the same as $\sigma$ except that $\sigma_k^{(i)}$ is replaced with $\tilde{\sigma}_k^{(i)}$, i.e,

$$\sigma_{\backslash(i,k)} = (\sigma_1^{(1)}, \cdots, \sigma_{n_1}^{(1)}, \cdots, \sigma_1^{(i)}, \cdots, \tilde{\sigma}_k^{(i)}, \cdots, \sigma_{n_1}^{(i)}, \cdots, \sigma_{n_{N_C}}^{(N_C)}).$$

We have the following bounded difference property for every possible choice of $(i, k)$:

$$
\begin{aligned}
diff_i &= \left| \left( T_f(\sigma) - T_{\tilde{f}}(\sigma) \right) - \left( T_f(\sigma_{\backslash(i,k)}) - T_{\tilde{f}}(\sigma_{\backslash(i,k)}) \right) \right| \\
&= \left| \left( \frac{\sigma_k^{(i)} - \tilde{\sigma}_k^{(i)}}{2} \right) \cdot \left[ \sum_{j \neq i} \sum_{x_n \in \mathcal{N}_j} \cdot \frac{1}{n_i n_j} \left( \ell(f^{(i)}, x_m, x_n) - \ell(\tilde{f}^{(i)}, x_m, x_n) \right) \right. \right. \\
&\quad \left. \left. + \sum_{j \neq i} \sum_{x_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \left( \ell(f^{(j)}, x_n, x_m) - \ell(\tilde{f}^{(j)}, x_n, x_m) \right) \right] \right| \\
&\leq \frac{2 \cdot (N_C - 1)}{n_i} \max_{i, \, x_m \in \mathcal{N}_i, \, x_n \notin \mathcal{N}_i} \left| \ell(f^{(i)}, x_m, x_n) - \ell(\tilde{f}^{(i)}, x_m, x_n) \right| \\
&\leq \frac{4 \cdot (N_C - 1)}{n_i} \cdot \phi_\ell \cdot \max_{i, \, x \in \mathcal{D}_{\mathcal{X}}} \left| f^{(i)}(x) - \tilde{f}^{(i)}(x) \right| \\
&= \frac{4 \cdot (N_C - 1)}{n_i} \cdot \phi_\ell \cdot \max_{i, \, x \in \mathcal{D}_{\mathcal{X}}} \left| \int_0^1 \left\langle \nabla \mathsf{soft}^{(i)} \left( \tau \cdot \tilde{s}(x) + (1 - \tau) \cdot s(x) \right), \, (s(x) - \tilde{s}(x)) \right\rangle d\tau \right| \\
&\leq \frac{4 \cdot (N_C - 1)}{n_i} \cdot \phi_\ell \cdot \left( \sup_x ||\nabla \mathsf{soft}^{(i)}(x)||_1 \right) \cdot \max_{z \in \mathcal{Z}_{\mathcal{D}}} |s(z) - \tilde{s}(z)| \\
&\leq \frac{2 \cdot (N_C - 1)}{n_i} \cdot \phi_\ell \cdot \max_{z \in \mathcal{Z}_{\mathcal{D}}} |s(z) - \tilde{s}(z)|.
\end{aligned}
$$

It is now easy to see that $C_G \cdot \left( T_f(\sigma) - T_{\tilde{f}}(\sigma) \right)$ satisfies the bounded difference property (Def.5). According to Lem.1, we could choose $v$ therein as:

$$v = \frac{1}{4} \sum_{i=1}^N C_G^2 \cdot diff_i^2 = \max_{z \in \mathcal{Z}_{\mathcal{D}}} |s(z) - \tilde{s}(z)|^2$$

This suffices to complete the proof based on Lem.1. $\qquad \square$

## E.4 Chaining Bounds

Now with the sub-Gaussian property proved in Lem.9, we can derive chaining upper bounds for the MAUC$^\downarrow$ Rademacher complexity based on the notion of covering number. Before presenting the results, let us start with the definition of $\epsilon$-covering and the covering number.

**Restate of Definition 3** ($\epsilon$-covering). *[39] Let $(\mathcal{H}, d)$ be a (pseudo)metric space, and $\Theta \in \mathcal{H}$. $\{h_1, \cdots, h_K\}$ is said to be an $\epsilon$-covering of $\Theta$ if $\Theta \in \bigcup_{i=1}^K \mathcal{B}(h_i, \epsilon)$, i.e., $\forall \theta \in \Theta, \exists i$ s.t. $d(\theta, h_i) \leq \epsilon$.*

**Restate of Definition 4** (Covering Number). *[39] Based on the notations in Def.3, the covering number of $\Theta$ with radius $\epsilon$ is defined as:*

$$\mathfrak{C}(\epsilon, \Theta, d) = \min \{n : \exists \epsilon - covering \ over \ \Theta \ with \ size \ n\}.$$

**Restate of Theorem 5** (Chaining Bounds for MAUC$^\downarrow$ Rademacher Complexity). *Suppose that the score function $s^{(i)}$ maps $\mathcal{X}$ onto a bounded interval $[-R_s, R_s]$, the following properties hold for $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H})$:*

(a) *For a decreasing precision sequence $\{\epsilon_k\}_{k=1}^K$, with $\epsilon_{k+1} = \frac{1}{2}\epsilon_k$, $k = 1, 2, \cdots, K-1$ and $\epsilon_0 \geq R_s$, we have:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}) \leq N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_K + 6 \cdot \sum_{k=1}^K \epsilon_k \phi_\ell (N_C - 1) \cdot \xi(\boldsymbol{Y}) \sqrt{\frac{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}}$$

(b) *There exists a universal constant $C$, such that:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}) \leq C \cdot \inf_{R_s \geq \alpha \geq 0} \left( N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \alpha + \phi_\ell \cdot (N_C - 1) \cdot \xi(\boldsymbol{Y}) \cdot \int_\alpha^{R_s} \sqrt{\frac{\log(\mathfrak{C}(\epsilon, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}} d\epsilon \right)$$

**Proof.** *First of all, we have:*

$$
\begin{aligned}
2C_G \hat{\mathfrak{R}}_{\mathsf{MAUC\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}) &= C_G \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{H}} T_f(\sigma) + \sup_{\tilde{f} \in \mathcal{H}} T_{\tilde{f}}(\sigma) \right] \\
&= C_G \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{H}} T_f(\sigma) + \sup_{\tilde{f} \in \mathcal{H}} T_{\tilde{f}}(-\sigma) \right] \\
&= C_G \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{H}} T_f(\sigma) + \sup_{\tilde{f} \in \mathcal{H}} -T_{\tilde{f}}(\sigma) \right] \\
&= C_G \mathbb{E}_\sigma \left[ \sup_{f, \tilde{f} \in \mathcal{H}} \left( T_f(\sigma) - T_{\tilde{f}}(\sigma) \right) \right]
\end{aligned}
\tag{44}
$$

*Now define $\hat{\mathcal{H}}$ as an $\epsilon$-cover of $\mathcal{H}$ with the pseudo-metric $d_{\infty,\mathcal{S}}$. Choose $\hat{f}, \hat{\tilde{f}} \in \hat{\mathcal{H}}$, such that $d_{\infty,\mathcal{S}}(f, \hat{f}) \leq \epsilon$ and $d_{\infty,\mathcal{S}}(\tilde{f}, \hat{\tilde{f}}) \leq \epsilon$. We then have the following result:*

$$
\begin{aligned}
& C_G \mathbb{E}_\sigma \left[ \sup_{f, \tilde{f} \in \mathcal{H}} \left( T_f(\sigma) - T_{\tilde{f}}(\sigma) \right) \right] \\
&= C_G \mathbb{E}_\sigma \left[ \sup_{f, \tilde{f} \in \mathcal{H}} \left( T_f(\sigma) - T_{\hat{f}}(\sigma) \right) + \left( T_{\hat{f}}(\sigma) - T_{\hat{\tilde{f}}}(\sigma) \right) + \left( T_{\hat{\tilde{f}}}(\sigma) - T_{\tilde{f}}(\sigma) \right) \right] \\
&\leq 2 C_G \mathbb{E}_\sigma \left[ \sup_{d_{\infty,\mathcal{N}}(\boldsymbol{s}, \hat{\boldsymbol{s}}) \leq \epsilon} \left( T_f(\sigma) - T_{\hat{f}}(\sigma) \right) \right] + C_G \mathbb{E}_\sigma \left[ \sup_{\hat{\boldsymbol{s}}, \hat{\tilde{\boldsymbol{s}}} \in \hat{\mathcal{H}}_\epsilon} \left( T_{\hat{f}}(\sigma) - T_{\hat{\tilde{f}}}(\sigma) \right) \right]
\end{aligned}
\tag{45}
$$

*Furthermore, let $\hat{\mathcal{H}}_k$ be an $\epsilon_k$-cover of $\mathcal{H}$. For each $k$, pick $\hat{\boldsymbol{s}}_k$, $\hat{\tilde{\boldsymbol{s}}}_k$ such that $d_{\infty,\mathcal{S}}(\hat{\boldsymbol{s}}_k, \boldsymbol{s}) \leq \epsilon_k$ and $d_{\infty,\mathcal{S}}(\hat{\tilde{\boldsymbol{s}}}_k, \tilde{\boldsymbol{s}}) \leq \epsilon_k$. Specifically, we take $\epsilon_K = \epsilon$, $\epsilon_0 \geq R_s$, which allows us to choose $\hat{\boldsymbol{s}}_0 = \hat{\tilde{\boldsymbol{s}}}_0$, $\epsilon_{k+1} = \frac{1}{2}\epsilon_k$, $\hat{f}_K = \hat{f}$, and $\hat{\tilde{f}}_K = \hat{\tilde{f}}$. With $\hat{f}_k = \mathsf{soft} \circ \hat{\boldsymbol{s}}_k$ and $\hat{\tilde{f}}_k = \mathsf{soft} \circ \hat{\tilde{\boldsymbol{s}}}_k$. Then, for $\boldsymbol{s}, \tilde{\boldsymbol{s}} \in \hat{\mathcal{H}}_\epsilon$, we can decompose $T_{\hat{f}}(\sigma)$ as*

$$T_{\hat{f}}(\sigma) = T_{\hat{f}_K}(\sigma) = T_{\hat{f}_0}(\sigma) + \sum_{i=1}^K \left( T_{\hat{f}_k}(\sigma) - T_{\hat{f}_{k-1}}(\sigma) \right) \quad and \quad T_{\hat{\tilde{f}}}(\sigma) = T_{\hat{\tilde{f}}_K}(\sigma) = T_{\hat{\tilde{f}}_0}(\sigma) + \sum_{i=1}^K \left( T_{\hat{\tilde{f}}_k}(\sigma) - T_{\hat{\tilde{f}}_{k-1}}(\sigma) \right).$$

*Based on the decompositions, we have:*

$$C_G \mathbb{E}_\sigma \left[ \sup_{\hat{\boldsymbol{s}}, \hat{\tilde{\boldsymbol{s}}} \in \hat{\mathcal{H}}_\epsilon} \left( T_{\hat{f}}(\sigma) - T_{\hat{\tilde{f}}}(\sigma) \right) \right] \leq 2 \cdot C_G \cdot \sum_{i=1}^K \mathbb{E}_\sigma \left[ \sup_{\substack{\hat{\boldsymbol{s}}_k \in \hat{\mathcal{H}}_k, \hat{\boldsymbol{s}}_{k-1} \in \hat{\mathcal{H}}_{k-1} \\ d_{\infty,\mathcal{S}}(\hat{\boldsymbol{s}}_k, \hat{\boldsymbol{s}}_{k-1}) \leq 3\epsilon_k}} \left( T_{\hat{f}_k}(\sigma) - T_{\hat{f}_{k-1}}(\sigma) \right) \right]. \tag{46}$$

*According to maximal inequality (Lem.D.1.2), we have:*

$$C_G \mathbb{E}_\sigma \left[ \sup_{\substack{\hat{\boldsymbol{s}}_k \in \hat{\mathcal{H}}_k, \hat{\boldsymbol{s}}_{k-1} \in \hat{\mathcal{H}}_{k-1} \\ d_{\infty,\mathcal{S}}(\hat{\boldsymbol{s}}_k, \hat{\boldsymbol{s}}_{k-1}) \leq 3\epsilon_k}} \left( T_{\hat{f}_k}(\sigma) - T_{\hat{f}_{k-1}}(\sigma) \right) \right] \leq 3\epsilon_k \sqrt{2 \log |\hat{\mathcal{H}}_k| \cdot |\hat{\mathcal{H}}_{k-1}|} \leq 6\epsilon_k \sqrt{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}. \tag{47}$$

*Following a similar derivation to Lem.9, we have:*

$$\mathbb{E}_\sigma \left[ \sup_{d_{\infty,\mathcal{N}}(\boldsymbol{s}, \hat{\boldsymbol{s}}) \leq \epsilon_K} \left( T_f(\sigma) - T_{\hat{f}}(\sigma) \right) \right] \leq N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_K \tag{48}$$

*Putting all together, we have:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H}) = \frac{1}{2} \cdot \mathbb{E}_\sigma \left[ \sup_{f,\tilde{f} \in \mathcal{H}} \left( T_f(\sigma) - T_{\tilde{f}}(\sigma) \right) \right] \le N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_K + \left( \frac{1}{C_G} \right) 6 \sum_{i=1}^{K} \epsilon_k \sqrt{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}$$

$$\le N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_K + 6 \sum_{i=1}^{K} \epsilon_k \phi_\ell \cdot (N_C - 1)\xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}}$$

*This ends the proof of (a).*

*Given the result of (a), we now turn to proof (b). First note that $\epsilon_k = 2(\epsilon_k - \epsilon_{k+1})$ and that $\log(\mathfrak{C}(\epsilon, \mathcal{F}, d_{\infty,\mathcal{S}})$ is non-increasing with respect to $\epsilon$. We then have:*

$$N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_K + 6 \sum_{i=1}^{K} \epsilon_k \phi_\ell \cdot (N_C - 1) \cdot \xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}}$$

$$\le 2N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \epsilon_{K+1} + 12 \sum_{i=1}^{K} (\epsilon_k - \epsilon_{k+1}) \cdot (N_C - 1)\phi_\ell \cdot \xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}}$$

*Furthermore, we have:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathcal{H}) \le 12\phi_\ell \left( N_C \cdot (N_C - 1) \cdot \epsilon_{K+1} + \sum_{i=1}^{K} (\epsilon_k - \epsilon_{k+1}) \cdot (N_C - 1) \cdot \xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon_k, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}} \right)$$

$$\le 12\phi_\ell \left( N_C \cdot (N_C - 1) \cdot \epsilon_{K+1} + \int_{\epsilon_{K+1}}^{R_s} (N_C - 1)\xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}} d\epsilon \right.$$

$$\left. + \int_{R_s}^{\epsilon_0} (N_C - 1)\xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}} d\epsilon \right)$$

$$= 12\phi_\ell \left( N_C \cdot (N_C - 1) \cdot \epsilon_{K+1} + \int_{\epsilon_{K+1}}^{R_s} \xi(\boldsymbol{Y})\sqrt{\frac{\log(\mathfrak{C}(\epsilon, \mathcal{F}, d_{\infty,\mathcal{S}}))}{N}} d\epsilon \right)$$

It then becomes clear that the universal constant should be chosen as $C = 12$. Set $\alpha = \epsilon_{K+1}$. By proper choices of $K$ and $\epsilon_0$ and let $\alpha = \frac{\epsilon_0}{2^{K+1}}$, we see that the inequality holds for all $\alpha \in [0, R_s]$, which thus completes the proof. $\qquad \square$

According to Thm.5, we have the following theorem as a result of a new minorization technique which appears in a recent work [65].

**Restate of Theorem 6** (Transformation Upper Bound). *Given the Hypothesis class*

$$\mathsf{soft} \circ \mathcal{F} = \left\{ \boldsymbol{g}(\boldsymbol{x}) = \mathsf{soft}(\boldsymbol{s}(\boldsymbol{x})) : \ \boldsymbol{s} \in \mathcal{F} \right\},$$

*where $\mathsf{soft}(\cdot)$ is the softmax function. Suppose that $\boldsymbol{s}(x) \subseteq [-R_s, R_s]^{N_C}$ and $\ell$ is $\phi_\ell$-Lipschitz continuous, the following inequality holds:*

$$\frac{\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{F})}{N_C(N_C - 1)} \le \phi_\ell \left( 2^9 \cdot \frac{1}{N_C} \cdot \sqrt{N_C} \cdot \xi(\boldsymbol{Y}) \cdot \log^{3/2}(e \cdot R_s \cdot N \cdot N_C) \cdot \hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F}) + \sqrt{\frac{1}{N}} \right), \tag{49}$$

*where the Rademacher complexity $\hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F})$ is defined as:*

$$\hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F})$$
$$= \mathbb{E}_\sigma \left[ \sup_{f=(f^{(1)}, \cdots, f^{(N_C)}) \in \mathcal{F}} \frac{1}{N \cdot N_C} \sum_{j=1}^{N_C} \sum_{i=1}^{N} \sigma_j^{(i)} \cdot f^{(j)}(\boldsymbol{x}_i) \right], \tag{50}$$

*where $\{\sigma_j^{(i)}\}_{(i,j)}$ is a sequence of independent Rademacher random variables.*

**Proof.** *With the help of the chaining bound in Thm.5-a), the proof follows an analogous spirit as [65, Prop.1], with the parameters therein set as $\lambda = 1$, $\theta = 0$, $q = N_C$, $n = N$.*

The result in Thm.6 directly relates the complicated pairwise Rademacher complexity $\hat{\mathfrak{R}}_{\mathsf{MAUC}^\downarrow,\mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{F})$ to the instance-wise Rademacher complexity $\hat{\mathfrak{R}}_{N \cdot N_C}$. This makes the derivation of the generalization upper bound much easier. Once a bound on the ordinary Rademacher complexity over the functional class $\mathcal{F}$ is available, we can directly plug it in to this theorem and find a resulting bound over the MAUC$^\downarrow$ complexity.

# F PRACTICAL GENERALIZATION BOUNDS

## F.1 Practical Results for Linear Models

In this section, we provide concrete generalization bounds for two practical models.

**Restate of Theorem 7.** *Based on the assumptions of Thm.4, we have the Given the hypothesis space for $\ell_p$ norm penalized linear model as:*

$$\mathcal{H}_{p,\gamma}^{Lin} = \big\{ f = (f^{(1)}, \cdots, f^{(n_C)}) : f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^{(i)}\boldsymbol{x},$$
$$||\boldsymbol{W}^{(i)}||_p \leq \gamma \big\},$$

*if $\ell$ is $\phi_\ell$-Lipschitz continuous, and the input features are sampled from $\mathcal{X} \subset \mathcal{R}^d$, and for all $\boldsymbol{x} \in \mathcal{X}$, we have $||\boldsymbol{x}||_2^2 \leq R$ with $0 < p < \infty$, $\frac{1}{p} + \frac{1}{\bar{p}} = 1$, for all $f \in \mathcal{H}_{p,\gamma}^{Lin}$, we have the following inequality holds with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\ell,\mathcal{S}}(f) + \mathcal{I}_{Lin}\Big( \chi(\boldsymbol{Y}), \xi(\boldsymbol{Y}), \delta \Big) \cdot \sqrt{\frac{1}{N}},$$

*where*

$$\mathcal{I}_{Lin}\Big( \chi(\boldsymbol{Y}), \xi(\boldsymbol{Y}), \delta \Big) = \frac{4R_\mathcal{X}\phi_\ell\gamma}{N_C - 1} \cdot \sqrt{\frac{2(\bar{p}-1)}{N_C}} \cdot \chi(\boldsymbol{Y})$$
$$+ \frac{5B}{N_C} \cdot \sqrt{2\log\left(\frac{2}{\delta}\right)} \cdot \xi(\boldsymbol{Y}).$$

**Proof.** *Follows Thm.4, Lem.10*

## F.2 Practical Results for Deep Fully-connected Neural Networks

**Settings**. We denote an $L$-layer deep fully-connected neural network with $N_C$-way output as :

$$f(\boldsymbol{x}) = \boldsymbol{W} f_{\boldsymbol{\omega},L}(\boldsymbol{x}) = \boldsymbol{W} s(\boldsymbol{\omega}_{L-2} \cdots s(\boldsymbol{\omega}_1 \boldsymbol{x})),$$

The notations are as follows: $s(\cdot)$ is the activation function; $n_{h_j}$ is the number of hidden neurons for the $j$-th layer; $\boldsymbol{\omega}_j \in \mathbb{R}^{n_{h_{j+1}} \times n_{h_j}}$, $j = 1, 2, \cdots, L-2$ are the weights for the first $L-1$ layers; $\boldsymbol{W} \in \mathbb{R}^{n_{h_{L-1}} \times N_C}$ is the weight for the output layer. Moreover, the $i$-th output of the network is defined as :

$$f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^{(i)^\top} f_{\boldsymbol{\omega},L}(\boldsymbol{x}).$$

where $\boldsymbol{W}^{(i)^\top}$ is the i-th row of $\boldsymbol{W}$. In the next theorem, we focus on a specific hypothesis class for such networks where the product of weight norms

$$\Pi_{\boldsymbol{W},\boldsymbol{\omega}} = ||\boldsymbol{W}||_F \cdot \prod_{j=1}^{L-2} ||\boldsymbol{\omega}_j||_F$$

are no more than $\gamma$. We denote such a hypothesis class as $\mathcal{H}_{\gamma,R_s,n_h}^{DNN}$:

$$\mathcal{H}_{\gamma,R_s,n_h}^{DNN} = \left\{ f : f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^{(i)^\top} f_{\boldsymbol{\omega},L}(\boldsymbol{x}), ||f^{(i)}||_\infty \leq R_s, \ i = 1, \cdots, N_C, \ \Pi_{\boldsymbol{W},\boldsymbol{\omega}} \leq \gamma \right\},$$

To obtain the final output, we perform a softmax operation over $f(\boldsymbol{x})$. Thus, a valid model under this setting could be chosen from the following hypothesis class:

$$\mathsf{soft} \circ \mathcal{H}_{\gamma,R_s,n_h}^{DNN} = \left\{ g : g^{(i)} = \frac{\exp(f^{(i)}(\boldsymbol{x}))}{\sum_{j=1}^{N_C} \exp(f^{(j)}(\boldsymbol{x}))}, \ f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN} \right\}.$$

**Restate of Theorem 8 (Practical Generalization Bounds for Deep Models).** *Based on the assumptions of Thm.4 where $s(\cdot)$ is a 1-Lipshiptz and positive homogeneous activation function. If $\ell$ is $\phi_\ell$-Lipschitz continuous, and the input features are sampled from $\mathcal{X} \subset \mathbb{R}^d$, and for all $\boldsymbol{x} \in \mathcal{X}$, we have $||\boldsymbol{x}||_2^2 \leq R_\mathcal{X}$, for all $f \in \mathsf{soft} \circ \mathcal{H}_{\gamma,R_s,n_h}^{DNN}$, we have the following inequality holds with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_\mathcal{S}(f) + \min\Big( \mathcal{I}_{DNN,1}, \ \mathcal{I}_{DNN,2} \Big) \cdot \sqrt{\frac{1}{N}}$$

*where $\chi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} \sum_{j \neq i} 1/\rho_i\rho_j}$, $\xi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} 1/\rho_i}$, $\rho_i = \frac{n_i}{N}$,*

$$\mathcal{I}_{DNN,1} = C_1 \frac{\sqrt{2}}{2}\phi_\ell \cdot \frac{\chi(\boldsymbol{Y})}{N_C - 1} + \left( \frac{\sqrt{2}C_1 R_\mathcal{X}\phi_\ell\gamma}{2} \cdot C_3 + \frac{C_2 B}{N_C} \cdot \sqrt{2\log\left(\frac{2}{\delta}\right)} \right) \cdot \xi(\boldsymbol{Y}),$$

$$\mathcal{I}_{DNN,2} = C_1 \phi_\ell \left( \frac{2^9}{N_C} \cdot \xi(\boldsymbol{Y}) \cdot \log^{3/2}(K \cdot N \cdot N_C) \gamma \cdot R_{\mathcal{X}} \cdot (\sqrt{2\log(2)L} + 1) + 1 \right) + C_2 \frac{B \cdot \sqrt{\log(2/\delta)} \cdot \xi(\boldsymbol{Y})}{N_C}$$

$C_1, C_2$ are universal constants as thm.4, $K = e \cdot R_s$, $C_3 = \frac{\sqrt{L \log 2} + \sqrt{N_C}}{\sqrt{N_C - 1}}$.

**Proof.** Follows Thm.4, Lem.11, Thm.6 and the fact that $\hat{\mathfrak{R}}_{N \cdot N_C}(\Pi \circ \mathcal{F}) \leq \frac{R_{\mathcal{X}}\gamma(\sqrt{2\log(2)L}+1)}{N \cdot N_C}$ from Thm.1 in [27].

### F.3 Practical Results for Deep Convolutional Neural Networks

Now we use the result in Thm.5 to derive a generalization bound for a general class of deep neural networks where fully-connected layers and convolutional layers coexist. In a nutshell, the result is essentially an application of Thm.5 to a recent idea appeared in [49].

**General Setting.** Now we are ready to introduce the setting of the deep neural networks employed in the forthcoming theoretical analysis, which is adopted from [49]. We focus on the deep neural networks with $N_{conn}$ fully-connected layers and $N_{conv}$ convolutional layers. The $i$-th convolutional layer has a kernel $\boldsymbol{K}^{(i)} \in \mathbb{R}^{k_i \times k_i \times c_{i-1} \times c_i}$. Recall that convolution is a linear operator. For a Given kernel $\boldsymbol{K}$, we denote its associated matrix as $op(\boldsymbol{K})$, such that $\boldsymbol{K}(\boldsymbol{x}) = op(\boldsymbol{K})\boldsymbol{x}$. Moreover, we assume that, each time, the convolution layer is followed by a componentwise non-linear activation function and an optional pooling operation. We assume that the activation functions and the pooling operations are all 1-Lipschitz. For the $i$-th fully-connected layer, we denote its weight as $\boldsymbol{V}^{(i)}$. Above all, the complete parameter set of a given deep neural network could be represented as $\boldsymbol{P} = \{\boldsymbol{K}^{(1)}, \cdots, \boldsymbol{K}^{(N_{conv})}, \boldsymbol{V}^{(1)}, \cdots, \boldsymbol{V}^{(N_{conn})}\}$. Again, we also assume that the loss function is $\phi_\ell$-Lipschitz and $Range\{\ell\} \subseteq [0, B]$. Finally, given two deep neural networks with parameters $\boldsymbol{P}$ and $\tilde{\boldsymbol{P}}$, we adopt a metric $d_{NN}(\cdot, \cdot)$ to measure their distance:

$$d_{NN}(\boldsymbol{P}, \tilde{\boldsymbol{P}}) = \sum_{i=1}^{N_{conv}} ||op(\boldsymbol{K}^{(i)}) - op(\tilde{\boldsymbol{K}}^{(i)})||_2 + \sum_{i=1}^{N_{conn}} ||\boldsymbol{V}^{(i)} - \tilde{\boldsymbol{V}}^{(i)}||_2$$

**Constraints Over the Parameters.** First, we define $\mathcal{P}_\nu^{(0)}$ as the class for <u>initialization of the parameters</u>:

$$\mathcal{P}_\nu^{(0)} = \left\{ \boldsymbol{P} : \left( \max_{i \in \{1, \cdots, N_{conv}\}} ||op(\boldsymbol{K}^{(i)})||_2 \right) \leq 1 + \nu, \ \left( \max_{j \in \{1, \cdots, N_{conn}\}} ||\boldsymbol{V}^{(j)}||_2 \right) \leq 1 + \nu \right\}$$

Now we further assume that the learned parameters should be chosen from a class denoted by $\mathcal{P}_{\beta,\nu}$, where the distance between the learned parameter and the fixed initialization residing in $\mathcal{P}_\nu^{(0)}$ is no bigger than $\beta$:

$$\mathcal{P}_{\beta,\nu} = \left\{ \boldsymbol{P} : d_{NN}(\boldsymbol{P}, \tilde{\boldsymbol{P}}_0) \leq \beta, \ \tilde{\boldsymbol{P}}_0 \in \mathcal{P}_\nu^{(0)} \right\}.$$

The following result then provides a generalization upper bound based on Thm.4 and Lem.F.4.3.

**Restate of Theorem 9.** *Based on the setting in Lem.F.4.3, given dataset $\mathcal{S} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^m$, where the instances are sampled independently, for all multiclass scoring functions $f \in$ soft $\circ \mathcal{F}_{\beta,\nu}$, if $Range\{\ell\} \subseteq [0, B]$, $\forall \delta \in (0, 1)$, the following inequalities hold with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\mathcal{S}}(f) + C_1 \cdot \left( \frac{\tilde{C} \cdot \phi_\ell \cdot R_s \cdot \xi(\boldsymbol{Y})}{N_C} \cdot \sqrt{\frac{N_{par}(\nu N_L + \beta + \log(3\beta R_{\mathcal{X}} N))}{N}} \right) + C_2 \cdot \frac{B \cdot \xi(\boldsymbol{Y})}{N_C} \cdot \sqrt{\frac{\log(2/\delta)}{N}}$$

*where $C_1, C_2, \tilde{C}$ are universal constants.*

### F.4 Essential Lemmas

*F.4.1 MAUC$^\downarrow$ Rademacher Complexity for p-norm penalized Linear Models*

**Lemma 10** (MAUC$^\downarrow$ Rademacher Complexity for Linear Model). *We have:*

$$\hat{\mathfrak{R}}_{MAUC^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H}_{p,\gamma}^{Lin}) \leq R_{\mathcal{X}} \phi_\ell \gamma \cdot \sqrt{2(\bar{p} - 1) \cdot N_C} \cdot \chi(\boldsymbol{Y}) \sqrt{\frac{1}{N}} \tag{51}$$

*where $\chi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} \sum_{j \neq i} \frac{1}{\rho_i, \rho_j}}$, $\rho_i = \frac{n_i}{N}$.*

**Proof.**

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}_{p,\gamma}^{Lin}) = \mathbb{E}_{\sigma}\left[\sup_{f\in\mathcal{H}} \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \frac{\sigma_m^{(i)}+\sigma_n^{(j)}}{2}\cdot\frac{\ell(f^{(i)}(\boldsymbol{x}_m)-f^{(i)}(\boldsymbol{x}_n))}{n_i n_j}\right]$$

$$= \mathbb{E}_{\sigma}\left[\sup_{\boldsymbol{W},||\boldsymbol{W}^{(i)}||_p\leq\gamma} \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \frac{\sigma_m^{(i)}+\sigma_n^{(j)}}{2}\cdot\frac{1}{n_i n_j}\cdot\ell(\boldsymbol{W}^{(i)\top}(\boldsymbol{x}_m-\boldsymbol{x}_n))\right]$$

$$\leq \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \mathbb{E}_{\sigma^{(i)}}\left[\sup_{\substack{\boldsymbol{W}^{(i)}\\||\boldsymbol{W}^{(i)}||_p\leq\gamma}} \sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \frac{\sigma_m^{(i)}}{2}\cdot\frac{1}{n_i n_j}\cdot\ell(\boldsymbol{W}^{(i)\top}(\boldsymbol{x}_m-\boldsymbol{x}_n))\right]$$

$$+ \sum_{i=1}^{N_C}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \mathbb{E}_{\sigma^{(j)}}\left[\sup_{\substack{\boldsymbol{W}^{(i)}\\||\boldsymbol{W}^{(i)}||_p\leq\gamma}} \sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \frac{\sigma_n^{(j)}}{2}\cdot\frac{1}{n_i n_j}\cdot\ell(\boldsymbol{W}^{(i)\top}(\boldsymbol{x}_m-\boldsymbol{x}_n))\right]$$

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}_{p,\gamma}^{Lin}) \overset{(*)}{\leq} \phi_\ell \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \mathbb{E}_{\sigma^{(i)}}\left[\sup_{\substack{\boldsymbol{W}^{(i)}\\||\boldsymbol{W}^{(i)}||_p\leq\gamma}} \sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \frac{\sigma_m^{(i)}}{2}\cdot\frac{1}{n_i n_j}\cdot\boldsymbol{W}^{(i)\top}(\boldsymbol{x}_m-\boldsymbol{x}_n)\right]$$

$$+ \phi_\ell \sum_{i=1}^{N_C}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \mathbb{E}_{\sigma^{(j)}}\left[\sup_{\substack{\boldsymbol{W}^{(i)}\\||\boldsymbol{W}^{(i)}||_p\leq\gamma}} \sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \frac{\sigma_n^{(j)}}{2}\cdot\frac{1}{n_i n_j}\cdot\boldsymbol{W}^{(i)\top}(\boldsymbol{x}_m-\boldsymbol{x}_n)\right]$$

*That is to say,*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}_{p,\gamma}^{Lin}) \leq \phi_\ell \cdot \gamma \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \mathbb{E}_{\sigma^{(i)}}\left[\left\|\sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \frac{\sigma_m^{(i)}}{2}\cdot\frac{1}{n_i n_j}\cdot(\boldsymbol{x}_m-\boldsymbol{x}_n)\right\|_{\bar{p}}\right]$$

$$+ \phi_\ell \cdot \gamma \sum_{i=1}^{N_C}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \mathbb{E}_{\sigma^{(j)}}\left[\left\|\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \frac{\sigma_n^{(j)}}{2}\cdot\frac{1}{n_i n_j}\cdot(\boldsymbol{x}_m-\boldsymbol{x}_n)\right\|_{\bar{p}}\right]$$

$$\overset{(**)}{\leq} \phi_\ell \cdot \gamma \sum_{i=1}^{N_C}\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \left(\mathbb{E}_{\sigma^{(i)}}\left[\left\|\sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \frac{\sigma_m^{(i)}}{2}\cdot\frac{1}{n_i n_j}\cdot(\boldsymbol{x}_m-\boldsymbol{x}_n)\right\|_{\bar{p}}^{\bar{p}}\right]\right)^{1/\bar{p}} \tag{52}$$

$$+ \phi_\ell \cdot \gamma \sum_{i=1}^{N_C}\sum_{\boldsymbol{x}_m\in\mathcal{N}_i} \mathbb{E}_{\sigma^{(j)}}\left[\left(\left\|\sum_{j\neq i}\sum_{\boldsymbol{x}_n\in\mathcal{N}_j} \frac{\sigma_n^{(j)}}{2}\cdot\frac{1}{n_i n_j}\cdot(\boldsymbol{x}_m-\boldsymbol{x}_n)\right\|_{\bar{p}}^{\bar{p}}\right)\right]^{1/\bar{p}},$$

*where* $(*)$ *is due to the Talagrand contraction lemma,* $(**)$ *follows that fact that* $()^{1/\bar{p}}$ *is concave for* $\bar{p} > 1$ *and Jensen's Inequality.*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathcal{H}_{p,\gamma}^{Lin}) \overset{(***)}{\leq} \phi_{\ell} \cdot \gamma \sqrt{\bar{p}-1} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \underset{\sigma^{(i)}}{\mathbb{E}} \left[ \left\| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (\boldsymbol{x}_m - \boldsymbol{x}_n) \right\|_2^2 \right]^{1/2}$$

$$+ \phi_{\ell} \cdot \gamma \cdot \sqrt{\bar{p}-1} \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \underset{\sigma^{(j)}}{\mathbb{E}} \left[ \left\| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_n^{(j)}}{2} \cdot \frac{1}{n_i n_j} \cdot (\boldsymbol{x}_m - \boldsymbol{x}_n) \right\|_2^2 \right]^{1/2}$$

$$\overset{(****)}{\leq} \phi_{\ell} \cdot \gamma \sqrt{\bar{p}-1} \sqrt{N \cdot N_C} \left( \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \underset{\sigma^{(i)}}{\mathbb{E}} \left[ \left\| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (\boldsymbol{x}_m - \boldsymbol{x}_n) \right\|_2^2 \right] \right.$$

$$\left. + \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \underset{\sigma^{(j)}}{\mathbb{E}} \left[ \left\| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_n^{(j)}}{2} \cdot \frac{1}{n_i n_j} \cdot (\boldsymbol{x}_m - \boldsymbol{x}_n) \right\|_2^2 \right] \right)^{1/2}$$

$$\leq R_{\mathcal{X}} \phi_{\ell} \gamma \sqrt{2\frac{\bar{p}-1}{N}} \cdot \sqrt{N_C \sum_{i=1}^{N_C} \cdot \sum_{j \neq i} \frac{1}{\rho_i \rho_j}},$$

*where* $(***)$ *is due to Lem.4, and* $(****)$ *is due to* $\forall \boldsymbol{x} \in \mathbb{R}^N \ \|\boldsymbol{x}\|_1 \leq \sqrt{N} \|\boldsymbol{x}\|_2$, *and the last inequality follows that :*

$$\underset{\sigma^{(i)}}{\mathbb{E}} \left[ \left\| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (\boldsymbol{x}_m - \boldsymbol{x}_n) \right\|_2^2 \right] \leq \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{1}{4n_i^2 n_j^2} \|\boldsymbol{x}_m - \boldsymbol{x}_n\|_2^2 \leq \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{R_{\mathcal{X}}^2}{n_i^2 n_j^2}$$

$$\underset{\sigma^{(j)}}{\mathbb{E}} \left[ \left\| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_n^{(j)}}{2} \cdot \frac{1}{n_i n_j} \cdot (\boldsymbol{x}_m - \boldsymbol{x}_n) \right\|_2^2 \right] \leq \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{4n_i^2 n_j^2} \|\boldsymbol{x}_m - \boldsymbol{x}_n\|_2^2 \leq \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{R_{\mathcal{X}}^2}{n_i^2 n_j^2}$$

### F.4.2 MAUC$^{\downarrow}$ Rademacher Complexity for A class of Deep Fully-connected Neural Networks

**Lemma 11.** *we have:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}) \leq \frac{\sqrt{2}}{2} \phi_{\ell} R_{\mathcal{X}} \gamma N_C \xi(\boldsymbol{Y}) \left( \sqrt{\frac{L \log 2(N_C-1)}{N}} + \sqrt{\frac{N_C(N_C-1)}{N}} \right)$$

$$+ \frac{\sqrt{2}}{2} \phi_{\ell} N_C \chi(\boldsymbol{Y}) \sqrt{\frac{1}{N}}$$

*where* $\chi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} \sum_{j \neq i} \frac{1}{\rho_i, \rho_j}}, \xi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} \frac{1}{\rho_i}}, \quad \rho_i = \frac{n_i}{N}.$

**Proof.** *Similar to Lem.10, we can achieve the bound:*

$$\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}) \leq \phi_{\ell} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \underset{\sigma^{(i)}}{\mathbb{E}} \left[ \sup_{g \in \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (g^{(i)}(\boldsymbol{x}_m) - g^{(i)}(\boldsymbol{x}_n)) \right]$$

$$+ \phi_{\ell} \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \underset{\sigma^{(j)}}{\mathbb{E}} \left[ \sup_{g \in \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_n^{(j)}}{2} \cdot \frac{1}{n_i n_j} \cdot (g^{(i)}(\boldsymbol{x}_m) - g^{(i)}(\boldsymbol{x}_n)) \right]$$

*By the sub-additivity of supremum, we have:*

$$
\hat{\mathfrak{R}}_{\mathsf{MAUC}^{\downarrow},\mathcal{S}}(\ell \circ \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}) \leq \underbrace{\phi_{\ell} \sum_{i=1}^{N_C}(N_C-1) \cdot \mathbb{E}_{\sigma^{(i)}}\left[\sup_{g \in \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i} \cdot (g^{(i)}(\boldsymbol{x}_m))\right]}_{(a)}
$$

$$
+ \underbrace{\phi_{\ell} \sum_{i=1}^{N_C} \mathbb{E}_{\sigma^{(j)}}\left[\sup_{g \in \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_n^{(j)}}{2} \cdot \frac{1}{n_j} \cdot (g^{(i)}(\boldsymbol{x}_n))\right]}_{(b)}
$$

(53)

$$
+ \underbrace{\phi_{\ell} \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \mathbb{E}_{\sigma^{(i)}}\left[\sup_{g \in \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}} (g^{(i)}(\boldsymbol{x}_n)) \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j}\right]}_{(c)}
$$

$$
+ \underbrace{\phi_{\ell} \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \mathbb{E}_{\sigma^{(j)}} \sup_{g \in \mathsf{soft} \circ \mathcal{H}_{\gamma}^{DNN,nh}}\left[(g^{(i)}(\boldsymbol{x}_m)) \cdot \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_n^{(j)}}{2} \cdot \frac{1}{n_i n_j}\right]}_{(d)}
$$

*We first derive the bound for $(a)+(b)$.*

*According Lem.5 and Lem.7, we know that:*

$$
(a)+(b) \leq \phi_{\ell} \sum_{i=1}^{N_C}(N_C-1) \cdot \mathbb{E}_{\sigma^{(i)}}\left[\sup_{f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN}} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i} \cdot (f^{(c)}(\boldsymbol{x}_m))\right]
$$

$$
+ \phi_{\ell} \sum_{i=1}^{N_C} \mathbb{E}_{\sigma^{(j)}}\left[\sup_{f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN}} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(j)}}{2} \cdot \frac{1}{n_j} \cdot (f^{(c)}(\boldsymbol{x}_n))\right]
$$

*Using the Lem.6 recursively towards hierarchical structure of DNN, we have:*

$$
\mathbb{E}_{\sigma^{(i)}}\left[\sup_{f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN}} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (f^{(c)}(\boldsymbol{x}_m))\right]
$$

$$
\leq \frac{1}{\lambda} \log\left(\exp\left(\lambda \cdot \mathbb{E}_{\sigma^{(i)}}\left[\sup_{f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN}} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (f^{(c)}(\boldsymbol{x}_m))\right]\right)\right)
$$

$$
\leq \frac{1}{\lambda} \log\left(\left(\mathbb{E}_{\sigma^{(i)}}\left[\sup_{f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN}} \exp\left(\lambda \cdot \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot (f^{(c)}(\boldsymbol{x}_m))\right)\right]\right)\right)
$$

$$
\leq \frac{1}{\lambda} \log\left(2^L \cdot \mathbb{E}_{\sigma^{(i)}}\left[\exp\left(\lambda\gamma \left\|\sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \cdot \boldsymbol{x}_m\right\|\right)\right]\right)
$$

$$
\mathbb{E}_{\sigma^{(j)}}\left[\sup_{f \in \mathcal{H}_{\gamma,R_s,n_h}^{DNN}} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(j)}}{2} \cdot \frac{1}{n_i n_j} \cdot (f^{(c)}(\boldsymbol{x}_n))\right]
$$

$$
\leq \frac{1}{\lambda} \log\left(2^L \cdot \mathbb{E}_{\sigma^{(j)}}\left[\exp\left(\lambda\gamma \left\|\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(j)}}{2} \cdot \frac{1}{n_i n_j} \cdot \boldsymbol{x}_n\right\|\right)\right]\right).
$$

*This suggests that:*

$$(a+b) \leq \phi_\ell \sum_{i=1}^{N_C} \frac{1}{\lambda} \log \left( 2^L \cdot \mathbb{E}_\sigma \left[ \exp \left( \lambda \gamma (N_C - 1) \left\| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i} \cdot \boldsymbol{x}_m \right\| \right. \right. \right.$$

$$\left. \left. \left. + \lambda \gamma \left\| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(j)}}{2} \cdot \frac{1}{n_j} \cdot \boldsymbol{x}_n \right\| \right) \right] \right),$$

*which is due to the property that $\mathbb{E}_{\boldsymbol{a},\boldsymbol{b}}(f(\boldsymbol{a}) \cdot g(\boldsymbol{a})) = \mathbb{E}_{\boldsymbol{a}}(f(\boldsymbol{a})) \cdot \mathbb{E}_{\boldsymbol{b}}(g(\boldsymbol{b}))$, if $\boldsymbol{a}$ and $\boldsymbol{b}$ are independent. Denote*

$$\boldsymbol{Z}_i = \gamma (N_C - 1) \left\| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i} \cdot \boldsymbol{x}_m \right\| + \gamma \left\| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(j)}}{2} \cdot \frac{1}{n_j} \cdot \boldsymbol{x}_n \right\| \tag{54}$$

*as a random variable with respect to the Rademacher random variables. We have:*

$$(a) + (b) \leq \phi_\ell \left( \sum_{i=1}^{N_C} \frac{L \log 2 + \log(\mathbb{E}_\sigma \exp(\lambda(\boldsymbol{Z}_i - \mathbb{E}(\boldsymbol{Z}_i))))}{\lambda} + \sum_{i=1}^{N_C} \mathbb{E}_\sigma(\boldsymbol{Z}_i) \right).$$

*Similar to the proof of Lem.10, we can bound*

$$\sum_{i=1}^{N_C} \mathbb{E}_\sigma(\boldsymbol{Z}_i) \leq \frac{\sqrt{2}}{2} \gamma \cdot R_\mathcal{X} \cdot (N_C)^{3/2} \left( (N_C - 1) \sum_{i=1}^{N_C} \cdot \sum_{j \neq i} \frac{1}{\rho_i} \right)^{1/2} \cdot \frac{1}{\sqrt{N}}. \tag{55}$$

*Utilizing the bounded difference inequality (Lem.1) of $\boldsymbol{Z}_i$, we come to the following result:*

$$\mathbb{E}_\sigma \exp \left( \lambda(\boldsymbol{Z}_i - \mathbb{E}(\boldsymbol{Z}_i)) \right) \leq \exp \left( \frac{\lambda^2 v_i}{2} \right), \quad v_i \leq \frac{N_C}{4} \frac{(N_C - 1)^2}{n_i} \gamma^2 R_\mathcal{X}^2 + \frac{N_C}{4} \sum_{j \neq i} \frac{\gamma^2 R_\mathcal{X}^2}{n_j}, \tag{56}$$

*which leads to the bound:*

$$\left( \sum_{i=1}^{N_C} \frac{L \log 2 + \log(\mathbb{E}_\sigma \exp(\lambda(\boldsymbol{Z}_i - \mathbb{E}(\boldsymbol{Z}_i))))}{\lambda} \right) \leq \left( \sum_{i=1}^{N_C} \frac{L \log 2 + (\frac{\lambda^2 v_i}{2})}{\lambda} \right) \tag{57}$$

*By choosing $\lambda = \sqrt{\frac{2 N_C L \log 2}{\sum_{i=1}^{N_C} v_i}}$, we reach the optimal bound as:*

$$\sum_{i=1}^{N_C} \frac{L \log 2 + \log(\mathbb{E}_\sigma \exp(\lambda(\boldsymbol{Z}_i - \mathbb{E}(\boldsymbol{Z}_i))))}{\lambda}$$

$$\leq \sqrt{2 N_C L \log 2 \cdot \sum_{i=1}^{N_C} v_i} \tag{58}$$

$$= R \cdot \gamma \cdot \sqrt{\frac{L \log 2}{2N}} \cdot \sqrt{\sum_{i=1}^{N_C} \frac{N_C - 1}{\rho_i}}.$$

*Putting all together, we have*

$$(a) + (b) \leq \frac{\sqrt{2}}{2} \phi_\ell R_\mathcal{X} \gamma \cdot (N_C(N_C - 1))^{1/2} \xi(\boldsymbol{Y}) \left( \sqrt{\frac{N_C L \log 2}{N}} + N_C \sqrt{\frac{1}{N}} \right) \tag{59}$$

*For (c) and (d), we have:*

$$(c) \leq \phi_\ell \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \mathbb{E}_{\sigma^{(i)}} \left[ \left( \sup_{f \in \mathcal{H}_{\gamma, R_s, n_h}^{DNN}} \left| f^{(i)}(\boldsymbol{x}_n) \right| \right) \cdot \left| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j} \right| \right] \tag{60}$$

$$\overset{(*)}{\leq} \phi_\ell \cdot \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \mathbb{E}_{\sigma^{(i)}} \left[ \left| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{\sigma_m^{(i)}}{2} \cdot \frac{1}{n_i n_j} \right| \right]$$

*Similarly, we have :*

$$(d) \overset{(**)}{\leq} \phi_\ell \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \underset{\sigma^{(j)}}{\mathbb{E}} \left[ \left| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\sigma_m^{(j)}}{2} \cdot \frac{1}{n_i n_j} \right| \right],$$

(61)

*where (\*) and (\*\*) are due to the fact that the softmax function is upper bounded by 1, i.e.,*

$$\sup_{f \in \mathcal{H}_{\gamma, R_s, n_h}^{DNN}} |f_{W,L}^{(i)}(\boldsymbol{x}_n)| \leq 1 \quad and \quad \sup_{f \in \mathcal{H}_{\gamma, R_s, n_h}^{DNN}} |f_{W,L}^{(i)}(\boldsymbol{x}_m)| \leq 1.$$

*This finishes the proof for Eq.(60) and Eq.(61).*

*Now we can reach that:*

$$
\begin{aligned}
(c) + (d) &\leq \phi_\ell \cdot \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \underset{\sigma^{(i)}}{\mathbb{E}} \left[ \left| \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \right| \right] \\
&\quad + \phi_\ell \cdot \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \underset{\sigma^{(j)}}{\mathbb{E}} \left[ \left| \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \right| \right] \\
&\leq \phi_\ell \cdot \sum_{i=1}^{N_C} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sqrt{ \underset{\sigma^{(i)}}{\mathbb{E}} \left[ \left( \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{c=1}^{N_C} \frac{\sigma_{m,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \right)^2 \right] } \\
&\quad + \phi_\ell \cdot \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sqrt{ \underset{\sigma^{(j)}}{\mathbb{E}} \left[ \left( \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \sum_{c=1}^{N_C} \frac{\sigma_{n,c}^{(i)}}{2} \cdot \frac{1}{n_i n_j} \right)^2 \right] }
\end{aligned}
$$

(62)

*Similar to the proof of 10, we have:*

$$(c) + (d) \leq \frac{\sqrt{2}}{2} \cdot \phi_\ell \cdot \chi(\boldsymbol{Y}) \cdot N_C \sqrt{\frac{1}{N}}$$

(63)

Now we provide an alternative upper bound by the chaining technique.

**Lemma 12.** *Under the setting of Thm.8.*

*For all $f \in \mathsf{soft} \circ \mathcal{H}_{\gamma, R_s, n_h}^{DNN}$, we have the following inequality holds with probability at least $1 - \delta$:*

$$R_\ell(f) \leq \hat{R}_{\mathcal{S}}(f) + \mathcal{I}_{DNN,2} \cdot \sqrt{\frac{1}{N}}$$

*where $\xi(\boldsymbol{Y}) = \sqrt{\sum_{i=1}^{N_C} 1/\rho_i}$, $\rho_i = \frac{n_i}{N}$,*

$$\mathcal{I}_{DNN,2} = C_1 \phi_\ell \left( \frac{2^9}{N_C} \cdot \xi(\boldsymbol{Y}) \cdot \log^{3/2} (K \cdot N \cdot N_C) \cdot \gamma \cdot R_{\mathcal{X}} \cdot (\sqrt{2 \log(2)L} + 1) + 1 \right) + C_2 \frac{B \cdot \sqrt{\log(2/\delta)} \cdot \xi(\boldsymbol{Y})}{N_C}$$

*$C_1, C_2$ are universal constants as thm.4, $K = e \cdot R_s$.*

**Proof.** *The result follows from Thm.6, and the fact that*

$$\hat{\mathfrak{R}}_{N \cdot N_C} (\Pi \circ \mathcal{F}) \leq \frac{R_{\mathcal{X}} \gamma \cdot (\sqrt{2 \log(2)L} + 1)}{\sqrt{N_C \cdot N}},$$

*which follows from Thm.1 in [27]*

*F.4.3* MAUC$^\downarrow$ *Rademacher Complexity for A Class of Deep Convolutional Neural Networks*

**Lemma 13.** *Denote the hypothesis class,*

$$\text{soft} \circ \mathcal{F}_{\beta,\nu} = \left\{ \boldsymbol{g}(\boldsymbol{x}) = \text{soft}(\boldsymbol{s}_{\boldsymbol{P}}(\boldsymbol{x})) : \ \boldsymbol{s}_{\boldsymbol{P}} \in \mathcal{F}_{\beta,\nu} \right\},$$

$$\mathcal{F}_{\beta,\nu} = \{s_{\boldsymbol{P}} : \mathbb{R}^{N_{N_L-1}} \to \mathbb{R}^{N_C} | \ \boldsymbol{P} \in \mathcal{P}_{\beta,\nu}, Range(s_{\boldsymbol{P}}) \subseteq [-R_s, R_s]^{N_C}\}$$

*Moreover, define* $\tilde{N} = \dfrac{1}{\sum_{i=1}^{N_C} \dfrac{1}{n_i}}$ *we assume that* $\sup_{x \in \mathcal{X}} ||vec(\boldsymbol{x})|| \le R_{\mathcal{X}}$ *and*

$$R_s > 1/\min\left\{ \sqrt{N}, \frac{\xi(\boldsymbol{Y})}{N_C} \cdot \sqrt{N_{par}\left(\nu N_L + \beta + \log(3R_{\mathcal{X}} \cdot \beta \cdot N)\right)} \right\},$$

*we have:*

$$\hat{\mathfrak{R}}_{\text{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \text{soft} \circ \mathcal{F}_{\beta,\nu}) \le \tilde{C}\left( \phi_\ell \cdot (N_C - 1) \cdot R_s \cdot \xi(\boldsymbol{Y}) \cdot \sqrt{\frac{N_{par}\left(\nu N_L + \beta + \log\left(3\beta R_{\mathcal{X}} N\right)\right)}{N}} \right), \tag{64}$$

*where* $N_L = N_{conv} + N_{conn}$, $N_{par}$ *is total number of parameters in the neural network, and* $\tilde{C}$ *is a universal constant.*

**Proof.** *From [49, Lem.3.4], we know that:*

$$\max_{z \in \mathcal{Z}_{\mathcal{D}}} |\boldsymbol{s}_{\boldsymbol{P}}(\boldsymbol{z}) - \tilde{\boldsymbol{s}}_{\tilde{\boldsymbol{P}}}(\boldsymbol{z})| \le \sup_{x \in \mathcal{X}} ||\boldsymbol{s}_{\boldsymbol{P}}(\boldsymbol{x}) - \tilde{\boldsymbol{s}}_{\tilde{\boldsymbol{P}}}(\boldsymbol{x})|| \le R_{\mathcal{X}} \cdot \beta \cdot \left( 1 + \nu + \frac{\beta}{N_L} \right)^{N_L} \cdot d_{NN}(\boldsymbol{P}, \tilde{\boldsymbol{P}}).$$

*Furthermore, according to [49, Lem.A.8], we have:*

$$\log(\mathfrak{C}(\epsilon, \mathcal{F}_{\beta,\nu}, d_{\infty,\mathcal{S}}) \le N_{par} \cdot \log\left( \frac{3C_L}{\epsilon} \right),$$

*where* $C_L = R_{\mathcal{X}} \cdot \beta \exp(\nu N_L + \beta)$. *Following Thm.5-b), we have:*

$$\hat{\mathfrak{R}}_{\text{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \mathcal{H}) \le C \cdot \inf_{R_s \ge \alpha \ge 0}\left( N_C \cdot (N_C - 1) \cdot \phi_\ell \cdot \alpha + \phi_\ell \cdot (N_C - 1) \cdot \int_\alpha^{R_s} \sqrt{\frac{N_{par}\left(\nu N_L + \beta + \log(3\beta R_{\mathcal{X}}/\epsilon)\right)}{\tilde{N}}} d\epsilon \right)$$

*Since* $R_s > 1/\sqrt{N}$, *we can choose* $\alpha = \sqrt{\frac{1}{N}}$, *the result follows from the inequality that:*

$$\int_\alpha^{R_s} \sqrt{\frac{N_{par}\left(\nu N_L + \beta + \log(3\beta R_{\mathcal{X}}/\epsilon)\right)}{\tilde{N}}} d\epsilon \le R_s \cdot \xi(\boldsymbol{Y}) \cdot \sqrt{\frac{N_{par}\left(\nu N_L + \beta + \log\left(3\beta R_{\mathcal{X}} N\right)\right)}{N}}.$$

*Together with the fact that* $R_s > 1/\left( \frac{\xi(\boldsymbol{Y})}{N_C} \cdot \sqrt{N_{par}\left(\nu N_L + \beta + \log(3R_{\mathcal{X}} \cdot \beta \cdot N)\right)} \right)$, *we have :*

$$\hat{\mathfrak{R}}_{\text{MAUC}^\downarrow, \mathcal{S}}(\ell \circ \text{soft} \circ \mathcal{F}_{\beta,\nu}) \le 2C\left( \phi_\ell \cdot (N_C - 1) \cdot R_s \cdot \xi(\boldsymbol{Y}) \cdot \sqrt{\frac{N_{par}\left(\nu N_L + \beta + \log\left(3\beta R_{\mathcal{X}} N\right)\right)}{N}} \right). \tag{65}$$

*The proof is completed by choosing* $\tilde{C} = 2C$. □

## G EFFICIENT COMPUTATION

In this section, we will propose acceleration algorithms for loss and gradient evaluation for three well-known losses: exponential loss $\ell_{exp}(\alpha, t) = \exp(-\alpha t)$, squared loss $\ell_{sq}(t) = (\alpha - t)^2$, and hinge loss $\ell_{hinge}(\alpha, t) = \max(\alpha - t, 0)$.

Generally, the empirical surrogate risk functions $\hat{R}_\ell$ have the following abstract form:

$$\hat{R}_\ell = \sum_{i=1}^{N_C} \sum_{j \ne i} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)),$$

$$\text{where} \quad f^{(i)}(\boldsymbol{x}) = g_i(\boldsymbol{W} h_{\boldsymbol{\theta}}(\boldsymbol{x})), \boldsymbol{W} = [\boldsymbol{w}^{(1)}, \cdots, \boldsymbol{w}^{(N_C)}]^\top.$$

We adopt this general form since it covers a lot of popular models. For example, when $g_i(\cdot)$ is defined as the activation function of the last layer of a neural network (say a softmax function), $\boldsymbol{w}^{(i)}$ are the weights of the last layer, and $h_{\boldsymbol{\theta}}(\cdot)$ is the neural net where the last layer is excluded, $f^{(i)}(\boldsymbol{x})$ becomes a deep neural network architecture with the last layer designed as a fully-connected layer. As for another instance, if both $g_i(\boldsymbol{x}) = \boldsymbol{x}$ and $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{x}$, then we reach a simple linear multiclass scoring function. Note that the scalability with respect to sample size only depends on the choice of $\ell$. The choice of $g_i(\cdot)$

and $h_{\boldsymbol{\theta}}(\cdot)$ only affects the instance-wise chain rule. This allows us to provide a general acceleration framework once the surrogate loss is fixed.

**Common Notations**. Now we introduce the common notations we adopted in this section. We assume that $\boldsymbol{w}^{(i)}$ be a vector such that $\boldsymbol{w}^{(i)} \in \mathbb{R}^{n_h \times 1}$. We let $\boldsymbol{W}$ be the matrix that concatenates $\boldsymbol{w}^{(i)}$ column-wisely, *i.e.* ,

$$\boldsymbol{W} = [\boldsymbol{w}^{(1)}, \cdots, \boldsymbol{w}^{(N_C)}]^\top, \quad \boldsymbol{W} \in \mathbb{R}^{N_C \times n_h}. \tag{66}$$

For the sake of convenience, we rearrange $\boldsymbol{\theta}$ as a vector in $\mathbb{R}^{d_\theta \times 1}$ such that it concatenates all the parameters in the model except $\boldsymbol{W}$. For example, in a deep neural network, $\boldsymbol{\theta}$ concatenates all the parameters except the weights of the last layer. Note that this does not affect the result of the loss and gradient evaluation. A simple reshape suffices to transfer our calculation of gradient to its expected result if $\boldsymbol{\theta}$ has complicated structures. $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ should be considered as a vector such that $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \in \mathbb{R}^{n_h \times 1}$. $\boldsymbol{H}$ is then defined as the matrix that concatenates all the $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ from every instance in the sense that:

$$\boldsymbol{H} = [h_{\boldsymbol{\theta}}(\boldsymbol{x}_1), \cdots, h_{\boldsymbol{\theta}}(\boldsymbol{x}_N)]^\top, \quad \boldsymbol{H} \in \mathbb{R}^{n_h \times N}. \tag{67}$$

For intermediate partial derivatives, we denote:

$$\partial_i f^{(j)}(\boldsymbol{x}) = \frac{\partial(f^{(j)}(\boldsymbol{x}))}{\partial \boldsymbol{w}^{(i)\top} h_{\boldsymbol{\theta}}(\boldsymbol{x})} \tag{68}$$

Furthermore, we have two compact expressions as:

$$\begin{aligned} \boldsymbol{\partial}_{i,j} &= [\partial_i f^{(j)}(\boldsymbol{x}_1), \cdots, \partial_i f^{(j)}(\boldsymbol{x}_N)]^\top, \quad \boldsymbol{\partial}_{i,j} \in \mathbb{R}^{N \times 1}, \\ \boldsymbol{\partial}_j(\boldsymbol{x}_k) &= [\partial_1[f^{(j)}(\boldsymbol{x}_k)], \cdots, \partial_{N_C}(f^{(j)}(\boldsymbol{x}_k))]^\top, \quad \boldsymbol{\partial}_j(\boldsymbol{x}_k) \in \mathbb{R}^{N_C \times 1}. \end{aligned} \tag{69}$$

Another variable is also essential when applying the chain rule:

$$\boldsymbol{U}^{(j)} = [\nabla_{\boldsymbol{\theta}} h_{\boldsymbol{\theta}}(\boldsymbol{x}_1) \boldsymbol{W}^\top \boldsymbol{\partial}_j(\boldsymbol{x}_1), \cdots, \nabla_{\boldsymbol{\theta}} h_{\boldsymbol{\theta}}(\boldsymbol{x}_N) \boldsymbol{W}^\top \boldsymbol{\partial}_j(\boldsymbol{x}_N)], \quad \boldsymbol{U}^{(j)} \in \mathbb{R}^{d_\theta \times N}. \tag{70}$$

Finally, we provide a compact notation of the weights $\frac{1}{n_i n_j}$ as:

$$\boldsymbol{D}^{(i)} = [\frac{1}{n_i n_{y_1}}, \cdots, \frac{1}{n_i n_{y_N}}]^\top, \quad \boldsymbol{D}^{(i)} \in \mathbb{R}^{N \times 1}. \tag{71}$$

## G.1 Exponential Loss

**Loss Evaluation**. For the exponential loss, we can simplify the computations by a factorization scheme:

$$\begin{aligned} \hat{R}_{exp} &= \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \exp\left(\alpha \cdot \left(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)\right)\right), \\ &= \sum_{i=1}^{N_C} \underbrace{\left(\sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \exp(\alpha \cdot f^{(i)}(\boldsymbol{x}_m))\right)}_{(a_i)} \cdot \underbrace{\left(\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \exp(-\alpha \cdot f^{(i)}(\boldsymbol{x}_n))\right)}_{(b_i)} \end{aligned}$$

From the derivation above, the loss evaluation could be done by first calculating $(a_i), (b_i)$ separately and then performing the multiplication, which only takes $O(N N_C T_\ell)$. This is a significant improvement compared with the original $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j N_C T_\ell)$ result.

**Gradient Evaluation**. Due to the factorization scheme, we can also accelerate the gradient evaluation in a similar spirit. According to. Eq.(66)-Eq.(70), we have:

$$\begin{aligned} \nabla_{\boldsymbol{w}^{(i)}}(a_j) &= \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \alpha \exp(\alpha \cdot f^{(j)}(\boldsymbol{x}_m)) \cdot \partial_i f^{(j)}(\boldsymbol{x}_m) \cdot h_{\boldsymbol{\theta}}(\boldsymbol{x}_m) \\ &= \alpha \boldsymbol{H}(\mathcal{I}_i \odot \hat{\boldsymbol{y}}_{exp}^{(i),+} \odot \boldsymbol{\partial}_{i,j}). \\ \nabla_{\boldsymbol{w}^{(i)}}(b_j) &= -\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \alpha \cdot \exp(-\alpha \cdot f^{(j)}(\boldsymbol{x}_n)) \cdot \partial_i f^{(j)}(\boldsymbol{x}_n) \cdot h_{\boldsymbol{\theta}}(\boldsymbol{x}_n) \\ &= -\alpha \boldsymbol{H}(\mathcal{I}_{\backslash i} \odot \boldsymbol{D}^{(i)} \odot \hat{\boldsymbol{y}}_{exp}^{(i),-} \odot \boldsymbol{\partial}_{i,j}), \\ \nabla_{h_{\boldsymbol{\theta}}(\boldsymbol{x}_k)}(a_j) &= \underbrace{\mathcal{I}_j(\boldsymbol{x}_k) \cdot \alpha \cdot \exp(\alpha \cdot f^{(j)}(\boldsymbol{x}_k))}_{q_j(\boldsymbol{x}_k)} \cdot \sum_{i=1}^{N_C} \partial_i f^{(j)}(\boldsymbol{x}_k) \cdot \boldsymbol{w}^{(i)}, \\ &= q_j(\boldsymbol{x}_k) \cdot \boldsymbol{W}^\top \boldsymbol{\partial}_j(\boldsymbol{x}_k) \\ \nabla_{h_{\boldsymbol{\theta}}(\boldsymbol{x}_k)}(b_j) &= -\underbrace{\mathcal{I}_j^c(\boldsymbol{x}_k) \cdot \boldsymbol{D}_k^{(i)} \cdot \alpha \cdot \exp(-\alpha \cdot f^{(j)}(\boldsymbol{x}_k))}_{\tilde{q}_j(\boldsymbol{x}_k)} \cdot \sum_{i=1}^{N_C} \partial_i f^{(j)}(\boldsymbol{x}_k) \cdot \boldsymbol{w}^{(i)}, \\ &= -\tilde{q}_j(\boldsymbol{x}_k) \cdot \boldsymbol{W}^\top \boldsymbol{\partial}_j(\boldsymbol{x}_k) \end{aligned}$$

where

$$\mathcal{I}_i = I\left[\boldsymbol{x}_i \in \mathcal{N}_i\right], \quad \mathcal{I}_{\setminus i} = I\left[\boldsymbol{x}_i \notin \mathcal{N}_i\right], \quad \mathcal{I}_j(\boldsymbol{x}) = \boldsymbol{I}\left[x \in \mathcal{N}_j\right], \quad \mathcal{I}_j^c(\boldsymbol{x}) = 1 - \mathcal{I}_j(\boldsymbol{x})$$
$$\hat{\boldsymbol{y}}_{exp}^{(i),+} = [\exp(\alpha \cdot f^{(i)}(\boldsymbol{x}_1)), \cdots, \exp(\alpha \cdot f^{(i)}(\boldsymbol{x}_N))]^\top,$$
$$\hat{\boldsymbol{y}}_{exp}^{(i),-} = [\exp(-\alpha \cdot f^{(i)}(\boldsymbol{x}_1)), \cdots, \exp(-\alpha \cdot f^{(i)}(\boldsymbol{x}_N))]^\top.$$

Furthermore, we define $\boldsymbol{q}_j = [q_j(\boldsymbol{x}_1), \cdots, q_j(\boldsymbol{x}_N)]^\top$ and $\tilde{\boldsymbol{q}}_j = [\tilde{q}_j(\boldsymbol{x}_1), \cdots, \tilde{q}_j(\boldsymbol{x}_N)]^\top$. From Eq.(70), we have:

$$\nabla_{\boldsymbol{\theta}}(a_j) = \sum_{k=1}^{N} \nabla_{\boldsymbol{\theta}}(h_{\boldsymbol{\theta}}(\boldsymbol{x}_k))\nabla_{h_{\boldsymbol{\theta}}(\boldsymbol{x}_k)}a_j = \boldsymbol{U}^{(j)}\boldsymbol{q}_j$$

$$\nabla_{\boldsymbol{\theta}}(b_j) = \sum_{k=1}^{N} \nabla_{\boldsymbol{\theta}}(h_{\boldsymbol{\theta}}(\boldsymbol{x}_k))\nabla_{h_{\boldsymbol{\theta}}(\boldsymbol{x}_k)}b_j = -\boldsymbol{U}^{(j)}\tilde{\boldsymbol{q}}_j$$

Then the gradient evaluation can be done by :

$$\nabla_{\boldsymbol{w}^{(i)}}\hat{R}_\ell = \sum_{j=1}^{N_C} \left(\nabla_{\boldsymbol{w}^{(i)}}(a_j)\right) \cdot (b_j) + (a_j) \cdot \left(\nabla_{\boldsymbol{w}^{(i)}}(b_j)\right)$$

$$= \alpha \boldsymbol{H}\left(\sum_{j=1}^{N_C} b_j \mathcal{I}_i \odot \hat{\boldsymbol{y}}_{exp}^{(i),+} \odot \boldsymbol{\partial}_{i,j} - a_j \mathcal{I}_{\setminus i} \odot \boldsymbol{D}^{(i)} \odot \hat{\boldsymbol{y}}_{exp}^{(i),-} \odot \boldsymbol{\partial}_{i,j}\right),$$

$$\nabla_{\boldsymbol{\theta}}\hat{R}_\ell = \sum_{j=1}^{N_C} \left(\nabla_{\boldsymbol{\theta}}(a_j)\right) \cdot (b_j) + (a_j) \cdot \left(\nabla_{\boldsymbol{\theta}}(b_j)\right)$$

$$= \sum_{j=1}^{N_C} \left[\boldsymbol{U}^{(j)}(b_j\boldsymbol{q}_j - a_j\tilde{\boldsymbol{q}}_j)\right].$$

From the compact forms shown above, the acceleration enjoys a $O(NN_CT_{grad})$ complexity instead of the original $O(N_C \sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j T_{grad})$ time.

## G.2 Hinge loss

**Loss Evaluation.** First we put down the hinge surrogate loss as:

$$\hat{R}_{hinge} = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \left(\alpha - \left(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)\right)\right)_+.$$

The key of our acceleration is to notice that the terms are non-zero only if $\left(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)\right) \leq \alpha$. Moreover, for these non-zero terms $\max(x, 0) = x$. This means that the hinge loss degenerates to an identity function for the activated nonzero terms, which enjoys efficient computation. So the key step in our algorithm is to find out the non-zero terms in an efficient manner.

Given a fixed class $i$ and an instance $\boldsymbol{x}_m \in \mathcal{N}_i$, we denote:

$$\mathcal{A}^{(i)}(\boldsymbol{x}_m) = \{\boldsymbol{x}_n \notin \mathcal{N}_i, \alpha > f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)\}.$$

With $\mathcal{A}^{(i)}(\boldsymbol{x}_m)$, one can reformulate $\hat{R}_{hinge}$ as:

$$\hat{R}_{hinge} = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \left[\left(\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j}\right) \cdot \left(\alpha - f^{(i)}(\boldsymbol{x}_m)\right) + \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} \cdot f^{(i)}(\boldsymbol{x}_n)\right],$$

$$= \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \left[\delta^{(i)}(\boldsymbol{x}_m) \cdot (\alpha - f^{(i)}(\boldsymbol{x}_m)) + \Delta^{(i)}(\boldsymbol{x}_m)\right],$$

where

$$\delta^{(i)}(\boldsymbol{x}_m) = \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j}, \quad \Delta^{(i)}(\boldsymbol{x}_m) = \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} f^{(i)}(\boldsymbol{x}_n).$$

According to this reformulation, once $\delta^{(i)}(\boldsymbol{x})$ and $\Delta^{(i)}(\boldsymbol{x})$ are all fixed, we can calculate the loss function within $O(N)$. So the efficiency bottleneck comes from the calculations of $\delta^{(i)}(\boldsymbol{x})$ and $\Delta^{(i)}(\boldsymbol{x})$.

---

**Algorithm 1** `CalIndex`

---

**Input:** $\boldsymbol{x}_0^\downarrow, \cdots, \boldsymbol{x}_{n_i-1}^\downarrow, \tilde{\boldsymbol{x}}_0^\downarrow, \cdots, \tilde{\boldsymbol{x}}_{N-n_i-1}^\downarrow, \alpha.$

1: $p \leftarrow 0$                      ▷ pointer for current positive instance

2: $q \leftarrow 0$                      ▷ pointer for current negative instance

3: $\mathcal{I}_{0:(n_i-1)} = 0.$             ▷ The index set that satisfies $\mathcal{A}_j^{(i)} = \{\tilde{\boldsymbol{x}}_k^{(i)\downarrow}\}_{k=1}^{\mathcal{I}_k}$

4: **while** $p < n_i, q < N - n_i$ **do**

5:   **if** $f^{(i)}(\boldsymbol{x}_p^{(i)\downarrow}) - f^{(i)}(\tilde{\boldsymbol{x}}_q^{(i)\downarrow}) < \alpha$ **then**

6:    $q+=1$          ▷ If the current pair is nonzero, check the next negative instance

7:   **else**

8:    $\mathcal{I}_p = \max(q-1, 0)$         ▷ The last element of $\mathcal{A}_p^{(i)}$ is $\tilde{\boldsymbol{x}}_{q-1}^{(i)\downarrow}$

9:    $p+=1$

10:   **end if**

11: **end while**

12: **if** $q = N - n_i$ **then**

13:   $\mathcal{I}_p = \max(q-1, 0)$         ▷ The last element of $\mathcal{A}_p^{(i)}$ is $\tilde{\boldsymbol{x}}_{N-n_i-1}^{(i)\downarrow}$

14: **end if**

15: **if** $p < n_i$ **then**      ▷ If $\mathcal{A}_p^{(i)}$ includes all the negative instances then for all $\mathcal{A}_{p'}^{(i)} = \mathcal{A}_p^{(i)}, p' > p$

16:   $\mathcal{I}_{(p+1):(n_i-1)} = \mathcal{I}_p$

17: **end if**

---

Next we propose an efficient way to calculate $\delta^{(i)}(\boldsymbol{x}_m), \Delta^{(i)}(\boldsymbol{x}_m), \mathcal{A}^{(i)}(\boldsymbol{x}_m)$. To do so, given a specific class $i$, we first sort the relevant (resp. irrelevant) instances descendingly according to their scores $f^{(i)}(\boldsymbol{x})$:

$$f^{(i)}(\boldsymbol{x}_0^\downarrow) \geq f^{(i)}(\boldsymbol{x}_1^\downarrow) \cdots, \geq f^{(i)}(\boldsymbol{x}_{n_i-1}^\downarrow), \quad \boldsymbol{x}_i^\downarrow \in \mathcal{N}_i,$$
$$f^{(i)}(\tilde{\boldsymbol{x}}_0^\downarrow) \geq f^{(i)}(\tilde{\boldsymbol{x}}_1^\downarrow) \cdots, \geq f^{(i)}(\tilde{\boldsymbol{x}}_{N-n_i-1}^\downarrow), \quad \tilde{\boldsymbol{x}}_i^\downarrow \notin \mathcal{N}_i.$$

It immediately follows that:

$$\mathcal{A}^{(i)}(\boldsymbol{x}_0^\downarrow) \subseteq \mathcal{A}^{(i)}(\boldsymbol{x}_1^\downarrow) \subseteq \cdots \subseteq \mathcal{A}^{(i)}(\boldsymbol{x}_{n_i-1}^\downarrow).$$

This allows us to find out all $\mathcal{A}^{(i)}(\boldsymbol{x}_k^\downarrow), k = 0, 1, \cdots, n_i - 1$, within a single pass of the whole dataset with an efficient dynamic programming.

Based on the construction of $\mathcal{A}_k^{(i)} = \mathcal{A}^{(i)}(\boldsymbol{x}_k^\downarrow)$, we provide the following recursive rules to obtain $\delta^{(i)}(\boldsymbol{x})$ and $\Delta^{(i)}(\boldsymbol{x})$:

$$\delta^{(i)}(\boldsymbol{x}_{k+1}^\downarrow) = \delta^{(i)}(\boldsymbol{x}_k^\downarrow) + \sum_{j \neq i} \sum_{\mathcal{N}_j \cap \mathcal{A}_{k+1}^{(i)} \setminus \mathcal{A}_k^{(i)}} \frac{1}{n_i n_j},$$

$$\Delta^{(i)}(\boldsymbol{x}_{k+1}^\downarrow) = \Delta^{(i)}(\boldsymbol{x}_k^\downarrow) + \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}_{k+1}^{(i)} \setminus \mathcal{A}_k^{(i)}} \frac{1}{n_i n_j} f^{(i)}(\boldsymbol{x}_n).$$

These recursive rules induce an efficient dynamic programming which only requires a single pass of the whole dataset. Putting all together, we then come to a $O(N_C \bar{N} T_l)$ time complexity speed-up for hinge loss evaluation. An implementation of this idea is detailed in Alg.1 and Alg.2.

**Gradient Evaluation**. Similar to the loss evaluation, we can provide an efficient acceleration method for the gradient evaluation as well. Based on a similar analysis, we have:

$$\nabla_{\boldsymbol{\Theta}} \hat{R}_{hinge} = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \left[ -\left( \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} \right) \cdot \nabla_{\boldsymbol{\Theta}} f^{(i)}(\boldsymbol{x}_m) + \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} \nabla_{\boldsymbol{\Theta}} f^{(i)}(\boldsymbol{x}_n) \right]$$

$$= \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \left( -\delta^{(i)}(\boldsymbol{x}_m) \cdot \nabla_{\boldsymbol{\Theta}} f^{(i)}(\boldsymbol{x}_m) + \Gamma^{(i)}(\boldsymbol{x}_m) \right)$$

where

$$\Gamma^{(i)}(\boldsymbol{x}_m) = \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}^{(i)}(\boldsymbol{x}_m)} \frac{1}{n_i n_j} \nabla_{\boldsymbol{\Theta}} f^{(i)}(\boldsymbol{x}_n).$$

For any $\boldsymbol{x} \in \mathcal{S}$, and any class $j$, based on Eq.(66)-Eq.(70), we can calculate $\nabla_{\boldsymbol{\Theta}} f^{(j)}(\boldsymbol{x})$ from:

$$\nabla_{\boldsymbol{w}^{(i)}} f^{(j)}(\boldsymbol{x}) = \partial_i f^{(j)}(\boldsymbol{x}) \cdot h_{\boldsymbol{\theta}}(\boldsymbol{x}),$$
$$\nabla_{h_{\boldsymbol{\theta}}(\boldsymbol{x})} f^{(i)}(\boldsymbol{x}) = \boldsymbol{W}^\top \boldsymbol{\partial}_j(\boldsymbol{x}),$$
$$\nabla_{\boldsymbol{\theta}} f^{(i)}(\boldsymbol{x}) = \nabla_{\boldsymbol{\theta}} h_{\boldsymbol{\theta}}(\boldsymbol{x}) \boldsymbol{W}^\top \boldsymbol{\partial}_j(\boldsymbol{x}).$$

---

**Algorithm 2** CalHingeLoss

---

**Input:** $\boldsymbol{X}, \boldsymbol{Y}, \alpha$.
**Output:** loss.

1: loss $\leftarrow 0$
2: **for** $i \leftarrow 1$ **to** $N_C$ **do**
3:
4:      Obtain $\boldsymbol{x}_0^{(i)\downarrow}, \cdots, \boldsymbol{x}_{n_i-1}^{(i)\downarrow}$ based on $f^{(i)}(\boldsymbol{x})$
5:      Obtain $\boldsymbol{x}_0^{(i)\downarrow}, \cdots, \boldsymbol{x}_{N-n_i-1}^{(i)\downarrow}$ based on $f^{(i)}(\boldsymbol{x})$
6:
7:      $f^{(i)}(\boldsymbol{x})_+ = [f^{(i)}(\boldsymbol{x}_0^{(i)\downarrow}), \cdots, f^{(i)}(\boldsymbol{x}_{n_i-1}^{(i)\downarrow})]$ .
8:      $\boldsymbol{D}^{(i)} = \left[ \dfrac{1}{n_i n_{y(\tilde{\boldsymbol{x}}_0^{(i)\downarrow})}}, \cdots, \dfrac{1}{n_i n_{y(\tilde{\boldsymbol{x}}_{N-1}^{(i)\downarrow})}} \right]$
9:
10:      $t_0 \leftarrow 0, \quad \delta_{\text{offset}}^{(i)} \leftarrow 0, \quad \Delta_{\text{offset}}^{(i)} \leftarrow 0, \quad loc = 0$
11:
12:      $\boldsymbol{\Delta}_{0:(n_i-1)}^{(i)} \leftarrow 0, \quad \boldsymbol{\delta}_{0:(n_i-1)}^{(i)} \leftarrow 0$
13:
14:      $\mathcal{I} \leftarrow$ CalIndex$(\boldsymbol{x}_0^{(i)\downarrow}, \cdots, \boldsymbol{x}_{n_i}^{(i)\downarrow}, \tilde{\boldsymbol{x}}_0^{(i)\downarrow}, \cdots, \tilde{\boldsymbol{x}}_{N-n_i-1}^{(i)\downarrow}, \alpha)$
15:
16:      **while** $k < n_i$ **do**
17:          **if** $\mathcal{I}_k \neq t_0 - 1$ **then**          $\triangleright$ continue scanning when $\mathcal{A}^{(k)} \neq \mathcal{A}^{(k-1)}$
18:              $t_1 \leftarrow \mathcal{I}_k + 1$          $\triangleright$ make sure that $\mathcal{A}_{k+1}^{(i)} \backslash \mathcal{A}_k^{(i)} = \{\tilde{\boldsymbol{x}}_u^{(i)\downarrow}\}_{u=t_0}^{t_1-1}$
19:              $\delta_k^{(i)} \leftarrow \delta_{\text{offset}}^{(i)} + \sum_{u=t_0}^{t_1-1} \boldsymbol{D}_u^{(i)}$          $\triangleright$ Recursion for $\delta^{(i)}(\boldsymbol{x}_k^{(i)\downarrow})$
20:              $\Delta_k^{(i)} \leftarrow \Delta_{\text{offset}}^{(i)} + \sum_{u=t_0}^{t_1-1} \boldsymbol{D}_u^{(i)} \cdot f^{(i)}(\tilde{\boldsymbol{x}}_k^{(i)\downarrow})$          $\triangleright$ Recursion for $\Delta^{(i)}(\boldsymbol{x}_k^{(i)\downarrow})$
21:              $t_0 \leftarrow t_1$          $\triangleright$ Update the initial element.
22:              $\delta_{\text{offset}}^{(i)} \leftarrow \delta_k^{(i)}, \quad \Delta_{\text{offset}}^{(i)} \leftarrow \Delta_k^{(i)}$          $\triangleright$ Update the offsets.
23:          **else if** $I_k \neq N - n_i - 1$ **then**          $\triangleright$ Check if $\mathcal{A}^{(k)} = \mathcal{A}^{(k-1)}$ and $\mathcal{A}^{(k)} \neq \mathcal{N}_{\backslash i}$
24:              $\delta_k^{(i)} = \delta_{\text{offset}}^{(i)}, \quad \Delta_k^{(i)} = \Delta_{\text{offset}}^{(i)}$          $\triangleright$ Copy the last updates
25:          **else**
26:              break          $\triangleright$ Stop the iteration when $\mathcal{A}^{(k)} = \mathcal{N}_{\backslash i}$ and $\mathcal{A}^{(k)} = \mathcal{A}^{(k-1)}$.
27:          **end if**
28:          $k + = 1$
29:      **end while**
30:      **if** $k < n_i$ **then**
31:          $\boldsymbol{\delta}_{k:(n_i-1)}^{(i)} = \delta_{\text{offset}}^{(i)}, \quad \boldsymbol{\Delta}_{k:(n_i-1)}^{(i)} = \Delta_{\text{offset}}^{(i)}$          $\triangleright$ Set $\mathcal{A}_x^{(i)} = \mathcal{A}_{k-1}^{(i)}, \ \forall x \geq k$.
32:      **end if**
33:
34:      loss$+ = (\alpha - f^{(i)}(\boldsymbol{x})_+)^\top \boldsymbol{\delta}^{(i)} + \boldsymbol{\Delta}^{(i)} \mathbf{1}$
35:
36: **end for**

---

We can also adopt a dynamic programming scheme similar to the loss evaluation to speed it up. The first step is exactly the same: obtain $\boldsymbol{x}_0^\downarrow, \cdots, \boldsymbol{x}_{n_i-1}^\downarrow$ and $\tilde{\boldsymbol{x}}_0^\downarrow, \cdots, \tilde{\boldsymbol{x}}_{N-n_i-1}^\downarrow$ with the same method as the loss evaluation and get $\mathcal{A}_k^{(i)} = \mathcal{A}^{(i)}(\boldsymbol{x}_k^\downarrow)$. Then, we have the following recursive expression for the new variable $\Gamma^{(i)}(\boldsymbol{x}_k^\downarrow)$:

$$\Gamma^{(i)}(\boldsymbol{x}_{k+1}^\downarrow) = \Gamma^{(i)}(\boldsymbol{x}_k^\downarrow) + \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j \cap \mathcal{A}_{k+1}^{(i)} \backslash \mathcal{A}_k^{(i)}} \frac{1}{n_i n_j} \nabla_{\boldsymbol{\Theta}} f^{(i)}(\boldsymbol{x}_n).$$

Now we can use a similar algorithm to evaluate the gradients, the details are omitted here due to its resemblance to Alg.2.

## G.3 Squared Loss

**Loss Evaluation**. For each fixed class $i$, we construct an affinity matrix $\boldsymbol{Aff}^{(i)}$ such that

$$\boldsymbol{Aff}_{m,n}^{(i)} = \begin{cases} \dfrac{1}{n_i n_{y_m}}, & n \in \mathcal{N}_i, m \notin \mathcal{N}_i, \\ \dfrac{1}{n_i n_{y_n}}, & m \in \mathcal{N}_i, n \notin \mathcal{N}_i, \\ 0, & \text{Otherwise,} \end{cases}$$

---

**Algorithm 3** `CalSqLoss`

---

**Input:** $\boldsymbol{X}, \boldsymbol{Y}, \alpha, \varphi(\cdot)$.
**Output:** loss.

1: loss $\leftarrow 0$
2: **for** $i \leftarrow 1$ **to** $N_C$ **do**
3:      Obtain $\boldsymbol{D}^{(i)}$
4:      $\boldsymbol{\Delta}^{(i)} \leftarrow (\boldsymbol{Y}^{(i)} - f^{(i)}(\boldsymbol{X}))$
5:      $\boldsymbol{\kappa}^{(i)} \leftarrow n_i \boldsymbol{D}^{(i)}(\mathbf{1} - \boldsymbol{Y}^{(i)}) + \dfrac{N_C - 1}{n_i}\boldsymbol{Y}^{(i)}$
6:      $\Delta_1^{(i)} \leftarrow \left(\boldsymbol{\Delta}_{sq}^{(i)\top}(\boldsymbol{D}^{(i)}\mathbf{1} - \boldsymbol{Y}^{(i)})\right)$
7:      $\Delta_2^{(i)} \leftarrow \boldsymbol{Y}^{(i)\top}\boldsymbol{\Delta}_{sq}^{(i)}$
8:      loss$+ = \frac{1}{2}\boldsymbol{\Delta}_{sq}^{(i)\top}(\boldsymbol{\kappa}^{(i)} \odot \boldsymbol{\Delta}_{sq}^{(i)}) - \Delta_1^{(i)} \cdot \Delta_2^{(i)}$
9: **end for**

---

which could be written in a matrix form,

$$\boldsymbol{Aff}^{(i)} = \boldsymbol{D}^{(i)}(\mathbf{1} - \mathbf{Y}^{(i)})\mathbf{Y}^{(i)\top} + \mathbf{Y}^{(i)}(\mathbf{1} - \mathbf{Y}^{(i)})^\top \boldsymbol{D}^{(i)},$$

Then the empirical risk function under the squared surrogate loss could be reformulated as:

$$\hat{R}_{sq} = \sum_{i=1}^{N_C} \boldsymbol{\Delta}_{sq}^{(i)\top} \mathcal{L}^{(i)} \boldsymbol{\Delta}_{sq}^{(i)},$$

where

$$\boldsymbol{\Delta}_{sq}^{(i)} = \mathbf{Y}^{(i)} - f^{(i)}(\boldsymbol{X}), \quad f^{(i)}(\boldsymbol{X}) = [f^{(i)}(\boldsymbol{x}_1), \cdots, f^{(i)}(\boldsymbol{x}_N)]^\top, \tag{72}$$

and $\mathcal{L}^{(i)} = diag(\boldsymbol{Aff}^{(i)}\mathbf{1}) - \boldsymbol{Aff}^{(i)}$ is the graph Laplacian matrix associated with $\boldsymbol{Aff}^{(i)}$. Based on this reformulation, one can find that the structure of $\mathcal{L}^{(i)}$ provides an efficient factorization scheme to speed-up loss evaluation. To see this, we have:

$$\hat{R}_{sq} = \sum_{i=1}^{N_C} \frac{1}{2}\boldsymbol{\Delta}_{sq}^{(i)\top}(\boldsymbol{\kappa}^{(i)} \odot \boldsymbol{\Delta}_{sq}^{(i)}) - \Delta_1^{(i)} \cdot \Delta_2^{(i)},$$

where

$$\Delta_2^{(i)} = \mathbf{Y}^{(i)\top}\boldsymbol{\Delta}_{sq}^{(i)}$$
$$\Delta_1^{(i)} = \boldsymbol{\Delta}_{sq}^{(i)\top}\boldsymbol{D}^{(i)}(\mathbf{1} - \mathbf{Y}^{(i)})$$
$$\boldsymbol{\kappa}^{(i)} = n_i\boldsymbol{D}^{(i)}(\mathbf{1} - \mathbf{Y}^{(i)}) + \frac{N_C - 1}{n_i}\mathbf{Y}^{(i)}.$$

Obviously both terms in the equation above require only $O(N_C N T_l)$ time. We summarize this factorization scheme in Alg.3.

**Gradient Evaluation**. From Alg.3, Eq.(66)-Eq.(70) and the chain rule, we have:

$$\nabla_{\boldsymbol{w}^{(i)}}\hat{R}_{sq} = \sum_{i=1}^{N_C} \frac{1}{2} \cdot \nabla_{\boldsymbol{w}^{(i)}}\boldsymbol{\Delta}_{sq}^{(i)\top}(\boldsymbol{\kappa}^{(i)} \odot \boldsymbol{\Delta}_{sq}^{(i)}) - \nabla_{\boldsymbol{w}^{(i)}}\Delta_1^{(i)} \cdot \Delta_2^{(i)}$$

where

$$\frac{1}{2}\nabla_{\boldsymbol{w}^{(i)}}\boldsymbol{\Delta}_{sq}^{(j)\top}(\boldsymbol{\kappa}^{(j)} \odot \boldsymbol{\Delta}_{sq}^{(j)}) = \boldsymbol{H}(\boldsymbol{\Delta}_{sq}^{(j)} \odot \boldsymbol{\kappa}^{(j)} \odot \boldsymbol{\partial}_{i,j}),$$
$$\nabla_{\boldsymbol{w}^{(i)}}\Delta_1^{(i)}\Delta_2^{(i)} = \boldsymbol{H}\left[\left(\Delta_2^{(i)} \cdot \boldsymbol{D}^{(i)}(\mathbf{1} - \mathbf{Y}^{(i)}) + \Delta_1^{(i)}\mathbf{Y}^{(i)}\right) \odot \boldsymbol{\partial}_{i,j}\right],$$
$$\frac{1}{2}\nabla_{\boldsymbol{\theta}}\boldsymbol{\Delta}_{sq}^{(j)\top}(\boldsymbol{\kappa}^{(j)} \odot \boldsymbol{\Delta}_{sq}^{(j)}) = \boldsymbol{U}^{(j)}(\boldsymbol{\kappa}^{(j)} \odot \boldsymbol{\Delta}^{(j)}),$$
$$\nabla_{\boldsymbol{\theta}}\Delta_1^{(i)}\Delta_2^{(i)} = \boldsymbol{U}^{(j)}\left[\Delta_2^{(i)} \cdot \boldsymbol{D}^{(i)}(\mathbf{1} - \mathbf{Y}^{(i)}) + \Delta_1^{(i)}\mathbf{Y}^{(i)}\right],$$

Again this shows that the gradient evaluation could be done within $O(N_C N T_{grad})$ time.

## G.4 Acceleration Scheme for general losses

Unfortunately, the lossless acceleration for general losses is impossible in general. In this subsection, following the spirit of [75], we provide a discussion on how to accelerate general loss functions approximately. However, we can, in turn, construct an approximation framework based on the Bernstein Polynomials. Note that the key difference between our work and [75] is that we do not need the minimax reformulation. Moreover, we also present a nice property of such approximations in terms of consistency.

First of all, we define the Berstein Polynomials as the following.

**Definition 7.** *The Berstein Polynomials of degree $m$ of the function $\varphi : [0, 1] \to \mathbb{R}$ are defined, for any $u \in [0, 1]$, by*

$$\mathbb{B}_m(\varphi, u) = \sum_{i=1}^{m} \varphi\left(\frac{k}{m}\right) \cdot \binom{m}{k} \cdot u^k \cdot (1-u)^{m-k} = \sum_{i=1}^{m} \binom{m}{k} \cdot \Delta^k \varphi(0) \cdot u^k. \tag{73}$$

*where $\Delta^k \varphi(0) = \sum_{j=0}^{k} (-1)^{k-j} \binom{k}{j} \varphi\left(\frac{j}{m}\right)$ is the forward difference operator on $\varphi$ at 0.*

Note that the domain of the scoring functions $f^{(i)}$ are restricted to $[0, 1]$, by choosing $\varphi(u) = \ell(2s - 1)$, we have:

$$\begin{aligned}
\ell(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)) &\approx \mathbb{B}_K\left(\varphi, \frac{1 + f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)}{2}\right) \\
&= \sum_{k=0}^{K} \binom{m}{k} \Delta^k \varphi(0) \left(\varphi, \frac{1 + f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)}{2}\right)^k \\
&= \frac{1}{K+1} \sum_{k=0}^{m} f_k^{(i)}(\boldsymbol{x}_m) \cdot \tilde{f}_k^{(i)}(\boldsymbol{x}_n)
\end{aligned} \tag{74}$$

where

$$\begin{aligned}
f_k^{(i)}(\boldsymbol{x}_m) &= \left(\frac{1}{2} + f^{(i)}(\boldsymbol{x}_m)\right)^k \\
f_k^{(i)}(\boldsymbol{x}_n) &= \sum_{j=k}^{m} \binom{m}{j} \cdot \binom{j}{k} \cdot \frac{(m+1) \cdot \Delta^k \varphi(0)}{2^j} \cdot \left(\frac{1}{2} - f^{(i)}(\boldsymbol{x}_n)\right)^{j-k}.
\end{aligned}$$

With the Bernstein Polynomials, we now construct generic acceleration methods for surrogate risk computation.

**Acceleration for General Loss Calculation**. Through the lens of the Bernstein polynomials, the loss function could be approximated by the following derivation:

$$\begin{aligned}
\hat{R}_\ell &= \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \ell\left(f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n)\right) \\
&\approx \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \frac{1}{K+1} \sum_{k=0}^{K} f_k^{(i)}(\boldsymbol{x}_m) \cdot \tilde{f}_k^{(i)}(\boldsymbol{x}_n) \\
&= \sum_{i=1}^{N_C} \frac{1}{K+1} \boldsymbol{e}_+^{(i)\top} \boldsymbol{e}_-^{(i)} \\
&\triangleq \hat{R}_\ell^K
\end{aligned}$$

where

$$\boldsymbol{e}_+^{(i)} = \left[\sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{f_1^{(i)}(\boldsymbol{x}_m)}{n_i}, \cdots, \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \frac{f_K^{(i)}(\boldsymbol{x}_m)}{n_i}\right]^\top,$$

$$\boldsymbol{e}_-^{(i)} = \left[\sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\tilde{f}_1^{(i)}(\boldsymbol{x}_n)}{n_j}, \cdots, \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{\tilde{f}_K^{(i)}(\boldsymbol{x}_n)}{n_j}\right]^\top.$$

Obviously, it only requires $O(K \cdot N_C \cdot N)$ time to finish the loss computation.

**Acceleration for General Gradient Calculation**. Similarly, the gradient calculation could be approximated by the following derivation:

$$\nabla \hat{R}_\ell^K = \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \nabla \ell \left( f^{(i)}(\boldsymbol{x}_m) - f^{(i)}(\boldsymbol{x}_n) \right)$$

$$\approx \sum_{i=1}^{N_C} \sum_{\boldsymbol{x}_m \in \mathcal{N}_i} \sum_{j \neq i} \sum_{\boldsymbol{x}_n \in \mathcal{N}_j} \frac{1}{n_i n_j} \cdot \frac{1}{K+1} \sum_{k=0}^{K} f_k^{(i)}(\boldsymbol{x}_m) \cdot \nabla \tilde{f}_k^{(i)}(\boldsymbol{x}_n)$$

$$= \sum_{i=1}^{N_C} \frac{1}{K+1} (\nabla \boldsymbol{e}_+^{(i)})^\top \boldsymbol{e}_-^{(i)} + \boldsymbol{e}_+^{(i)\top}(\nabla \boldsymbol{e}_-^{(i)})$$

$\nabla \boldsymbol{e}_+^{(i)}$ and $\nabla \boldsymbol{e}_-^{(i)}$ could be obtained in a similar spirit to the calculation of the exponential loss. It thus requires $O(K \cdot N_C \cdot N \cdot T_{grad})$ time to finish the gradient computation.

**Consistency of the Approximated Loss function** Now, we proceed to present an elegant property of the approximated risk functions. Specifically, we show that $\hat{R}_\ell^K$ is, in itself, a consistent loss function if $\ell$ is consistent. To do this, we need the following theorem about the Bernstein polynomials.

**Theorem 10.** *If $f \in C^k[0,1]$, for some $k \geq 0$, then :*

$$m \leq \frac{d^k(f(x))}{dx^k} \leq M, \ \Rightarrow \ c_k \cdot m \leq \frac{d^k(\mathbb{B}_K(f,x))}{dx^k} \leq c_k \cdot M$$

Then the following corollary gives the final result.

**Corollary 2.** *For all surrogate loss function $\ell \in C^2[0,1]$,*

$$\hat{R}_\ell \text{ is MAUC}^\downarrow \text{ consistent } \Rightarrow \hat{R}_\ell^K \text{ is MAUC}^\downarrow \text{ consistent}$$

**Proof.** *Again, by choosing $\varphi(u) = \ell(2s - 1)$, we have $\varphi(u) \in C^2[0,1]$. The result follows directly from Thm.10, and Thm.3.* $\square$

## G.5  Summary and Concluding Remarks

TABLE 6
Acceleration for three losses , where $\bar{N} = \sum_{i=1}^{N_C} n_i \log n_i + (N - n_i) \log(N - n_i)$

| Algorithms | loss | gradient | requirement |
|---|---|---|---|
| exp + acceleration | $O(N_C \cdot N \cdot T_\ell)$ | $O(N_C \cdot N \cdot T_{grad})$ | $\min_i n_i \gg 2$ |
| squared + acceleration | $O(N_C \cdot N \cdot T_\ell)$ | $O(N_C \cdot N \cdot T_{grad})$ | $e^{\frac{1}{2}(N-n_i)} \gg n_i \gg N - e^{\frac{1}{2}n_i}$ |
| hinge + acceleration | $O(N_C \cdot \bar{N} \cdot T_\ell)$ | $O(N_C \cdot \bar{N} \cdot T_{grad})$ | $\min_i n_i \gg 2$ |
| general + acceleration | $O(K \cdot N_C \cdot N \cdot T_\ell)$ | $O(K \cdot N_C \cdot N \cdot T_{grad})$ | $e^{\frac{1}{2}(N-n_i)} \gg n_i \gg N - e^{\frac{1}{2}n_i}$ |
| w/o acceleration | $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j \cdot T_\ell)$ | $O(\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j \cdot T_{grad})$ | \ |

$$\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j - N_C \cdot N = \sum_{i=1}^{N_C} [n_i(N - n_i) - (n_i + (N - n_i))].$$

To ensure $N_C \cdot N < \sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j$, one must let

$$n_i(N - n_i) > (n_i + (N - n_i)), \quad \forall i.$$

Obviously, this requirement could be met when $\min_i n_i \gg 2$. Similarly, we have:

$$\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j - N_C \cdot \bar{N} = \sum_{i=1}^{N_C} [n_i(N - n_i) - n_i \log(n_i) - (N - n_i) \log(N - n_i)].$$

In this sense $N_C \cdot \bar{N} < \sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j$ could be satisfied if

$$\frac{1}{2}(N - n_i) > \log(n_i), \ \ \frac{1}{2}(n_i) > \log(N - n_i), \ \ \forall i.$$

This means that

$$\exp(\frac{1}{2}(N - n_i)) \gg n_i \gg N - \exp(\frac{1}{2}n_i), \ \ \forall i.$$

should be guaranteed to ensure a significant improvement. It is easy to see that all the requirements are essentially loose restricts on the data distribution which could be satisfied for the majority of real-world datasets. As another remark, we

do not provide a detailed calculation of $\nabla_{\boldsymbol{\theta}} h_{\boldsymbol{\theta}}(\boldsymbol{x})$, since it follows the chain rule and could be obtained from off-the-shelf toolboxes such as tensorflow and pytorch where auto differentiation is supported. Moreover, our acceleration schemes are based either on changing the order of operations or on reformulating the original loss. Hence, when the model is implemented in such auto differentiation platforms, the gradient acceleration does not need to be explicitly implemented. One just needs to implement the accelerated loss evaluation to fix the computational graph, the gradient evaluation then follows the graph naturally.

## H EXPERIMENTS

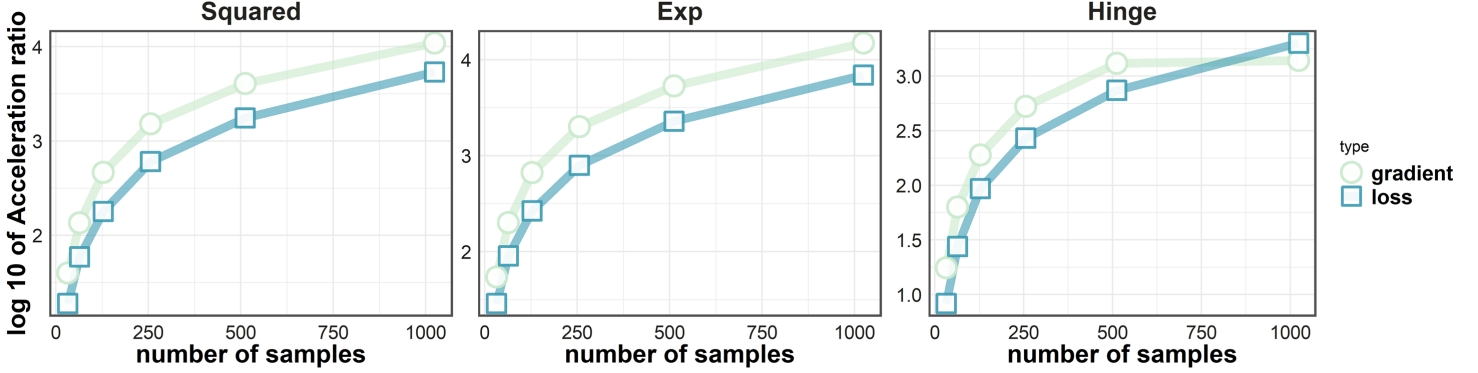### H.1 Numerical Validations for the Acceleration Algorithms



Fig. 3. **Acceleration Ratio vs. Sample Scale**

**Setting**. As a warm-up experiment, we first validate the effectiveness of our proposed speed-up methods detailed in Sec.6. Note that our algorithms only concern loss and gradient evaluations, which do not depend on the optimization algorithms in question. Hence, to be fair, we will perform stand-alone running time tests where no optimization algorithms are involved. Following this spirit, our comparisons are based on the results from 30 repetitions of pure loss and gradient evaluations. To do this, we randomly generate a series of dataset $\boldsymbol{Z}_i = (\boldsymbol{X}_i, \boldsymbol{Y}_i)$, $i = 1, 2, \cdots, 5$, where $\boldsymbol{X}_i \in \mathbb{R}^{N_k \times d}, \boldsymbol{Y}_i \in \mathbb{R}^{N_i \times N_C}$ are the input feature matrix and the output one-hot label matrix. Specifically, we set $N_i = \{32, 64, 128, 256, 512, 1024\}$, $d = 100$, $N_C = 5$. $\boldsymbol{X}_i$ is generated from the uniform distribution within $[0, 1]$. Denote $\rho_i = n_i/N$, then $\boldsymbol{Y}$ is generated such that $\rho_1 = 0.2, \rho_2 = 0.1, \rho_3 = 0.2, \rho_4 = 0.4, \rho_5 = 0.1$. Moreover, we fix a linear scoring function $f(\cdot)$ with softmax output $f^{(i)}(\boldsymbol{x}) = \text{softmax}(\boldsymbol{W}^{(i)}\boldsymbol{x})$, where the weight matrix $\boldsymbol{W} = [\boldsymbol{W}^{(1)}, \cdots, \boldsymbol{W}^{(5)}]$ is generated in the same way as $\boldsymbol{X}$. Then we test the running time of the naive evaluation algorithm and our speed-up algorithms respectively over different $\boldsymbol{Z}_i$.

**Results and Analysis**. To show the efficiency improvements, we plot the mean of the acceleration ratio in Fig.3, which is defined as:

$$\text{acceleration ratio} = \frac{\text{Running time of the naive method}}{\text{Running time of the acceleration method}}.$$

over 30 repetitions against the scale of the datasets. We have the following two findings: (a) The proposed algorithms generate significant speed-up over all sample scales being tested. They produce $10 \sim 100$x speed-up when the sample size remains as small as 32 while the result gradually climbs-up as the sample size enlarges and finally reaches 10000x speed-up for the squared and the exponential loss when the sample size becomes 1024. This shows that the proposed algorithms could both work for mini-batch-based and full-batch-based optimization methods. (b) The curves suggest that the acceleration ratios are roughly scaled as $O(N)$, which is consistent with the complexity analysis shown in Tab.6. To see this, for squared and exponential loss, we have:

$$\text{acceleration ratio} \approx \frac{\sum_{i=1}^{N_C} \sum_{j \neq i} n_i n_j}{N_C \cdot N}$$
$$= \frac{\sum_{i=1}^{N_C} \sum_{j \neq i} \rho_i \cdot \rho_j}{N_C} \cdot N$$

Note that we keep $\rho_i$ and $N_C$ as constants in our experiment, this suggests that the acceleration ratio should be $O(N)$. For hinge loss, there are extra factors from $\log n_i$ and $\log(N - n_i)$. However, these factors mostly remain at a constant level considering that the number of samples being tested is not large. This is why we also observe an almost linear scale-up of the acceleration ratio of hinge loss. Moreover, for hinge loss, we see that the increase of the acceleration ratio becomes slower for the gradient evaluations when the number of samples surpasses 500. A proper reason for this phenomenon is that the evaluation algorithms for hinge loss are partially implemented by Cython. For a large scale dataset, the contribution of

a `Cython` implementation should be more significant for loss than for gradient. This is because that the matrix operations in gradient evaluation could be accelerated by proper and compact use of `numpy` even without `Cython`, while the loss evaluations, containing mostly elementary operations, enjoy a shaper speed-up from `Cython`.

## H.2 Dataset Description

TABLE 7
Basic Information of the Datasets

| Dataset | #samples | #classes | #features | $r_\chi$ |
|---|---|---|---|---|
| Balance | 625 | 3 | 4 | 5.88 |
| Dermatology | 357 | 6 | 34 | 5.55 |
| Ecoli | 336 | 8 | 7 | 71.50 |
| New-thyroid | 215 | 3 | 5 | 5.00 |
| Pageblocks | 548 | 5 | 10 | 164.00 |
| SegmentImb | 749 | 7 | 18 | 20.31 |
| Shuttle | 2,175 | 5 | 9 | 853.00 |
| Svmguide2 | 391 | 3 | 20 | 4.17 |
| Yeast | 1,484 | 10 | 8 | 92.60 |
| CIFAR-100-Imb | 23,350 | 100 | 2,048 | 50.00 |
| User-Imb | 24,400 | 12 | 21,527 | 20.00 |
| iNaturalList2017 | 675,170 | 5,089 | 2,048 | 276.4 |

TABLE 8
$n_i$ for User-Imb

| class | sample |
|---|---|
| F23- | 800 |
| F24-26 | 400 |
| F27-28 | 400 |
| F29-32 | 800 |
| F33-42 | 1,600 |
| F43+ | 400 |
| M22- | 1,600 |
| M23-26 | 8,000 |
| M27-28 | 800 |
| M29-31 | 1,600 |
| M32-38 | 8,000 |
| M39+ | 8,000 |

TABLE 9
The label distribution for CIFAR-100-Imb Dataset, where each cell in the table presents a specific $n_i$. The row title in the left shows the range of id in the row. The class ids are numbered in the same way as the original dataset.

| number of samples in each class | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 ∼ 10 | 100 | 50 | 360 | 20 | 20 | 300 | 420 | 20 | 100 | 200 |
| 11 ∼ 20 | 100 | 100 | 600 | 50 | 100 | 50 | 50 | 100 | 100 | 100 |
| 21 ∼ 30 | 50 | 100 | 600 | 100 | 50 | 600 | 100 | 20 | 200 | 20 |
| 31 ∼ 40 | 600 | 100 | 100 | 300 | 50 | 100 | 20 | 600 | 100 | 600 |
| 41 ∼ 50 | 20 | 600 | 100 | 300 | 20 | 300 | 420 | 600 | 480 | 100 |
| 51 ∼ 60 | 20 | 20 | 50 | 100 | 20 | 420 | 600 | 50 | 20 | 100 |
| 61 ∼ 70 | 50 | 100 | 200 | 360 | 100 | 100 | 100 | 50 | 100 | 50 |
| 71 ∼ 80 | 50 | 150 | 600 | 50 | 300 | 20 | 150 | 20 | 20 | 100 |
| 81 ∼ 90 | 20 | 480 | 50 | 50 | 100 | 20 | 50 | 20 | 300 | 50 |
| 91 ∼ 100 | 50 | 360 | 100 | 150 | 20 | 100 | 100 | 20 | 100 | 50 |

Here, we provide a more detailed description for traditional datasets from LIBSVM and KEEL.

(a) **LIBSVM Datasets**. https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/

- **Abalone⋆**. The task of this dataset is to predict the age of abalone from physical measurements. The features include basic information such as the gender, length, height of a given subject. To construct a multi-class problem from this dataset. We regard each possible age as a class. On top of this, we merge the first three classes, and we also merge all the classes greater than 20.
- **Shuttle**. The shuttle dataset contains 9 features all of which are numerical. The first one being time. The last column is the class which has been coded as follows: 1 Rad Flow,,2 Fpv Close,3 Fpv Open,4 High,5 Bypass,6 Bpv Close,7 Bpv Open.
- **Svmguide-2**. [34] released this dataset when the authors are providing a guidance for the `LIBSVM` library. This dataset comes from a bacteria protein subcellular localization dataset which was proposed in [26].

(b) **KEEL Datasets**. http://www.keel.es/

- **Balance**. This data set was generated to model psychological experimental results. Each example is classified as having the balance scale tip to the right, tip to the left, or be balanced. The attributes are the left weight, the left distance, the right weight, and the right distance.
- **Dermatology**. The differential diagnosis of erythemato-squamous diseases is a real problem in dermatology. They all share the clinical features of erythema and scaling, with few differences. The diseases in this group are psoriasis, seborrheic dermatitis, lichen planus, pityriasis rosea, chronic dermatitis, and pityriasis rubra pilaris. Usually a biopsy is necessary for the diagnosis, but, unfortunately, these diseases share many histopathological features as well. Another difficulty for the differential diagnosis is that a disease may show the features of another disease at the beginning stage and may have the characteristic features at the following stages. Patients were first evaluated clinically with 12

features. Afterwards, skin samples were taken for the evaluation of 22 histopathological features. The values of the histopathological features are determined by an analysis of the samples under a microscope. In the dataset constructed for this domain, the family history feature has the value 1 if any of these diseases has been observed in the family, and 0 otherwise. The age feature simply represents the age of the patient. Every other feature (clinical and histopathological) was given a degree in the range of 0 to 3. Here, 0 indicates that the feature was not presented, 3 indicates the largest amount possible, and 1, 2 indicate the relative intermediate values.

- **Ecoli**. The task of this dataset is to predict protein localization sites, where the features are given as a series of scores evaluated by standard detection methods.
- **Glass**. A Glass type identification dataset with 9 features.
- **New Thyroid**. This is a Thyroid Disease (New Thyroid) dataset with three classes: normal hyper and hypo.
- **Page Blocks**. The task of this dataset is to classify all the blocks of the page layout of a document that has been detected by a segmentation process, where each of the observed blocks has 10 features describing its basic information.
- **Wine**. A Wine classification dataset where the model should classify a given data point to one of the types of wines based on 13 dimensions of features.

(c) **Other Datasets**:

(1) **CIFAR-100-Imb**. We sampled an imbalanced version of the CIFAR-100 [9] dataset, which originally contains 100 image classes each with 600 instances. The 100 classes are encoded as $1, \cdots 100$. For each class $i$, we randomly sampled $n_i$ instances shown in Tab.9 to gain an imbalanced version of the dataset.

(2) **User-Imb**. The original dataset is collected from TalkingData, a famous third-party mobile data platform from China, which predicts mobile users' demographic characteristics based on their app usage records. The dataset is collected for the Kaggle Competition named Talking Data Mobile User Demographics. [10] The raw features include logged events, app attributes, and device information. There are 12 target classes 'F23-', 'F24-26','F27-28','F29-32', 'F33-42', 'F43+', 'M22-', 'M23-26', 'M27-28', 'M29-31', 'M32-38', 'M39+', which describe the demographics (gender and age) of users. In our experiments, we sample an imbalanced subset of the original dataset to leverage a class-skewed dataset. The number of instances in each class after resampling is listed as Tab.8.

(3) **iNaturalist2017.** iNaturalist Challenge 2017 dataset[11] is a large-scale image classification benchmark with 675,170 images covering 5,089 different species of plants and animals. We split the dataset into the training set, validation set, and test set at a ratio of 0.7:0.15:0.15. Since directly training deep models in such a large-scale dataset is time-consuming, we instead generate 2048-d features with a ResNet-50 model pre-trained on ImageNet for each image and train models with three fully-connected layers for all methods. We utilize Adam optimizer to train the models, with an initial learning rate of $10^{-5}$. To ensure all categories are covered in a mini-batch, the batch size is set to 8196. Other hyperparameters are the same as those in the CIFAR-100-Imb dataset.

## H.3  Implementation Details

**Infrastructure**. All the experiments are carried out on a ubuntu 16.04.6 server equipped with Intel(R) Xeon(R) CPU E5-2620 v4 cpu and a TITAN RTX GPU. The codes are implemented via `python 3.6.7`, the basic dependencies are: `pytorch` (v-1.1.0), `sklearn` (v-0.21.3), `numpy` (v-1.16.2). For traditional datasets, we implement our proposed algorithms with the help of the `sklearn` and `numpy`. For hinge loss, we use Cython to accelerate the dynamic programming algorithm. For the deep learning datasets, our proposed algorithms are implemented with `pytorch`.

**Evaluation Metric**. Given a trained scoring function $f = (f^{(1)}, \cdots, f^{(N_C)})$, all the forthcoming results are evaluated with reward reversion of MAUC, where we denote it as $\mathsf{MAUC}^{\uparrow}$:

$$\mathsf{MAUC}^{\uparrow} = \frac{1}{N_C \cdot (N_C - 1)} \cdot \sum_{i=1}^{N_C} \sum_{j \neq i} \frac{\left|\left\{ (\boldsymbol{x}_1, \boldsymbol{x}_2) : \boldsymbol{x}_1 \in \mathcal{N}_i, \ \boldsymbol{x}_2 \in \mathcal{N}_j, f^{(i)}(\boldsymbol{x}_1) > f^{(i)}(\boldsymbol{x}_2) \right\}\right|}{n_i n_j}$$

in this subsection.

**Traditional Datasets**. For all the experiments, hyper-parameters are tuned based on the training and validation set while the performances are evaluated over a test set, where the training, valid, and test set accounts for $80\%$, $10\%$, $10\%$ of the original dataset, respectively. To remain the label distribution of the overall dataset, the dataset splits are generated with a stratified sampling strategy with respect to classes. The experiments are done with 15 repetitions for each involved algorithm. Since all the datasets fall in this class contains simple features, we use a linear model composited with a softmax transformation as the scoring function. More formally, given an instance $\boldsymbol{x}$, we have $f(\boldsymbol{x}) = \mathsf{softmax}(\boldsymbol{W}\boldsymbol{x})$, where $\boldsymbol{W} = [\boldsymbol{\omega}^{(1)}, \cdots, \boldsymbol{\omega}^{(N_C)}] \in \mathbb{R}^{d \times N_C}$ is the weight of the linear function. For the sampling-based competitors, the corresponding sampling algorithm is first performed on the dataset, then a scoring function is trained based on the multiclass cross-entropy loss. For LR, the scoring function is directly trained based on the multiclass cross-entropy loss without sampling. For our proposed algorithms, we directly wrap the $\mathsf{MAUC}^{\downarrow}$ surrogate losses on the scoring function $f$ and perform the training process without any simpling procedure. The hyper-parameters we adopted to produce the performance of our algorithms are shown in Tab.10.

9. https://www.cs.toronto.edu/~kriz/cifar.html
10. https://www.kaggle.com/c/talkingdata-mobile-user-demographics/data
11. https://www.kaggle.com/c/inaturalist-challenge-at-fgvc-2017

TABLE 10
Best parameters for $(\lambda, \alpha)$ over traditonal datasets

| Dataset | Ours1 | Ours2 | Ours3 |
|---|---|---|---|
| balance | $(1e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.7)$ |
| dermatology | $(1e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.3)$ |
| ecoli | $(6e\text{-}4, 0.6)$ | $(4e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.2)$ |
| new-thyroid | $(1e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.5)$ | $(1e\text{-}4, 0.9)$ |
| pageblocks | $(6e\text{-}4, 0.8)$ | $(2e\text{-}4, 0.9)$ | $(6e\text{-}4, 0.5)$ |
| segmentImb | $(2e\text{-}4, 0.9)$ | $(4e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.3)$ |
| shuttle | $(2e\text{-}4, 0.8)$ | $(1e\text{-}4, 0.7)$ | $(1e\text{-}4, 0.5)$ |
| svmguide2 | $(1e\text{-}4, 0.9)$ | $(1e\text{-}4, 0.9)$ | $(2e\text{-}4, 0.4)$ |
| yeast | $(6e\text{-}4, 0.3)$ | $(9e\text{-}3, 0.1)$ | $(1e\text{-}4, 0.2)$ |

TABLE 11
The Common Architecture for CIFAR-100-Imb
Dataset

| layers | $n_{units}$ | activation | $BN$ |
|---|---|---|---|
| $fc_1$ | $1,024$ | relu | ✓ |
| $fc_2$ | $256$ | relu | ✓ |
| $output$ | $100$ | softmax | ✗ |

TABLE 12
The Common Architecture for User-Imb Dataset

| layers | $n_{units}$ | activation | $BN$ |
|---|---|---|---|
| $fc_1$ | $128$ | relu | ✓ |
| $fc_2$ | $64$ | relu | ✓ |
| $output$ | $12$ | softmax | ✗ |

**Deep Learning Datasets**. For all the experiments, hyper-parameters are tuned based on the training and validation set while the performances are evaluated over a test set, where the training, valid, and test set accounts for $80\%$, $10\%$, $10\%$ of the original dataset, respectively. To remain the label distribution of the overall dataset, the dataset splits are generated with a stratified sampling strategy with respect to classes. For these two datasets, all the models are constructed with a deep neural network. The scoring functions share a common form as $f(\boldsymbol{x}) = \mathsf{softmax}(fc_2(fc_1(\boldsymbol{x})))$, where $fc_1, fc_2$ are fully-connected layers. For CIFAR-100-Imb, we adopt the ResNet-50 pre-trained features from `layer4` after average pooling as the input for each method. For User-Imb, the models are trained from scratch. The architecture we use for CIFAR-100-Imb and User-Imb are shown in Tab.11 and Tab.12, respectively. For the sampling-based competitors, the corresponding sampling algorithm is first performed on the dataset, then a scoring function is trained based on the multiclass cross-entropy loss together with a three-layered deep neural network. For LR, the scoring function is directly trained based on the multiclass cross-entropy loss and the neural net without sampling. For our proposed algorithms, we directly wrap the MAUC$^{\downarrow}$ surrogate losses on the neural net $f$ and perform the training process without any simpling procedure. The hyper-parameters we adopted to produce the performance of our algorithms are shown in Tab.13 and Tab.14, respectively for CIFAR-100-Imb and User-Imb.

TABLE 13
Hyperparameters for CIFAR-100-Imb

| | batch size | lr | weight decay | lr decay | epoch | gamma | #epochs before lr decay |
|---|---|---|---|---|---|---|---|
| Ours1 | 1000 | 0.0012 | 0.000005 | 0.97 | 50.00 | 1.00 | 5.00 |
| Ours2 | 1000 | 0.001 | 0.000001 | 0.99 | 48.00 | 4.00 | 5.00 |
| Ours3 | 1000 | 0.001 | 0.00005 | 0.99 | 50.00 | 4.00 | 3.00 |

TABLE 14
Hyperparameters for User-Imb

| | batch size | lr | weight decay | lr decay | gamma | #epochs before lr decay |
|---|---|---|---|---|---|---|
| Ours1 | 32 | 0.005 | .0001 | 0.97 | 0.5 | 1.00 |
| Ours2 | 32 | 0.005 | .0001 | 0.97 | 2.0 | 1.00 |
| Ours3 | 32 | 0.005 | .0001 | 0.97 | 2.0 | 1.00 |

## H.4 Running Time Comparison for Deep Learning Datasets.

To validate the effectiveness of our acceleration method, we report the average running time per epoch in Tab.15 here for all three deep learning datasets. *Accelerated* refers to our proposed method. For the *Selected* method, the sample indexes of different classes are cached in different tensors. Naive refers to a straightforward implementation.

TABLE 15
Running Time Comparison for Different Implementation Schemes, - Refers to the Fact that the Running Time is Longer than 12 h.

| loss | implementation | CIFAR-100 | User | iNaturalist |
|---|---|---|---|---|
| Square | Accelerated | 0.95 | 8.59 | 121.8347 |
| | Selected | 3111.1 | 89.35 | - |
| | Naïve | 51747.03 | 360.36 | - |
| EXP | Accelerated | 0.97 | 9.41 | 117.4387 |
| | Selected | 2392.55 | 72.31 | - |
| | Naïve | 56138.42 | 352.77 | - |
| Hinge | Accelerated | 0.83 | 8.42 | 115.4508 |
| | Selected | 2703.15 | 75.15 | - |
| | Naïve | 51302.48 | 352.45 | - |