

---

# Weighted ROC Curve in Cost Space: Extending AUC to Cost-Sensitive Learning

---

Huiyang Shao<sup>1,2</sup>      Qianqian Xu<sup>1\*</sup>      Zhiyong Yang<sup>2</sup>  
Peisong Wen<sup>1,2</sup>      Peifeng Gao<sup>2</sup>      Qingming Huang<sup>2,1,3\*</sup>

<sup>1</sup> Key Lab. of Intelligent Information Processing, Institute of Computing Tech., CAS

<sup>2</sup> School of Computer Science and Tech., University of Chinese Academy of Sciences

<sup>3</sup> BDKM, University of Chinese Academy of Sciences

shaohuiyang21@mails.ucas.ac.cn    xuqianqian@ict.ac.cn

wenpeisong20z@ict.ac.cn    gaopeifeng21@mails.ucas.ac.cn

yangzhiyong21@ucas.ac.cn    qmhuang@ucas.ac.cn

## Abstract

In this paper, we aim to tackle flexible cost requirements for long-tail datasets, where we need to construct a (1) cost-sensitive and (2) class-distribution robust learning framework. The misclassification cost and the area under the ROC curve (AUC) are popular metrics for (1) and (2), respectively. However, limited by their formulations, models trained with AUC are not well-suited for cost-sensitive decision problems, and models trained with fixed costs are sensitive to the class distribution shift. To address this issue, we present a new setting where costs are treated like a dataset to deal with arbitrarily unknown cost distributions. Moreover, we propose a novel weighted version of AUC where the cost distribution can be integrated into its calculation through decision thresholds. To formulate this setting, we propose a novel bilevel paradigm to bridge weighted AUC (WAUC) and cost. The inner-level problem approximates the optimal threshold from sampling costs, and the outer-level problem minimizes the WAUC loss over the optimal threshold distribution. To optimize this bilevel paradigm, we employ a stochastic optimization algorithm (SACCL) which enjoys the same convergence rate ( $O(\epsilon^{-4})$ ) with the SGD. Finally, experiment results show that our algorithm performs better than existing cost-sensitive learning methods and two-stage AUC decisions approach.

## 1 Introduction

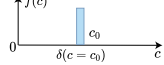
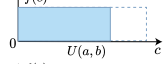

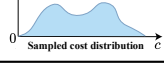
Receiver Operating Characteristics (ROC) is a popular tool to describe the trade-off between the True Positive Rate (TPR) and False Positive Rate (FPR) of a scoring function. AUC is defined by the area under the ROC curve [17, 18]. This metric naturally measures the average classification performance under different thresholds and is widely used (*e.g.*, disease prediction [19], and anomaly detection [29]). Compared with accuracy, AUC is insensitive to the threshold and cost [7], making it be a popular metric for long-tail learning [32] and achieve remarkable success [24, 44, 26].

Similar to AUC optimization, cost-sensitive learning is a common data mining method [10, 2, 4]. The main goal is to incorporate the misclassification costs in the model, which is more compatible with realistic scenarios (*e.g.*, the cost of misdiagnosing a disease as healthy is greater than the counterexample). Over the past two decades, researchers have pointed out that the ROC curve can be transferred to cost space by utilizing a threshold choice method, this is equivalent to computing the area under the convex hull of the ROC curve [21]. In this way, AUC can be seen as the performance of the model with a uniform cost distribution [16]. However, AUC considers all situations, which can not focus more on hard samples, Partial AUC (PAUC) is proposed as an extension of AUC with

---

\*Corresponding authors.

Table 1: Comparison with existing classification settings. Cost distribution represents the cost condition of each setting.

Different setting	Formulation	Attr.1	Attr.2	Cost distribution
Cost learning	$\mathbb{E}_x[c \cdot \pi \cdot p(1 x) + (1 - c) \cdot (1 - \pi) \cdot p(0 x)]$	×	×	
AUC/PAUC	$\mathbb{E}_\tau[\text{TPR}(\tau)\text{FPR}'(\tau)] \quad \tau \sim U(a, b)$	✓	×	
WAUC	$\mathbb{E}_\tau[\text{TPR}(\tau)\text{FPR}'(\tau)W(\tau)] \quad W(\tau) \sim \text{beta}(a, b)$	✓	×	
Our method	$\mathbb{E}_{\tau \sim \tau^*}[\text{TPR}(\tau)\text{FPR}'(\tau)] \quad \tau^* \in \arg \min_\tau \mathcal{L}_{\text{COST}}$	✓	✓	

truncated uniform cost distribution [31]. Recently, some studies extend PAUC using parameterized cost distributions and propose WAUC to fit real-world applications [16, 30].

Whether we use AUC or cost learning, our main purpose is to train models with these attributes:

**Attr.1:** The trained model can be robust to class distribution shift in the test **without class prior**.

**Attr.2:** The trained model can be robust to cost distribution in the test **without cost prior**.

However, to the best of our knowledge, there are few methods can train a model to have both of these attributes. According to Tab. 1, both AUC-related methods and cost-sensitive learning require a strong prior knowledge of the cost distribution; (1) The cost learning mainly considers a specified cost  $c$  and class imbalanced radio  $\pi$ . Models trained under this method are sensitive to class distribution, which does not apply to the scenario where test data distribution with offset. (2) AUC (PAUC) assumes that the cost distribution belongs to (truncated) uniform cost distribution  $U(a, b)$ . Models trained with them will have poor performance when the true cost distribution is not uniform [16]. (3) WAUC considers optimizing models based on more complex forms of cost distribution, such as  $\text{beta}(a, b)$ . However, we can not obtain the cost prior in real problem scenarios, e.g., financial market prediction [14]. Considering the weakness comes from the existing settings, we will explore the following question in this paper:

*Can we bridge AUC and complicated cost distribution to training robust model on desired cost-sensitive and arbitrary class imbalanced decision scenarios?*

To answer this question, we propose a view that, in some real applications [14], the cost, like the instance data, is not available prior but can be obtained by sampling. Therefore, we choose to sample desired cost to approximate the true cost distribution. Different from previous settings, ours is closer to real world, the main process can be divided into three parts:

**Step.1 Cost Sampling:** Firstly, we sample some desired costs to construct the empirical cost set.

**Step.2 Data Sampling:** Next, we sample some instance data to construct the empirical dataset.

**Step.3 Build Formulation:** Finally, we construct the appropriate formulation to maximize the performance in different desired costs and ensure model is robust to distribution shift.

It is natural for us to ask the question: Can we use the existing methods to realize this process? It's clear the answer is no. For AUC-related methods, they can not perform Step.1, and for cost-sensitive learning, they fail to achieve robust distribution shift and multiple costs in Step.3 (as shown in Fig. 1 (orange line)). Hence, we propose a novel bilevel formulation combining the advantages of WAUC and cost learning. The inner-level process calculates the optimal threshold from sampling costs, and the outer-level process minimizes the WAUC loss over the optimal threshold distribution. The method can help the model improve robustness to class distribution in cost-sensitive decision problems. The main process is shown in Fig. 1 (green line). We summarize our contributions below:

- We propose a setting that focuses on the robustness of the model to the class distribution and cost distribution simultaneously. This setting treats cost as data that can be sampled, not as prior information, which is closer to the real-world cost-sensitive scenario.

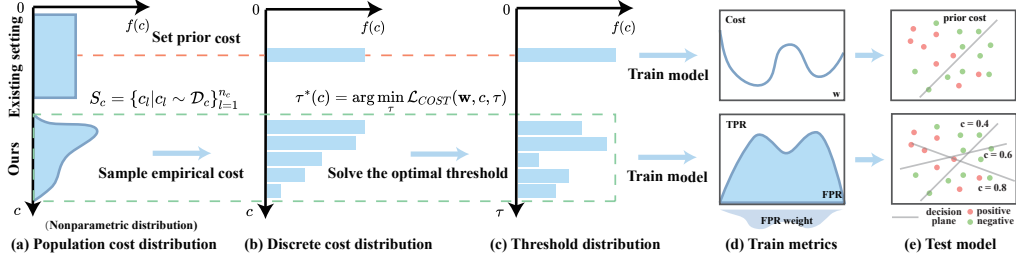


Figure 1: The comparison of our proposed setting with the previous setting. The orange line represents the previous cost-sensitive learning approach, and the green line represents our method.

- We present a bilevel paradigm where the inner cost function is an inner constraint of outer WAUC optimization. For sake of optimization, we reformulate this paradigm into a nonconvex-strongly convex bilevel form. Moreover, we employ a stochastic optimization algorithm for WAUC (SACCL), which can solve this problem efficiently.
- We conduct extensive experiments on multiple imbalanced cost-sensitive classification tasks. The experimental results speak to the effectiveness of our proposed methods.

## 2 Observation and Motivation

In Tab. 1, we compare the existing methods with ours from different views. However, the table comparison does not have a very visual presentation. In this section, we will analyze the disadvantages of existing settings and explain our motivation. We train the model with all methods on a Cifar-10-Long-Tail training set under the different imbalanced ratios and cost distribution. We visualize the feature representation (the last layer’s output of the trained model) in test data by t-SNE. The blue point represents negative samples predicted by the optimal threshold, and the orange point represents positive samples. The smaller the overlap between them, the better the performance. From Fig. 2, we can make the following remarks: (1) According to Fig. 2 (a), AUC is robust to changes in the imbalance ratio but completely not applicable with the cost distribution. (2) According to Fig. 2 (b), cost learning can process different cost distributions, but is sensitive to imbalance ratios.

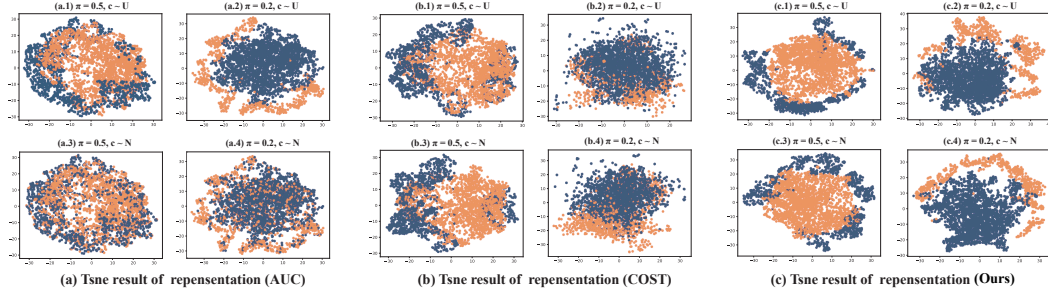


Figure 2: The feature representations comparison among different methods. (a) AUC optimization. (b) Cost learning. (c) Ours result. We solve the optimal threshold with the  $\hat{\mathcal{L}}_{COST}$  (defined in Sec. 3).  $\pi$  denotes the probability of positive class,  $c$  denotes the cost ratio of misclassification ( $U$  denotes Uniform,  $N$  denotes Normal). For example,  $\pi = 0.5, c \sim N$  means model tested on dataset which has imbalanced ratio  $\pi$  and cost set sampled from  $N$ .

Hence, our motivation is to propose a new approach to solve the problems in AUC and Cost-sensitive learning. As shown in Fig. 2 (c), our proposed method can have better learning results in various complex cost distributions and imbalance ratios. That means our method can overcome the shortcomings of traditional AUC and cost-learning, which perfectly fits the proposed setting.

## 3 Preliminaries

**Notations.** In this section, we give definitions and preliminaries about AUC. We denote  $(x, y)$  be an instance, where  $x$  is drawn from feature space  $\mathcal{X} \subseteq \mathbb{R}^d$  ( $d$  is feature number) and  $y$  is drawn from

label space  $\mathcal{Y} = \{0, 1\}$ . Let  $\mathcal{D}_P$  ( $\mathcal{D}_N$  *resp.*) be positive (negative *resp.*) instance distribution. Let  $\mathbf{x}^+ \sim \mathcal{D}_P$  ( $\mathbf{x}^- \sim \mathcal{D}_N$  *resp.*) be positive (negative *resp.*) instance. We denote  $S_+ = \{(\mathbf{x}_i^+, y_i)\}_{i=1}^{n_+}$  ( $S_- = \{(\mathbf{x}_j^-, y_j)\}_{j=1}^{n_-}$  *resp.*) as a set of training data drawn from  $\mathcal{D}_P$  ( $\mathcal{D}_N$  *resp.*), where  $n_+$  ( $n_-$  *resp.*) denotes the instance number of  $S_+$  ( $S_-$  *resp.*). Let  $\mathbb{I}_{(\cdot)}$  be the indicator function, which returns 1 when the condition is true and 0 otherwise. In this paper, we focus on the deep neural network scoring function  $s(\mathbf{w}, \mathbf{x}) : \mathcal{X} \mapsto [0, 1]$ , where parameterized by  $\mathbf{w}$  on an input  $\mathbf{x}$ .

**AUC & WAUC.** For specified threshold  $\tau$ , the TPR of a classifier derived from  $s(\mathbf{w}, \mathbf{x})$  measures the likelihood that it accurately predicts a positive instance when getting a random positive instance from  $\mathcal{D}_P$ . Formally, we have:

$$(\text{Pop.}) \text{TPR}_s(\tau) = \mathbb{P}_{\mathbf{x}^+ \sim \mathcal{D}_P}[s(\mathbf{w}, \mathbf{x}^+) > \tau] \quad (\text{Emp.}) \widehat{\text{TPR}}_s(\tau) = \frac{1}{n_+} \sum_{i=1}^{n_+} \mathbb{I}_{s(\mathbf{w}, \mathbf{x}_i^+) > \tau}. \quad (1)$$

In a similar spirit, the classifier's FPR on threshold  $\tau$  refers to the probability that it predicts positive when it gets a negative instance from  $\mathcal{D}_N$ .

$$(\text{Pop.}) \text{FPR}_s(\tau) = \mathbb{P}_{\mathbf{x}^- \sim \mathcal{D}_N}[s(\mathbf{w}, \mathbf{x}^-) > \tau] \quad (\text{Emp.}) \widehat{\text{FPR}}_s(\tau) = \frac{1}{n_-} \sum_{j=1}^{n_-} \mathbb{I}_{s(\mathbf{w}, \mathbf{x}_j^-) > \tau}. \quad (2)$$

AUC measures a scoring function's trade-off between TPR and FPR under uniform thresholds. Denote  $\tau$  drawn from the distribution  $\mathcal{D}_\tau$ , WAUC utilizes the threshold distribution explicitly based on AUC.  $\text{FPR}'_s$  denotes the probability density function of  $s(\mathbf{w}, \mathbf{x}^-)$ .

$$\text{AUC} = \int_{-\infty}^{\infty} \text{TPR}_s(\tau) \text{FPR}'_s(\tau) d\tau \quad (3a)$$

$$\text{WAUC} = \int_{-\infty}^{\infty} \text{TPR}_s(\tau) \text{FPR}'_s(\tau) p(\tau) d\tau, \quad (3b)$$

**Cost function [2].** In some real application scenarios, we need to consider the misclassification cost. We denote  $c_{(\cdot)}$  as misclassification cost for class  $(\cdot)$ , cost  $c$  drawn from  $\mathcal{D}_c$ . Since we could not obtain the cost distribution  $\mathcal{D}_c$ , we sample empirical set  $S_c = \{c_l\}_{l=1}^{n_c}$ ,  $n_c$  denotes the sample number of cost  $c$ . Given a scoring function  $s$  and parameter  $\mathbf{w}$ , the cost function  $\mathcal{L}_{COST}$  is (the empirical version of cost function,  $\hat{\mathcal{L}}_{COST}$  contains the empirical forms of TPR and FPR):

$$\mathcal{L}_{COST}(\mathbf{w}, c, \tau^*(c)) = c \cdot \pi \cdot (1 - \text{TPR}_s(\tau^*(c))) + (1 - c) \cdot (1 - \pi) \cdot \text{FPR}_s(\tau^*(c)), \quad (4)$$

where  $\pi = n_+ / (n_+ + n_-)$ ,  $c = c_+ / (c_+ + c_-)$  and  $\tau^*(c)$  is optimal threshold for score function  $s$  under specified  $c$  [21], the sample number  $n_\tau = n_c$ .

## 4 Problem Formulation

In this section, we introduce how to link the ROC curve to the cost space. First, we reformulate Eq.(3) into expectation:

$$\text{AUC} = \mathbb{E}_{\tau \sim U} [\text{TPR}_s(\tau) \cdot \text{FPR}'_s(\tau)] \quad (5a)$$

$$\text{WAUC} = \mathbb{E}_{\tau \sim \mathcal{D}_\tau} [\text{TPR}_s(\tau) \cdot \text{FPR}'_s(\tau)]. \quad (5b)$$

If threshold  $\tau$  is drawn from the uniform distribution  $U$ , WAUC will degrade to the standard AUC formulation. However, AUC only describes the global mean performance under all possible costs. If we want to extend AUC to cost-sensitive problems, maybe lifting the restriction on the uniform distribution of  $\tau$  is a good solution. Hence, we release the  $\mathcal{D}_\tau$ 's restriction to make it belongs to complicated distribution (*e.g.*, normal distribution, exponential distribution). Then we can extend AUC to WAUC. However, using WAUC raises another question: how do we get  $\mathcal{D}_\tau$ ? We find that  $\tau^*(c)$  is one of parameters of  $\mathcal{L}_{COST}(\mathbf{w}, c, \tau^*(c))$ . A natural idea is to use  $\mathcal{L}_{COST}$  to solve for the optimal  $\tau^*(c)$  and to combine the  $\tau^*(c)$  solved for at different  $c$  to obtain  $\mathcal{D}_\tau$ .

$$\tau^*(c) = \arg \min_{\tau} \mathcal{L}_{COST}(\mathbf{w}, \tau, c) = c \cdot \pi \cdot (1 - \text{TPR}_s(\tau)) + (1 - c) \cdot (1 - \pi) \cdot \text{FPR}_s(\tau), \quad (6)$$

If we couple Eq.(5b) and Eq.(6) together so that WAUC can enjoy the optimal threshold distribution in  $\mathcal{L}_{COST}$ , then we can break the barrier between the ROC curve and the cost space. With the help of the threshold as a bridge, we can extend the AUC metric to achieve the WAUC cost-sensitive learning. Then we give the problem formulation (intuitively, from the result of 2 (c), (OP0) satisfies both Attr.1 and Attr.2 simultaneously):

$$\begin{aligned} \text{(OP0) (outer.)} \quad & \text{WAUC} = \mathbb{E}_{\tau \sim \tau^*} [\text{TPR}_s(\tau) \cdot \text{FPR}'_s(\tau)] \\ \text{(inner.)} \quad & \tau^* = \{\tau^*(c) = \arg \min_{\tau} \mathcal{L}_{COST}(\mathbf{w}, \tau, c) | c \sim \mathcal{D}_c\} \end{aligned} \quad (7)$$

Nevertheless, there are still three main challenges in WAUC cost-sensitive learning:

- (1) Given the scoring function  $s$  and negative dataset  $S_-$ , how to estimate  $\text{FPR}'_s(\tau)$  in WAUC?
- (2) The inner problem is nonconvex, which is hard to give a theoretical convergence guarantee.
- (3) How to design a formulation that can bridge WAUC and  $\mathcal{L}_{COST}$  so that WAUC can be optimized over the cost distribution of the desired problem scenario?

We will address the challenge (1) in Sec. 5.1, challenge (2) in Sec. 5.2 and challenge (3) in Sec. 5.3.

## 5 Methodology

### 5.1 The Estimation of False Positive Rate

For challenge (1), we choose the kernel density estimation (KDE) to estimate  $\text{FPR}'_s(\tau)$  and denote it as  $K(x)$  (please see definition in Sec. C.1). Then we can address the density estimation problem. However, Eq.(5b) still exists non-differentiable and non-smooth term  $\mathbb{I}_{(\cdot)}$ , which is hard to optimize. Hence, we propose the following smooth and differentiable WAUC estimator to approximate Eq.(5b).

**Definition 5.1.** Denote  $K(x)$  be statistics kernel with bandwidth  $m$  and  $S_-^{\mathbf{w}} = \{s(\mathbf{w}, \mathbf{x}_j^-)\}_{j=1}^{n_-}$ . With Lemma 5.2, we have the approximate estimator and loss function for WAUC:

$$\widehat{\text{WAUC}} = \int_{-\infty}^{\infty} \text{TPR}_s(\tau) \mathcal{K}(S_-^{\mathbf{w}}, \tau) p(\tau) d\tau, \quad \widehat{\mathcal{L}}_{\text{WAUC}}(\mathbf{w}, \boldsymbol{\tau}) = \frac{1}{n_{\tau}} \sum_{l=1}^{n_{\tau}} \hat{h}(\mathbf{w}, \tau_l) \quad (8)$$

where  $\boldsymbol{\tau} = \{\tau_l\}_{l=1}^{n_{\tau}}$  and the point loss  $\hat{h}$  is defined by

$$\hat{h}(\mathbf{w}, \tau) = 1 - \frac{1}{n_+ n_-} \sum_{i=1}^{n_+} \sum_{j=1}^{n_-} \sigma(s(\mathbf{w}, \mathbf{x}_i^+) - \tau_l) \cdot K((s(\mathbf{w}, \mathbf{x}_j^-) - \tau_l)/m)/m. \quad (9)$$

$\sigma(x) = 1/(1 + \exp(-\beta x))$ ,  $\beta$  is smooth parameter and we have  $\sigma(x) \xrightarrow{\beta \rightarrow \infty} \mathbb{I}_x$ .

**Lemma 5.2.** Given a scoring function  $s$ , if  $\boldsymbol{\tau}$  is known, when the number of instances is large enough,  $\widehat{\text{WAUC}}$  almost surely converges to WAUC.

$$\lim_{n \rightarrow \infty} |\widehat{\text{WAUC}} - \text{WAUC}| \xrightarrow{a.s.} 0. \quad (10)$$

With KDE consistency [41], when the negative sample size is large enough, Lemma 5.2 provides theoretical approximation guarantees for our proposed WAUC estimator in Prop. 5.1.

### 5.2 The Estimation of Threshold Weighting

For challenge (2), a natural idea is to use  $\widehat{\mathcal{L}}_{COST}$  to solve for the optimal threshold set  $\hat{\tau}^*$  when given the cost set  $S_c$  and the scoring function  $s$ . Then we can use the optimal threshold set  $\hat{\tau}^*$  to calculate  $\widehat{\text{WAUC}}$ . Firstly, we define the solution for  $\hat{\tau}^*$  be

$$\hat{\tau}^* = \left\{ \hat{\tau}^*(c) | \hat{\tau}^*(c) \in \arg \min_{\tau} \widehat{\mathcal{L}}_{COST}(\mathbf{w}, c), c \in S_c \right\}. \quad (11)$$

However, it's noticed that the  $\arg \min_{\tau} \hat{\mathcal{L}}_{COST}$  in Eq.(11) is non-convex. As we analyzed before,  $\tau^*$  is the inner constraint of  $\widehat{\text{WAUC}}$ . To the best of our knowledge, there are few studies on optimizing two coupled non-convex problems simultaneously with theoretical convergence guarantees. Most studies on coupled optimization assume that the inner problems have good properties, such as strong convexity. Hence, we propose the approximated convex formulation of the inner problem for  $\hat{\tau}^*$ .

**Theorem 5.3.** *When we set  $\kappa, M$  are large positive numbers and  $M'^2 < M^2 \frac{6\kappa^2 e^{3\kappa}}{(e^\kappa + 1)^6}$ , then we have the approximated convex formulation for  $\hat{\mathcal{L}}_{COST}$*

$$\begin{aligned} \min_{\tau, \mathbf{P} \in \mathbb{R}^{n_+}, \mathbf{N} \in \mathbb{R}^{n_-}} \quad & \hat{\mathcal{L}}_{eq}(\mathbf{w}, \tau, c) := c \cdot \pi \cdot \left(1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i\right) + (1 - c) \cdot (1 - \pi) \cdot \left(\frac{1}{n_-} \sum_{j=1}^{n_-} N_j\right) \\ & + \frac{1}{n_+} \sum_{i=1}^{n_+} M' \psi(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) - P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) + M\psi(P_i - 1) + M\psi(\tau - 1) \\ & + \frac{1}{n_-} \sum_{j=1}^{n_-} M' \psi(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) - N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) + M\psi(N_j - 1) \quad 0 \leq \tau, P_i, N_j \end{aligned} \quad (12)$$

where  $\psi(x) = \log(1 + \exp(\kappa x))/\kappa$ .  $\hat{\mathcal{L}}_{eq}$  in Eq.(12) is  $\mu_g$ -strongly convex w.r.t.  $\tau$ . Eq.(12) has same solution as  $\min_{\tau} \hat{\mathcal{L}}_{COST}$  when the parameters satisfy the conditions of the penalty.

Thm. 5.3 provides an optimization method with good properties. Eq. (12) adopts the penalty function to convert inequality constraints into a part of the objective function. When these inequality constraints are not satisfied, the objective function will increase to infinity. Otherwise, we will get  $P_i = \mathbb{I}[s(\mathbf{w}, \mathbf{x}_i^+) > \tau]$  and  $N_j = \mathbb{I}[s(\mathbf{w}, \mathbf{x}_j^-) > \tau]$ , then we will get the same formulation as  $\hat{\mathcal{L}}_{COST}$ . When the parameters meet the requirements, Eq.(12) has the same solution as  $\hat{\mathcal{L}}_{COST}$ . We give the proof of Thm. 5.3 and the definition  $\mu$  in Sec. C.4. Moreover, we give the analysis of approximation error between  $\hat{\mathcal{L}}_{COST}$  and Thm. 5.3 in Sec. B.7.

### 5.3 Bilevel Optimization for WAUC learning

After answering questions in challenge (1) and (2), we have solved most of the problems in WAUC cost-sensitive learning. However, there remains a challenge (3) in optimization: How do we design learning paradigms to solve the coupled optimization problem of WAUC and  $\mathcal{L}_{COST}$ ? In recent years, bilevel optimization has achieved remarkable success. This approach can combine two related optimization problems to form a coupled optimization formulation. Hence, with Prop. 5.1 and Thm.5.3, we propose a bilevel paradigm to formulate this coupled optimization problem.

$$\begin{aligned} (OP1) \quad & \text{(outer.)} \quad \min_{\mathbf{w}} \hat{F}(\mathbf{w}) := \hat{f}(\mathbf{w}, \tau^*) := \hat{\mathcal{L}}_{\text{WAUC}}(\mathbf{w}, \hat{\tau}^*) \\ & \text{(inner.)} \quad \hat{\tau}^* = \arg \min_{\tau, \mathbf{P}_a, \mathbf{N}_a} \hat{g}(\mathbf{w}, \tau) := \frac{1}{n_{\tau}} \sum_{l=1}^{n_{\tau}} \hat{\mathcal{L}}_{eq}(\mathbf{w}, \tau_l, c_l), \end{aligned} \quad (13)$$

where  $\mathbf{P}_a \in \mathbb{R}^{n_{\tau} \times n_+}$  and  $\mathbf{N}_a \in \mathbb{R}^{n_{\tau} \times n_-}$ . (OP1) describes a bilevel optimization formulation for WAUC cost-sensitive learning, where the inner-level provides a threshold optimization process, and the outer-level minimizes the WAUC loss over the optimal threshold distribution. Moreover, this formulation is consistent with the mainstream bilevel optimization problem (outer-level is smooth and non-convex, inner-level is convex and smooth), which enjoys a faster convergence rate.

## 6 Optimization Algorithm

In this section, we focus on optimizing (OP1) in an end-to-end manner. Hence, we propose a stochastic algorithm for WAUC cost-sensitive learning shown in Alg. 1, which is referred to SAACL.

### 6.1 Main Idea of SAACL

We provide some intuitive explanations of our algorithm. At each iteration  $k$ , SAACL alternates between the inner-level gradient update on  $\tau$  and the outer-level gradient update on  $\mathbf{w}$ . During



---

**Algorithm 1** Stochastic Algorithm for WAUC Cost-sensitive Learning

---

**Input:** training data  $S_+$  and  $S_-$ , iteration numbers  $K$  and  $T$ , batch size  $B$ .

**Initialize:** parameters  $\mathbf{w}_0 \in \mathbb{R}^n$ ,  $\boldsymbol{\tau}_0 \in \mathbb{R}^{n_\tau}$ , stepsizes  $\alpha_k, \beta_k$ .

**for**  $k = 0$  **to**  $K$  **do**

  set  $\boldsymbol{\tau}_{k,0} = \boldsymbol{\tau}_k$ .

**for**  $t = 0$  **to**  $T$  **do**

    drawn  $\mathcal{B}_t = \{(\mathbf{x}_b, y_b)\}_{b=1}^B$  from  $S_+$  and  $S_-$  uniformly.

$\forall c_l \in S_c, \boldsymbol{\tau}_{k,t+1}^l = \boldsymbol{\tau}_{k,t}^l - \beta_k \nabla_{\boldsymbol{\tau}} \hat{g}(\mathbf{w}_k, c_l; \mathcal{B}_t)$ .

**end for**

  set  $\boldsymbol{\tau}_{k+1} = \boldsymbol{\tau}_{k,T}$

  drawn  $\mathcal{B}_k = \{(\mathbf{x}_b, y_b)\}_{b=1}^B$  from  $S_+$  and  $S_-$  uniformly.

$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k [\nabla_{\mathbf{w}} \hat{f}(\mathbf{w}_k, \boldsymbol{\tau}_{k+1}; \mathcal{B}_k) - \nabla_{\mathbf{w}}^2 \hat{g}(\mathbf{w}_k, \boldsymbol{\tau}_{k+1}; \mathcal{B}_k) \cdot$   
     $\left[ \frac{N}{L_{g,1}} \prod_{n=1}^{N'} \left( I - \frac{1}{L_{g,1}} \nabla_{\boldsymbol{\tau}}^2 \hat{g}(\mathbf{w}, \boldsymbol{\tau}_{k+1}; \mathcal{B}_k)^{-1} \right) \right] \nabla_{\boldsymbol{\tau}} \hat{f}(\mathbf{w}_k, \boldsymbol{\tau}_{k+1}; \mathcal{B}_k)]$

**end for**

---

iteration  $k$ , we update  $\boldsymbol{\tau}_{k,t}$  with standard SGD  $T$  steps to ensure that  $\boldsymbol{\tau}_{k+1}$  is as optimal as possible. After updating inner-level variables, we perform outer-level optimization with  $\boldsymbol{\tau}_{k+1}$  as the parameter to update  $\mathbf{w}_k$ . Notice that  $T$  will not take a large value to ensure the validity of the coupling update of  $\boldsymbol{\tau}$  and  $\mathbf{w}$ . Let  $\alpha_k$  and  $\beta_k$  be stepsizes of  $\mathbf{w}$  and  $\boldsymbol{\tau}$  that have the same decrease rate as SGD. We denote  $n$  be the number of elements in deep neural network parameters  $\mathbf{w}$ .

## 6.2 Convergence Analysis of SAACL

In this subsection, we present the convergence analysis for SAACL. We give some Lipschitz continuity assumptions that are common in bilevel optimization problems [11, 28].

**Assumption 6.1. (Lipschitz continuity)** Assume that  $f$ ,  $\nabla f$ , and  $\nabla g$  are respectively  $L_{f,0}, L_{f,1}, L_{g,1}$ -Lipschitz continuous.

**Assumption 6.2. (Bounded stochastic derivatives)** The variance of stochastic derivatives  $\nabla f(\mathbf{w}, \boldsymbol{\tau}; \mathcal{B})$  and  $\nabla g(\mathbf{w}, \boldsymbol{\tau}; \mathcal{B})$  are bounded by  $\sigma_{f,1}^2, \sigma_{g,1}^2$ , respectively.

Based on Assumption 6.1 and Assumption 6.2, following [5], Thm. 6.3 indicates that we can optimize (OP1) with the same convergence rate as the traditional SGD algorithm.

**Theorem 6.3.** Suppose Assumption 6.1 and 6.2 hold. We define

$$\bar{\alpha}_1 = \frac{1}{2L_F + 4L_f L_y + 2L_f L_{yx}/(L_y \eta)}, \quad \bar{\alpha}_2 = \frac{16T\mu L_{g,1}}{(\mu + L_{g,1})^2 (8L_f L_y + 2\eta L_{yx} \tilde{C}_f^2 \bar{\alpha}_1)}, \quad (14)$$

where  $\eta = L_F/L_y$ ,  $L_F, L_f, L_y$  and  $\tilde{C}_f^2$  come from Lem. 2 and Lem. 4 in [5]. We select the following stepsize as

$$\alpha_k = \min \left\{ \bar{\alpha}_1, \bar{\alpha}_2, \frac{1}{\sqrt{K}} \right\} \quad \beta_k = \frac{8L_f L_y + 2\eta L_{yx} \tilde{C}_f^2 \bar{\alpha}_1}{4T\mu} \alpha_k \quad (15)$$

For any  $T \geq 1$ , the iteration sequence  $\{\mathbf{w}_k\}_{k=1}^K$  and  $\{\boldsymbol{\tau}_k\}_{k=1}^K$  generated by Algorithm 1 satisfy

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla F(\mathbf{w}_k)\|^2] \leq \gamma \left( \frac{3M\kappa e^\kappa / (e^\kappa + 1)^2 + L_{g,1}}{24M\kappa e^\kappa / (e^\kappa + 1)^2 L_{g,1}} \right)^2 \frac{1}{T\sqrt{K}} + O\left(\frac{1}{\sqrt{K}}\right). \quad (16)$$

where  $\gamma = 2\alpha\sigma_{g,1}^2 \frac{L_f}{L_y} \left( 1 + 5L_f L_y \bar{\alpha}_1 + \frac{\eta L_{yx} \tilde{C}_f^2}{4} \bar{\alpha}_1^2 \right) (8L_f L_y + 2\eta L_{yx} \tilde{C}_f^2 \bar{\alpha}_1)^2$

**Remark 6.4.** When  $\kappa$  and  $M$  are large enough positive integers, according to Eq. (16), Alg. 1 is still guaranteed to find a  $\epsilon$ -stationary point within  $O(\epsilon^{-4})$  iterations ( $\epsilon$  is error tolerance).

## 7 Experiments

In this section, we conduct a series of experiments for WAUC cost-sensitive learning on common long-tail benchmark datasets. Due to space limitations, please refer to Sec. B for the details of our experiments. The source code is available in supplemental materials.

## 7.1 Dataset Details

We use three datasets: **Binary CIFAR-10-Long-Tail Dataset** [23], **Binary CIFAR-100-Long-Tail Dataset** [23], and **Jane Street Market Prediction** [14]. Binary CIFAR-10-Long-Tail Dataset and Binary CIFAR-100-Long-Tail Dataset are common datasets in long-tail learning, and we construct their cost distributions. Jane Street Market Prediction is data from real cost-sensitive learning application scenarios. For all datasets, we divide them into the training set, validation set, and test set with a proportion 0.7:0.15:0.15. All image data is normalized to ensure a more stable training process.

Table 2: Performance comparisons on benchmark datasets with different metrics. The first and second best results are highlighted with **bold text** and underline, respectively.

dataset	type	methods	Subset1		Subset2		Subset3		$\widehat{AUC} \uparrow$		
			$\widehat{WAUC} \uparrow$	$\widehat{\mathcal{L}}_{COST} \downarrow$	$\widehat{WAUC} \uparrow$	$\widehat{\mathcal{L}}_{COST} \downarrow$	$\widehat{WAUC} \uparrow$	$\widehat{\mathcal{L}}_{COST} \downarrow$	Subset1	Subset2	Subset3
CIFAR-10-LT	Competitors	BCE	0.525	0.027	0.533	0.015	0.318	0.029	<u>0.822</u>	<u>0.960</u>	<b>0.870</b>
		ExAUC	0.516	0.029	0.518	0.013	0.366	0.028	<b>0.845</b>	<b>0.963</b>	0.858
		SqAUC	0.407	0.028	0.548	0.012	0.327	0.031	0.811	0.933	<u>0.867</u>
		NWAUC	0.565	0.030	0.574	0.017	0.396	0.027	0.786	0.885	0.827
		PAUC-exp	0.549	0.029	0.508	0.015	0.354	0.028	0.650	0.801	0.736
		PAUC-poly	0.526	0.028	0.470	0.015	0.354	0.029	0.661	0.812	0.742
		PAUCI	0.516	0.027	0.520	0.015	0.382	0.028	0.704	0.847	0.734
		CS-hinge	0.566	0.026	0.633	<b>0.010</b>	0.377	0.022	0.675	0.782	0.762
		AdaCOS	0.576	0.025	0.559	0.014	0.391	0.023	0.758	0.873	0.742
		ECL	0.589	0.026	0.561	0.014	0.388	0.020	0.694	0.918	0.762
	Our method	WAUC-Gau	<b>0.679</b>	<u>0.024</u>	<u>0.660</u>	0.012	<u>0.467</u>	<u>0.015</u>	0.787	0.934	0.843
		WAUC-Log	<u>0.653</u>	<b>0.023</b>	<b>0.674</b>	<u>0.011</u>	<b>0.468</b>	<b>0.014</b>	0.820	0.958	<u>0.869</u>
CIFAR-100-LT	Competitors	BCE	0.556	0.022	0.463	0.012	0.512	0.019	<u>0.912</u>	<u>0.957</u>	0.806
		ExAUC	0.522	0.019	0.502	0.011	0.506	0.017	<b>0.933</b>	<b>0.967</b>	0.833
		SqAUC	0.483	0.024	0.367	0.015	0.474	0.018	0.889	0.955	<u>0.855</u>
		NWAUC	0.654	0.025	0.511	0.016	0.631	0.019	0.867	0.925	0.807
		PAUC-exp	0.464	0.020	0.282	0.014	0.469	0.016	0.826	0.811	0.787
		PAUC-poly	0.461	0.022	0.262	0.017	0.473	0.017	0.828	0.887	0.791
		PAUCI	0.549	0.018	0.439	0.016	0.514	0.018	0.812	0.843	<u>0.822</u>
		CS-hinge	0.523	0.017	0.457	0.010	0.515	0.014	0.734	0.910	0.716
		AdaCOS	0.590	0.018	0.474	0.011	0.587	0.016	0.769	0.919	0.727
		ECL	0.583	0.017	0.497	0.009	0.595	0.015	0.863	0.939	0.794
	Our method	WAUC-Gau	<b>0.745</b>	<u>0.015</u>	<b>0.589</b>	<u>0.005</u>	<u>0.728</u>	<u>0.013</u>	0.842	0.928	0.745
		WAUC-Log	<u>0.719</u>	<b>0.012</b>	<u>0.560</u>	<b>0.003</b>	<b>0.745</b>	<b>0.010</b>	0.906	0.960	<b>0.875</b>

Table 3: Performance comparisons on benchmark datasets in real world cost-sensitive problem. Profit means represents the money earned by the model over the entire trading period.

Methods	BCE	ExAUC	SqAUC	NWAUC	PAUC-exp	PAUC-poly	PAUCI	CS-hinge	AdaCOS	ECL	WAUC-Gau	WAUC-Log
$\widehat{WAUC} \uparrow$	0.5427	0.594 $\pm$ .003	0.508 $\pm$ .005	0.562 $\pm$ .002	0.576 $\pm$ .004	0.481 $\pm$ .005	0.529 $\pm$ .006	0.592 $\pm$ .002	0.6527 $\pm$ .005	0.625 $\pm$ .004	<b>0.698</b> $\pm$ .002	0.675 $\pm$ .001
$\widehat{\mathcal{L}}_{COST} \downarrow$	0.254 $\pm$ .004	0.269 $\pm$ .005	0.246 $\pm$ .003	0.251 $\pm$ .002	0.270 $\pm$ .004	0.246 $\pm$ .004	0.243 $\pm$ .001	0.237 $\pm$ .006	0.229 $\pm$ .004	0.226 $\pm$ .007	<b>0.209</b> $\pm$ .003	0.213 $\pm$ .002
$\widehat{AUC} \uparrow$	0.528 $\pm$ .005	<b>0.539</b> $\pm$ .004	0.526 $\pm$ .005	0.520 $\pm$ .004	0.529 $\pm$ .002	0.519 $\pm$ .005	0.510 $\pm$ .003	0.522 $\pm$ .005	0.5246 $\pm$ .002	<u>0.530</u> $\pm$ .003	0.526 $\pm$ .002	0.5235 $\pm$ .003
Profit $\uparrow$	4955 $\pm$ 20.14	5468 $\pm$ 17.90	5183 $\pm$ 30.91	5395 $\pm$ 22.48	5418 $\pm$ 14.06	4862 $\pm$ 28.04	4963 $\pm$ 15.09	5583 $\pm$ 30.05	5839 $\pm$ 34.92	5764 $\pm$ 25.09	<b>6526</b> $\pm$ 15.98	6308 $\pm$ 16.09

## 7.2 Overall Performance

In Tab. 2 and Tab. 3, we collect all the methods' performance on test sets of three types of datasets. For cost distribution of  $c$ , we sample some data from a normal distribution  $\mathcal{N}(0.5, 1)$  to construct a dataset  $S_c$  (we clip all data to  $[0, 1]$ ). We also conduct numerous experiments for other types of distribution of  $c$ , and please see Appendix B for the details. From the results, we make the following observations:

(1) For  $\widehat{WAUC}$  and  $\widehat{\mathcal{L}}_{COST}$  metric, Our proposed algorithm achieves superior performance in most benchmark datasets compared to other methods. This demonstrates that our proposed WAUC cost-sensitive learning can extend the ROC curve into the cost space. Models trained with our proposed bilevel optimization formulation can enjoy high WAUC and cost-related metrics.

(2) AUC and cost-related metrics are inconsistent. From the high-performing heatmap of Tab. 2, it can be noticed that  $\widehat{\mathcal{L}}_{COST}$  and  $\widehat{AUC}$  have two completely different highlight regions. This indicates that the assumption of uniform distribution of AUC does not match the realistic scenario.

(3) We also find that AUC-related and traditional classification algorithms do not perform well in cost-sensitive problems. This means that if we first train the model with the classification algorithm, subsequently using the cost function to solve for the optimal threshold for decision does not work well. Meanwhile, the algorithm that can learn from scratch has better scalability. Two-stage decision method Therefore, designing one-stage algorithms for WAUC cost-sensitive learning is necessary.



### 7.3 Sensitivity Analysis

In this subsection, we show the sensitivity of  $\beta$ ,  $T$ , and bandwidth on test data.

**Effect of  $\beta$ .** In Fig. 3 (a) and (d), we observe that for both WAUC and cost metrics, when  $\beta$  closes to 7, the model will have the largest performance improvement and the lowest variance. This can be explained in two ways: (1) When the  $\beta$  is too small, the error between the  $\sigma(x)$  and the 0-1 loss function is large, resulting in a large approximation error between the WAUC and the WAUC. (2) When the  $\beta$  is too large, the gradient also tends to be 0. Therefore, choosing a beta value that trades off the approximation error and the gradient is essential.

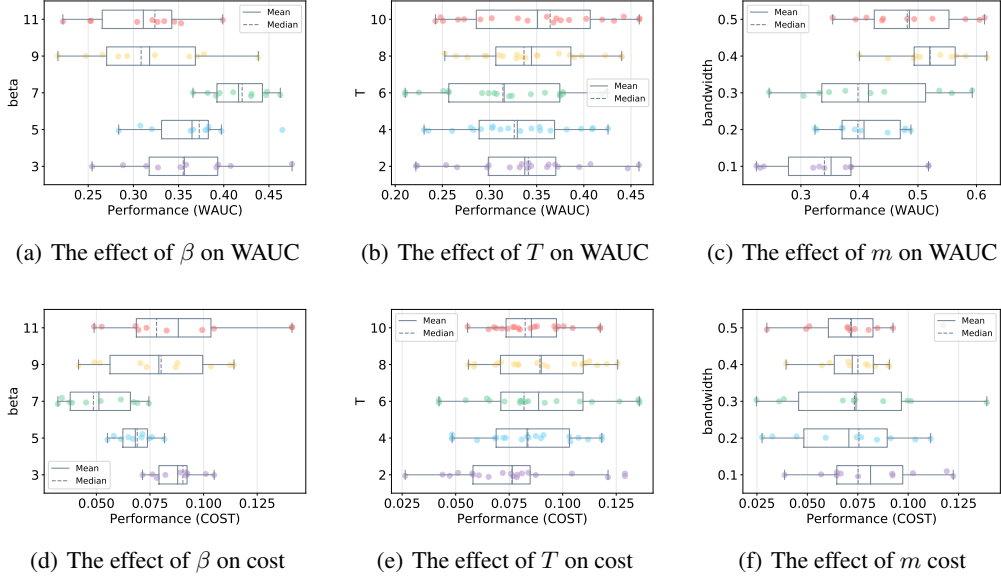


Figure 3: Sensitivity analysis on test data where WAUC and cost for WAUC-Gau with respect to  $\beta$ ,  $T$ , and bandwidth. The other two variables are fixed for each box in the plots, and the scattered points along the box show the variation.

**Effect of  $T$ .** As we mentioned in the Section 6.1, choosing a smaller  $T$  can effectively improve the performance of the model. However, as shown in Fig. 3(e), a larger  $T$  value can reduce the variance. Hence, as set in our experiments,  $T = 3$  is a good choice to ensure the average performance and variance of the model.

**Effect of  $m$ .** From Fig. 3(c), we find that the kernel’s bandwidth strongly influences the model’s performance. The model’s bandwidth and performance are almost proportional; the closer the bandwidth is to  $[0.4, 0.5]$ , the better the effect; otherwise, the effect is worse. This indicates that our proposed method is sensitive to the bandwidth parameter, which also compounds the bandwidth characteristics in the KDE method.

## 8 Conclusion

This paper focuses on extending the traditional AUC metric to associate with misclassification costs. Restricted by the assumption of cost distribution, existing settings could not describe the model’s performance in the complicated cost-sensitive scenario. To address this problem, we propose a novel setting that treats the cost as sampled data. We employ the WAUC metric and propose a novel estimator to approximate it. With the help of threshold weighting, we establish the correspondence between WAUC and the cost function. To describe this connection, we present a bilevel optimization formulation to couple them, where the inner-level problem provides a threshold optimization process, and the outer-level minimizes the WAUC loss based on the inner thresholds. This paradigm ensures that the WAUC can always be optimized at the optimal threshold value based on the complicated cost distribution in reality. Moreover, we propose a stochastic algorithm to optimize this formulation. We prove that our algorithm enjoys the same convergence rate as standard SGD. Finally, numerous

experiments have shown that our method can extend AUC to cost-sensitive scenarios with significant performance.

## Acknowledgements

This work was supported in part by the National Key R&D Program of China under Grant 2018AAA0102000, in part by National Natural Science Foundation of China: 62236008, U21B2038, U2001202, 61931008, 6212200758, 61976202, and 62206264, in part by the Fundamental Research Funds for the Central Universities, in part by Youth Innovation Promotion Association CAS, in part by the Strategic Priority Research Program of Chinese Academy of Sciences (Grant No. XDB28000000) and in part by the Innovation Funding of ICT, CAS under Grant No. E000000.

## References

- [1] J A, Hanley, B J, and McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 1982.
- [2] Daniel Andrade and Yuzuru Okajima. Efficient bayes risk estimation for cost-sensitive classification. In *The 22nd international conference on artificial intelligence and statistics*, pages 3372–3381. PMLR, 2019.
- [3] Jerome Bracken and James T McGill. Mathematical programs with optimization problems in the constraints. *Operations Research*, 21(1):37–44, 1973.
- [4] Nontawat Charoenphakdee, Zhenghang Cui, Yivan Zhang, and Masashi Sugiyama. Classification with rejection based on cost-sensitive classification. In *International Conference on Machine Learning*, pages 1507–1517. PMLR, 2021.
- [5] Tianyi Chen, Yuejiao Sun, and Wotao Yin. Tighter analysis of alternating stochastic gradient method for stochastic nested problems. *arXiv preprint arXiv:2106.13781*, 2021.
- [6] Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of operations research*, 153(1):235–256, 2007.
- [7] Corinna Cortes and Mehryar Mohri. AUC optimization vs. error rate minimization. In *Advances in Neural Information Processing Systems*, pages 313–320, 2003.
- [8] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9268–9277, 2019.
- [9] Tom Fawcett. An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874, 2006.
- [10] Alberto Fernández, Salvador García, Mikel Galar, Ronaldo C Prati, Bartosz Krawczyk, Francisco Herrera, Alberto Fernández, Salvador García, Mikel Galar, Ronaldo C Prati, et al. Cost-sensitive learning. *Learning from Imbalanced Data Sets*, pages 63–78, 2018.
- [11] Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pages 1568–1577. PMLR, 2018.
- [12] Yoav Freund, Raj Iyer, Robert E Schapire, and Yoram Singer. An efficient boosting algorithm for combining preferences. *Journal of Machine Learning Research*, 4:933–969, 2003.
- [13] Thore Graepel, Klaus Obermayer, et al. Large margin rank boundaries for ordinal regression. In *Advances in Large Margin Classifiers*, pages 115–132. 2000.
- [14] Jane Street Group. Jane street market prediction. <https://www.kaggle.com/competitions/jane-street-market-prediction/overview>, 2021.
- [15] Zhishuai Guo, Mingrui Liu, Zhuoning Yuan, Li Shen, Wei Liu, and Tianbao Yang. Communication-efficient distributed auc maximization with deep neural networks. In *International Conference on Machine Learning*, pages 3864–3874, 2020.

- [16] David J Hand. Measuring classifier performance: a coherent alternative to the area under the roc curve. *Machine learning*, 77(1):103–123, 2009.
- [17] James A Hanley and Barbara J McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982.
- [18] James A Hanley and Barbara J McNeil. A method of comparing the areas under receiver operating characteristic curves derived from the same cases. *Radiology*, 148(3):839–843, 1983.
- [19] Huaying Hao, Huazhu Fu, Yanwu Xu, Jianlong Yang, Fei Li, Xiulan Zhang, Jiang Liu, and Yitian Zhao. Open-narrow-synechia anterior chamber angle classification in as-oct sequences. *arXiv preprint arXiv:2006.05367*, 2020.
- [20] José Hernández-Orallo, Peter Flach, and César Ferri Ramírez. A unified view of performance metrics: Translating threshold choice into expected classification loss. *Journal of Machine Learning Research*, 13:2813–2869, 2012.
- [21] José Hernández-Orallo, Peter Flach, and César Ferri. Roc curves in cost space. *Machine learning*, 93(1):71–91, 2013.
- [22] Thorsten Joachims. A support vector method for multivariate performance measures. In *International Conference on Machine Learning*, pages 377–384, 2005.
- [23] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [24] Yunwen Lei and Yiming Ying. Stochastic proximal auc maximization. *JMLP*, 22(61):1–45, 2021.
- [25] Jialiang Li and Jason P Fine. Weighted area under the receiver operating characteristic curve and its application to gene selection. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 59(4):673–692, 2010.
- [26] Mingrui Liu, Xiaoxuan Zhang, Zaiyi Chen, Xiaoyu Wang, and Tianbao Yang. Fast stochastic auc maximization with  $o(1/n)$ -convergence rate. In *ICML*, pages 3189–3197. PMLR, 2018.
- [27] Mingrui Liu, Zhuoning Yuan, Yiming Ying, and Tianbao Yang. Stochastic AUC maximization with deep neural networks. In *International Conference on Learning Representations*, 2020.
- [28] Risheng Liu, Pan Mu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A generic first-order algorithmic framework for bi-level programming beyond lower-level singleton. In *International Conference on Machine Learning*, pages 6305–6315. PMLR, 2020.
- [29] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future frame prediction for anomaly detection—a new baseline. In *CVPR*, pages 6536–6545, 2018.
- [30] Andreas Maurer and Massimiliano Pontil. Estimating weighted areas under the roc curve. *Advances in Neural Information Processing Systems*, 33:7733–7742, 2020.
- [31] Donna Katzman McClish. Analyzing a portion of the roc curve. *Medical decision making*, 9(3): 190–195, 1989.
- [32] Charles E Metz. Basic principles of roc analysis. In *Seminars in nuclear medicine*, volume 8, pages 283–298. Elsevier, 1978.
- [33] Harikrishna Narasimhan and Shivani Agarwal. Support vector algorithms for optimizing the partial area under the roc curve. *Neural Computation*, 29(7):1919–1963, 2017.
- [34] Sakrapee Paisitkriangkrai, Chunhua Shen, and Anton Van Den Hengel. Efficient pedestrian detection by directly optimizing the partial area under the roc curve. In *Proceedings of the IEEE international conference on computer vision*, pages 1057–1064, 2013.
- [35] Margaret Sullivan Pepe and Mary Lou Thompson. Combining diagnostic test results to increase accuracy. *Biostatistics*, 1(2):123–140, 2000.

- [36] Alain Rakotomamonjy. Support vector machines and area under roc curve. *PSI-INSA de Rouen: Technical Report*, 2004.
- [37] Shoham Sabach and Shimrit Shtern. A first order method for solving convex bilevel optimization problems. *SIAM Journal on Optimization*, 27(2):640–660, 2017.
- [38] HuiYang Shao, Qianqian Xu, Zhiyong Yang, Shilong Bao, and Qingming Huang. Asymptotically unbiased instance-wise regularized partial auc optimization: Theory and algorithm. In *Advances in Neural Information Processing Systems*, 2022.
- [39] Luis N Vicente and Paul H Calamai. Bilevel and multilevel programming: A bibliography review. *Journal of Global optimization*, 5(3):291–306, 1994.
- [40] Sam Wieand, Mitchell H Gail, Barry R James, and Kang L James. A family of nonparametric statistics for comparing diagnostic markers with paired or unpaired data. *Biometrika*, 76(3): 585–592, 1989.
- [41] Dominik Wied and Rafael Weißbach. Consistency of the kernel density estimator: a survey. *Statistical Papers*, 53(1):1–21, 2012.
- [42] Lian Yan, Robert H Dodier, Michael Mozer, and Richard H Wolniewicz. Optimizing classifier performance via an approximation to the wilcoxon-mann-whitney statistic. In *International Conference on Machine Learning*, pages 848–855, 2003.
- [43] Tianbao Yang and Yiming Ying. AUC maximization in the era of big data and ai: A survey. *arXiv preprint arXiv:2203.15046*, 2022.
- [44] Zhenhuan Yang, Wei Shen, Yiming Ying, and Xiaoming Yuan. Stochastic auc optimization with general loss. *CPAA*, 19(8), 2020.
- [45] Zhiyong Yang, Qianqian Xu, Shilong Bao, Yuan He, Xiaochun Cao, and Qingming Huang. When all we need is a piece of the pie: A generic framework for optimizing two-way partial auc. In *International Conference on Machine Learning*, pages 11820–11829, 2021.
- [46] Yao Yao, Qihang Lin, and Tianbao Yang. Large-scale optimization of partial auc in a range of false positive rates. *arXiv preprint arXiv:2203.01505*, 2022.
- [47] Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online auc maximization. *Advances in Neural Information Processing Systems*, 29:451–459, 2016.
- [48] Zhuoning Yuan, Zhishuai Guo, Nitesh Chawla, and Tianbao Yang. Compositional training for end-to-end deep auc maximization. In *International Conference on Learning Representations*, 2021.
- [49] Adriano Z Zambom and Dias Ronaldo. A review of kernel density estimation with applications to econometrics. *International Econometric Review*, 5(1):20–42, 2013.
- [50] Dixian Zhu, Gang Li, Bokun Wang, Xiaodong Wu, and Tianbao Yang. When auc meets dro: Optimizing partial auc for deep learning with non-convex convergence guarantee. *arXiv preprint arXiv:2203.00176*, 2022.

## A Related Work

**ROC curve & cost-sensitive learning.** ROC curve and cost curve are two statistical tools frequently used in machine learning applications. Over the past two decades, some studies have explored their relationship and achieved significant success. [16] found that AUC implicitly uses a threshold weighting function corresponding to a cost weighting function. When these weighting functions output constant, we can infer that the cost function is a linear transformation of AUC. [21] shown that utilizing a natural threshold choice method can transfer ROC curves to cost space. [20] proposed a unified view of performance metrics. With the help of ROC convex hull [9], they give a clear interpretation of the threshold choice of the ROC curve. However, all of the above studies are based on the assumption of a uniform distribution in costs and thresholds. Moreover, their method can not be applied to end-to-end learning.

**AUC.** Since AUC offers some excellent properties in classification, it has become one of the standard performance metrics for binary imbalanced learning [1, 31, 43]. A partial list of the related studies includes [13, 7, 42, 22, 35, 12, 36, 47]. In deep age, there are some studies [27, 15, 48] focused on applying the AUC metric to deep end-to-end learning.

**PAUC.** The concept of PAUC was first introduced by [31], mainly used in disease diagnosis and biology. Prior to the deep learning era, earlier studies focused on optimizing PAUC in negative cost-sensitive scenarios. A partial list of the related studies includes [34, 33, 45, 50, 46, 38]. However, both AUC nor PAUC optimization does not consider the relationship between the ROC curve and cost space, which is hard to be applied in reality.

**WAUC.** The idea of weighting thresholds in AUC is first described by [40]. [25] constructed a framework for ROC analysis that incorporates the specificity distribution (*e.g.*, normal distribution). [30] proved exponential bounds on the estimation error of their proposed WAUC estimator and given conditions of the weight function. However, the weighting functions of these works are pre-selected, which are problem-independent, do not relate to the cost, and are too far from the practical application.

**Bilevel optimization.** Bilevel optimization is a classical algorithm for operations research. This formulation focuses on coupled optimization problems, where the inner optimization problem is the constraint of the outer optimization [3, 39, 6]. Recently, many studies have focused on designing gradient-based first-order algorithms to solve bilevel optimization problems [37, 11, 28, 5]. Most of them assume that the inner optimization problem has good properties, *e.g.*, strong convexity.

## B Experiment details

### B.1 Main Idea of Experiments

Our experiments mainly explore the following three problems:

- Our method versus the recent SOTA cost-sensitive learning algorithm. Direct optimization of cost function leads to the model is sensitive to the class distribution, while AUC has the advantage of being robust to the class distribution. Compared with previous cost-sensitive learning methods, we combine the advantages of AUC and cost metrics, allowing the model to enjoy a higher WAUC value while minimizing the cost. Models trained with our method can be applied to cost-sensitive decision problems (*e.g.*, financial market prediction, higher WAUC to guarantee decision profit)
- Traditional AUC is inconsistent with the cost-related metrics and cannot be used in cost-sensitive learning scenarios. From the experimental results in our paper, we can see that most AUC optimization methods do not minimize the misclassification cost. This indicates that in practical applications, using AUC optimization may maximize revenue without considering the cost. Ultimately, the misclassification cost of the decision is not acceptable.
- Our method versus the model trained with AUC/PAUC/CE first and then solves the optimal threshold with  $\hat{\mathcal{L}}_{COST}$  for a decision. We want to prove that model trained with AUC/PAUC/CE first and then getting the posterior threshold performs poorly in cost-sensitive learning. Hence, developing a novel one-stage method to address this problem is necessary.

## B.2 Dataset Details

**Binary CIFAR-10-Long-Tail Dataset.** The CIFAR-10 dataset [23] consists of 60,000 images divided into ten categories. By choosing one superclass as positive and the other as negative, we construct a long-tail binary version of CIFAR-10. For scalability, we generate three subsets of CIFAR-10 which are composed of the different superclasses, including 1) birds, 2) automobiles, and 3) cats.

**Binary CIFAR-100-Long-Tail Dataset.** Different from CIFAR-10, We choose a set of classes in CIFAR-100 [23] as positive and the rest of the classes as negative. Similarly, we construct three subsets and the positive set containing 1) insects, 2) vegetables and fruits, and 3) large omnivores and herbivores.

**Jane Street Market Prediction.** In reality, some applications like financial markets prediction [14] involve investment issues. Every investment has a cost involved. Developing trading strategies to identify and take advantage of inefficiencies is challenging. We adopt the actual financial markets data [14] to test all methods on real cost situation.

## B.3 Implementation Details

We conducted all experiments on a Ubuntu 16.04.1 server equipped with an Intel(R) Xeon(R) Silver 4110 CPU and four RTX 3090 GPUs. We implement all algorithms code in Python 3.8 and pytorch 1.8.2 environment. For the fairness of the experiment, we adopt the ResNet-18 model as the backbone of all competitors. The model’s output will be scaled into  $[0, 1]$  with a Sigmoid function. We set the batch size as 256 (64 per GPU) and epochs as 50. We employ `torch.optim.SGD` as the basic optimizer and `torch.nn.DataParallel` as tools for parallel computing. We set  $n_c$  as 50.

## B.4 Parameter Tuning

We tune the learning rate of all methods in the range  $[10^{-2}, 10^{-5}]$  and the weight decay in  $[10^{-3}, 10^{-5}]$ . Following the original paper, the warm-up steps for PAUC-poly and PAUC-exp is tuned in  $[3, 20]$ . Specifically, control parameter  $\gamma$  in PAUC-poly is searched in  $\{0.03, 0.05, 0.08, 0.1, 1, 3, 5\}$ . For PAUC-exp,  $\gamma$  is searched in  $[8, 30]$ . For PAUCI, stepsize related parameter  $k$  is tuned in  $[1, 10]$ ,  $\nu, \lambda, c_1, c_2$  are searched in  $[0, 1]$ .  $m, \kappa, \omega$  is tuned in  $[10, 100]$ ,  $[2, 6]$ ,  $[0, 4]$  respectively. For our method, we set inner loop iterations  $T = 3$ , tune bandwidth  $m$  in  $[0.1, 0.5]$ , smooth parameter  $\beta$  in  $[1, 10]$  and  $M, M', \kappa$  in  $[32, 64, 128, 256, 512, 1024]$ .

## B.5 Competitors

We implement two types of our methods, including Gaussian kernel  $\frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$  and Logistics kernel  $\frac{1}{e^x + 2 + e^{-x}}$  respectively. We denote them as WAUC-Gau and WAUC-Log. For other competitors, we classify them into three categories

- (1) Common methods for binary classification problems, including the class balanced CE loss (BCE) [8]; the Exp loss of AUC (ExAUC); the square loss of AUC (SqAUC); the naive WAUC estimator (NWAUC).
- (2) Recent SOTA methods of PAUC, including approximated PAUC estimator PAUC-poly [45] (poly calibrated weighting function) and PAUC-exp [45] (exp calibrated weighting function); asymptotically unbiased instance-wise regularized PAUC optimization [38] which is denoted as PAUCI.
- (3) The cost-sensitive learning algorithms, including cost-sensitive classification with rejection which is denoted as CS-hinge [4], Bayes risk estimator for cost-sensitive classification [2] which is denoted as AdaCOS, expected cost loss which is denoted as ECL.

## B.6 Experiment Results

We give the details of the dataset used in the experiment in the following table. All serial numbers and category names correspond to the information in the original dataset.

We also conduct experiments to study our method’s performance under other categories of cost distributions. We sampled the cost from the truncated uniform distribution  $U[0.5, 1]$  and beta



Table 4: Details of dataset.

Dataset	Pos. Class ID	Pos. Class Name	# Pos	#Neg
CIFAR-10-LT-1	2	birds	1,508	8,907
CIFAR-10-LT-2	1	automobiles	2,517	7,898
CIFAR-10-LT-3	3	birds	904	9,511
CIFAR-100-LT-1	6,7,14,18,24	insects	1,928	13,218
CIFAR-100-LT-2	0,51,53,57,83	fruits and vegetables	885	14,261
CIFAR-100-LT-3	15,19,21,32,38	large omnivores herbivores	1,172	13,974

distribution  $\text{beta}(2, 4)$ . We then conduct experiments on all competitors in subsets of CIFAR-10-Long-Tail and CIFAR-100-Long-Tail. All configurations are consistent with the sec.7 except for the cost data set  $S_c$ .

Table 5: Performance comparisons on benchmark datasets with different metrics (cost sampled from uniform distribution). The first and second best results are highlighted with **bold text** and underline, respectively.

dataset	type	methods	Subset1		Subset2		Subset3		$\widehat{\text{AUC}} \uparrow$		
			$\widehat{\text{WAUC}} \uparrow$	$\widehat{\mathcal{L}}_{\text{COST}} \downarrow$	$\widehat{\text{WAUC}} \uparrow$	$\widehat{\mathcal{L}}_{\text{COST}} \downarrow$	$\widehat{\text{WAUC}} \uparrow$	$\widehat{\mathcal{L}}_{\text{COST}} \downarrow$	Subset1	Subset2	Subset3
CIFAR-10-LT	Competitors	BCE	0.524	0.078	0.633	0.054	0.387	0.046	0.776	0.908	0.810
		ExAUC	0.472	0.076	0.709	0.098	0.456	0.050	0.787	0.736	0.856
		SqAUC	0.438	0.075	0.642	0.082	0.316	0.044	0.802	0.941	0.852
		NWAUC	0.573	0.063	0.714	0.072	0.469	0.041	0.777	0.882	0.809
		PAUC-exp	0.424	0.087	0.697	0.122	0.352	0.057	0.754	0.810	0.814
		PAUC-poly	0.431	0.087	0.629	0.094	0.328	0.054	0.763	0.760	0.766
		PAUCI	0.439	0.069	0.651	0.085	0.339	0.050	0.786	0.768	0.847
		CS-hinge	0.482	0.087	0.607	0.096	0.404	0.054	0.734	0.748	0.777
		AdaCOS	0.545	0.071	0.631	0.086	0.365	0.048	0.747	0.903	0.798
		ECL	0.469	0.087	0.749	0.110	0.307	0.042	0.697	0.665	0.842
	Our method	WAUC-Gau	0.642	0.065	0.781	0.040	0.511	0.042	0.790	0.942	0.833
		WAUC-Log	0.668	0.062	0.752	0.037	0.509	0.040	0.819	0.952	0.864
CIFAR-100-LT	Competitors	BCE	0.555	0.049	0.382	0.017	0.340	0.045	0.868	0.947	0.765
		ExAUC	0.601	0.069	0.517	0.011	0.428	0.039	0.710	0.956	0.906
		SqAUC	0.570	0.042	0.438	0.015	0.304	0.038	0.905	0.957	0.845
		NWAUC	0.490	0.051	0.552	0.016	0.384	0.038	0.860	0.954	0.842
		PAUC-exp	0.466	0.079	0.416	0.040	0.319	0.051	0.811	0.499	0.776
		PAUC-poly	0.424	0.067	0.429	0.028	0.319	0.051	0.751	0.836	0.784
		PAUCI	0.468	0.050	0.471	0.017	0.327	0.045	0.852	0.890	0.747
		CS-hinge	0.485	0.065	0.436	0.029	0.307	0.051	0.763	0.826	0.693
		AdaCOS	0.599	0.065	0.412	0.022	0.330	0.049	0.762	0.921	0.693
		ECL	0.503	0.078	0.504	0.031	0.318	0.051	0.845	0.787	0.722
	Our method	WAUC-Gau	0.687	0.047	0.547	0.016	0.409	0.043	0.869	0.932	0.763
		WAUC-Log	0.663	0.035	0.534	0.010	0.479	0.026	0.906	0.960	0.891

Under the condition of beta cost distribution, the performance of all methods in Tab. 6 is similar to Tab. 2. We can get a similar analysis result. Under the condition of uniform cost distribution, in Tab. 5, WAUC will degenerate to AUC. Therefore the AUC-related optimization algorithm will perform well and the gap with our proposed method will become smaller. Since PAUC is a special version of AUC based on the assumption of truncated uniform distribution of costs, the related algorithm has a clear advantage. However, although these algorithms have improved performance on the WAUC metric, the results on  $\widehat{\mathcal{L}}_{\text{COST}}$  are not good. Our proposed WAUC cost-sensitive learning can enjoy high WAUC metrics and cost metrics on both different kinds of cost distributions.

For the calculation of the  $\widehat{\text{WAUC}}$  metric needs to involve the solution of the optimal threshold  $\hat{\tau}^*$ . However, this optimization problem is non-convex and it is difficult to find the optimal solution quickly using existing optimization methods. Therefore, we propose an algorithm that can be solved within  $O(n_+ + n_-)$  iterations to obtain the optimal threshold, with the following details

## B.7 Additional Experiment

In Thm. 5.3, we propose a convex formulation which can approximate the  $\widehat{\mathcal{L}}_{\text{COST}}$ . According to the conditions of the penalty, we need  $\kappa, M \rightarrow \infty$  and  $M'^2 < M^2 \frac{6\kappa^2 e^{3\kappa}}{(e^\kappa + 1)^6}$  to ensure that Thm. 5.3 is approximately equivalent to  $\widehat{\mathcal{L}}_{\text{COST}}$ .

Table 6: Performance comparisons on benchmark datasets with different metrics (cost sampled from beta distribution). The first and second best results are highlighted with **bold text** and underline, respectively.

dataset	type	methods	Subset1		Subset2		Subset3		$\widehat{AUC} \uparrow$		
			WAUC $\uparrow$	$\widehat{L}_{COST} \downarrow$	WAUC $\uparrow$	$\widehat{L}_{COST} \downarrow$	WAUC $\uparrow$	$\widehat{L}_{COST} \downarrow$	Subset1	Subset2	Subset3
CIFAR-10-LT	Competitors	BCE	0.316	0.033	0.474	0.032	0.412	0.024	0.821	0.915	0.806
		ExAUC	0.389	0.035	0.474	0.037	0.491	0.021	0.841	<b>0.967</b>	<u>0.867</u>
		SqAUC	0.270	0.033	0.495	0.039	0.425	0.022	<b>0.858</b>	0.933	0.855
		NWAUC	0.326	0.034	0.538	0.033	0.462	0.023	0.817	0.922	0.826
		PAUC-exp	0.071	0.033	0.490	0.059	0.476	0.027	0.783	0.715	0.740
		PAUC-poly	0.137	0.035	0.477	0.045	0.445	0.026	0.746	0.830	0.703
		PAUCI	0.185	0.038	0.526	0.038	0.476	0.022	0.796	0.844	0.787
		CS-hinge	0.365	0.032	0.574	0.035	0.556	0.025	0.756	0.897	0.771
		AdaCOS	0.347	0.038	0.567	0.036	0.516	0.024	0.803	0.900	0.808
		ECL	0.303	0.034	<u>0.584</u>	<u>0.030</u>	0.537	0.022	<u>0.847</u>	0.925	0.856
	Our method	WAUC-Gau	<b>0.413</b>	<u>0.031</u>	<b>0.640</b>	<u>0.030</u>	<b>0.610</b>	<u>0.020</u>	0.815	0.943	0.830
		WAUC-Log	<u>0.393</u>	<b>0.029</b>	0.578	<b>0.024</b>	<u>0.575</u>	<b>0.019</b>	<u>0.852</u>	<u>0.959</u>	<b>0.869</b>
CIFAR-100-LT	Competitors	BCE	0.547	0.025	0.423	0.016	0.190	0.027	0.869	0.930	0.757
		ExAUC	0.646	0.023	0.459	0.015	0.184	0.021	<b>0.929</b>	0.948	<b>0.900</b>
		SqAUC	0.460	0.023	0.396	0.009	0.157	0.025	0.892	<u>0.949</u>	0.854
		NWAUC	<u>0.652</u>	0.024	0.494	0.009	0.151	0.026	0.889	0.941	0.819
		PAUC-exp	0.566	0.034	0.464	0.015	0.184	0.028	0.800	0.835	0.793
		PAUC-poly	0.269	0.029	0.410	0.012	0.167	0.028	0.788	0.885	0.740
		PAUCI	0.530	0.024	0.449	<u>0.007</u>	0.193	0.025	0.827	0.895	0.705
		CS-hinge	0.503	0.026	0.423	0.010	0.272	<u>0.017</u>	0.847	0.899	0.749
		AdaCOS	0.605	0.027	0.452	0.011	0.244	<u>0.017</u>	0.852	0.915	0.696
		ECL	0.503	0.025	0.424	0.010	0.220	<b>0.016</b>	0.872	0.930	0.799
	Our method	WAUC-Gau	<b>0.747</b>	<u>0.022</u>	<b>0.658</b>	0.008	<u>0.275</u>	0.023	0.869	0.925	0.757
		WAUC-Log	0.647	<b>0.018</b>	<u>0.645</u>	<b>0.005</b>	<b>0.290</b>	<u>0.017</u>	<u>0.911</u>	<b>0.961</b>	<u>0.865</u>

**Algorithm 2** Algorithm for Solving the Optimal Threshold

0

**Input:** test data  $S_+^t$  and  $S_-^t$ , cost dataset  $S_c$ . **Initialize:** parameters  $\hat{\tau}^* = \{0\}_{l=1}^{n_c}$   
**for**  $l = 0$  **to**  $n_c$  **do**  
     $\hat{\tau}^*[l] = \arg \min_{\tau \in \{S_+^t, S_-^t\}} \widehat{L}_{COST}$   
**end for**

In this subsection, we study the optimal value gap between original  $\widehat{L}_{COST}$  problem and Thm. 5.3. We optimize the  $\widehat{L}_{COST}$  and Eq. (12) in Cifar-10-Long-Tail training set. We set cost distribution as Uniform and  $\pi = 0.5$ . We calculate the mean optimal value ( $p^* = \min \widehat{L}_{COST}$ ,  $d^* = \min$  Eq. (12)) of them over the cost distribution. Finally, we calculate the error between them  $((p^* - d^*)^2)$  and list the results in Tab. 7.

Table 7: Optimal value gap between original  $\widehat{L}_{COST}$  optimization problem and Thm. 5.3

	$M = 64$	$M = 128$	$M = 256$	$M = 512$	$M = 1024$	$M = 2048$
	$\kappa = 64$	$\kappa = 128$	$\kappa = 256$	$\kappa = 512$	$\kappa = 1024$	$\kappa = 2048$
$M' = 32$	0.058	0.050	0.048	0.045	0.042	0.041
$M' = 64$	0.049	0.042	0.039	0.036	0.034	0.032
$M' = 128$	0.038	0.036	0.032	0.029	0.027	0.025
$M' = 256$	0.025	0.023	0.020	0.018	0.016	0.015
$M' = 512$	0.018	0.015	0.013	0.011	0.009	0.008
$M' = 1024$	0.012	0.010	0.008	0.007	0.007	0.005

In Tab. 7, we find that when  $\kappa$ ,  $M$ ,  $M'$  grows, the optimal value gap will increase quickly. Moreover, when  $\kappa = M = 64$  and  $M' = 32$ , the error is small enough. Hence, in real-world application, the approximation error and effect of hyparameters is acceptable.

## C Proofs for Section 4

### C.1 KDE Definition

Here, we give the notation and definition of KDE.

**Definition C.1 (Kernel density estimation [49]).** Denote  $x$  as a random variable with probability density function  $f_x$ . Given a dataset  $S = \{x_i\}_{i=1}^{n_x}$  and threshold  $\tau$ , we denote

$$\mathcal{K}(S, \tau) = \frac{1}{|S|m} \sum_{x_i \in S} K\left(\frac{x_i - \tau}{m}\right), \quad (17)$$

as an estimation of  $f_x$ , where  $m$  is bandwidth. The non-negative real-valued integrable function  $K$  satisfies

$$(1) \int_{-\infty}^{\infty} K(x)dx = 1, \quad (2) K(x) = K(-x). \quad (18)$$

## C.2 Proof of Proposition 5.1

**Restatement of Proposition 5.1.** Denote  $K(x)$  be statistics kernel with bandwidth  $m$  and  $S_-^w = \{s(\mathbf{w}, \mathbf{x}_j^-)\}_{j=1}^{n_-}$ . With Lemma 5.2, we have the approximate estimator and loss function for WAUC:

$$\widehat{\text{WAUC}} = \int_{-\infty}^{\infty} \text{TPR}_s(\tau) \mathcal{K}(S_-^w, \tau) p(\tau) d\tau, \quad \hat{\mathcal{L}}_{\text{WAUC}}(\mathbf{w}, \boldsymbol{\tau}) = \frac{1}{n_\tau} \sum_{l=1}^{n_\tau} \hat{h}(\mathbf{w}, \tau_l) \quad (19)$$

where  $\boldsymbol{\tau} = \{\tau_l\}_{l=1}^{n_\tau}$  and the point loss  $\hat{h}$  is defined by

$$\hat{h}(\mathbf{w}, \tau) = 1 - \frac{1}{n_+ n_-} \sum_{i=1}^{n_+} \sum_{j=1}^{n_-} \sigma(s(\mathbf{w}, \mathbf{x}_i^+) - \tau_l) \cdot K((s(\mathbf{w}, \mathbf{x}_j^-) - \tau_l)/m)/m. \quad (20)$$

$\sigma(x) = 1/(1 + \exp(-\beta x))$ ,  $\beta$  is smooth parameter and we have  $\sigma(x) \xrightarrow{\beta \rightarrow \infty} \mathbb{I}_x$ .

*Proof.*

$$\begin{aligned} \widehat{\text{WAUC}} &= \int_{-\infty}^{\infty} \text{TPR}(\tau) \mathcal{K}(S_-, \tau) p(\tau) d\tau \\ &= \mathbb{E}_\tau \left[ \mathbb{P}_{\mathbf{x}^+} [s(\mathbf{w}, \mathbf{x}^+) > \tau] \cdot \mathbb{E}_{\mathbf{x}^-} \left[ \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m}\right) \right] \right] \\ &= \mathbb{E}_\tau \left[ \mathbb{E}_{\mathbf{x}^+} [\mathbb{I}_{s(\mathbf{w}, \mathbf{x}^+) > \tau}] \cdot \mathbb{E}_{\mathbf{x}^-} \left[ \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m}\right) \right] \right] \\ &= \mathbb{E}_{\tau, \mathbf{x}^+} \left[ \mathbb{I}_{s(\mathbf{w}, \mathbf{x}^+) > \tau} \cdot \mathbb{E}_{\mathbf{x}^-} \left[ \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m}\right) \right] \right] \\ &= \mathbb{E}_{\tau, \mathbf{x}^+, \mathbf{x}^-} \left[ \mathbb{I}_{s(\mathbf{w}, \mathbf{x}^+) > \tau} \cdot \left[ \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m}\right) \right] \right]. \end{aligned} \quad (21)$$

Replacing the  $\mathbb{I}_{(\cdot)}$  function with  $\sigma(x) = 1/(1 + \exp(-\beta x))$  and change it to empirical formulation.

$$\hat{\mathbb{E}}_{\mathbf{x}^+ \sim S_+, \mathbf{x}^- \sim S_-, \tau \sim \boldsymbol{\tau}} \left[ \sigma(s(\mathbf{w}, \mathbf{x}^+) - \tau) \cdot \left[ \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m}\right) \right] \right] \quad (22)$$

Since we want to maximize the WAUC metric, we employ  $1 - \widehat{\text{WAUC}}$  as loss function, then we have

$$\begin{aligned} (\text{Pop.}) \quad \mathcal{L}_{\text{WAUC}} &= \mathbb{E}_{\mathbf{x}^+ \sim \mathcal{D}_P, \mathbf{x}^- \sim \mathcal{D}_N} \left[ 1 - \frac{1}{n_\tau} \sum_{l=1}^{n_\tau} \sigma(s(\mathbf{w}, \mathbf{x}^+) - \tau_l) \cdot \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau_l}{m}\right) \right] \\ (\text{Emp.}) \quad \hat{\mathcal{L}}_{\text{WAUC}} &= \hat{\mathbb{E}}_{\mathbf{x}^+ \sim S_+, \mathbf{x}^- \sim S_-} \left[ 1 - \frac{1}{n_\tau} \sum_{l=1}^{n_\tau} \sigma(s(\mathbf{w}, \mathbf{x}^+) - \tau_l) \cdot \frac{1}{m} K\left(\frac{s(\mathbf{w}, \mathbf{x}^-) - \tau_l}{m}\right) \right]. \end{aligned} \quad (23)$$

Reformulating the  $\hat{\mathcal{L}}_{\text{WAUC}}$ , we have:

$$\hat{\mathcal{L}}_{\text{WAUC}}(\mathbf{w}, \boldsymbol{\tau}) = \frac{1}{n_\tau} \sum_{l=1}^{n_\tau} \hat{h}(\mathbf{w}, \tau_l) \quad (24)$$

where

$$\hat{h}(\mathbf{w}, \tau) = 1 - \frac{1}{n_+ n_-} \sum_{i=1}^{n_+} \sum_{j=1}^{n_-} \sigma(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) \cdot K((s(\mathbf{w}, \mathbf{x}_j^-) - \tau)/m)/m. \quad (25)$$

$\sigma(x) = 1/(1 + \exp(-\beta x))$ ,  $\beta$  is smooth parameter and we have  $\sigma(x) \xrightarrow{\beta \rightarrow \infty} \mathbb{I}_x$  to approximate  $\widehat{\text{TPR}}$  and  $\widehat{\text{FPR}}$ .  $\square$

### C.3 Proof of Lemma 5.2

**Restatement of Lemma 5.2.** *Given a scoring function  $s$ , if  $\tau$  is known, when the number of instances is large enough,  $\widehat{\text{WAUC}}$  almost surely converges to WAUC.*

$$\lim_{n_- \rightarrow \infty} |\widehat{\text{WAUC}} - \text{WAUC}| \xrightarrow{a.s.} 0. \quad (26)$$

*Proof.*

$$\begin{aligned} & \lim_{n_- \rightarrow \infty} |\widehat{\text{WAUC}} - \text{WAUC}| \\ &= \lim_{n_- \rightarrow \infty} \left| \hat{\mathbb{E}}_{\tau, \mathbf{x}^+} \left[ \sigma(s(\mathbf{w}, \mathbf{x}^+) - \tau) \cdot \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] \right] - \mathbb{E}_{\tau, \mathbf{x}^+} [\mathbb{I}_{s(\mathbf{w}, \mathbf{x}^+) > \tau} \cdot \text{FPR}'_s(\tau)] \right| \\ &\leq \lim_{n_- \rightarrow \infty} \sup_{\tau, \mathbf{x}^+} \left| \sigma(s(\mathbf{w}, \mathbf{x}^+) - \tau) \cdot \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] - \mathbb{I}_{s(\mathbf{w}, \mathbf{x}^+) > \tau} \cdot \text{FPR}'_s(\tau) \right| \\ &\leq \lim_{n_- \rightarrow \infty} \sup_{\tau} \max \left\{ \underbrace{\left| \sigma(\delta) \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] - \text{FPR}'_s(\tau) \right|}_{\sigma(\delta) \rightarrow 1 \text{ when } \beta \rightarrow \infty}, \underbrace{\left| \sigma(-\delta) \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] \right|}_{\sigma(-\delta) \rightarrow 0 \text{ when } \beta \rightarrow \infty} \right\} \\ &\leq \lim_{n_- \rightarrow \infty} \sup_{\tau} \left| \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] - \text{FPR}'_s(\tau) \right| \\ &\stackrel{(a)}{\leq} \underbrace{\lim_{n_- \rightarrow \infty} \sup_{\tau} \left| \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] - \mathbb{E} \left[ \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] \right] \right|}_{(c) \text{ Kernel density function consistency lemma}} \\ &\quad + \underbrace{\lim_{n_- \rightarrow \infty} \sup_{\tau} \left| \mathbb{E} \left[ \hat{\mathbb{E}}_{\mathbf{x}^-} \left[ \frac{1}{m} K \left( \frac{s(\mathbf{w}, \mathbf{x}^-) - \tau}{m} \right) \right] \right] - \text{FPR}'_s(\tau) \right|}_{(d) \text{ Law of large numbers}} \\ &= 0 \end{aligned} \quad (27)$$

where (a) comes from triangle inequality and  $\delta > 0$ . We assume that  $\mathbb{E}_{\tau, \mathbf{x}^+} [\mathbb{I}_{s(\mathbf{w}, \mathbf{x}^+) = \tau}] = 0$ . For terms (c) and (d), please see the proof of [41].  $\square$

### C.4 Proof of Theorem 5.3

**Restatement of Theorem 5.3.** *When we set  $\kappa, M$  are large positive numbers and  $M'^2 < M^2 \frac{6\kappa^2 e^{3\kappa}}{(e^\kappa + 1)^6}$ , then we have the approximated convex formulation for  $\hat{\mathcal{L}}_{\text{COST}}$*

$$\begin{aligned} & \min_{\tau, \mathbf{P} \in \mathbb{R}^{n_+}, \mathbf{N} \in \mathbb{R}^{n_-}} \hat{\mathcal{L}}_{eq}(\mathbf{w}, \tau, c) := c \cdot \pi \cdot \left( 1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i \right) + (1 - c) \cdot (1 - \pi) \cdot \left( \frac{1}{n_-} \sum_{j=1}^{n_-} N_j \right) \\ & + \frac{1}{n_+} \sum_{i=1}^{n_+} M' \psi(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) - P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) + M \psi(P_i - 1) + M \psi(\tau - 1) \\ & + \frac{1}{n_-} \sum_{j=1}^{n_-} M' \psi(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) - N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) + M \psi(N_j - 1) \quad 0 \leq \tau, P_i, N_j \end{aligned} \quad (28)$$

where  $\psi(x) = \log(1 + \exp(\kappa x)) / \kappa$ .  $\hat{\mathcal{L}}_{eq}$  in Eq.(12) is  $\mu_g$ -strongly convex w.r.t.  $\tau$ . Eq.(12) has same solution as  $\min_{\tau} \hat{\mathcal{L}}_{COST}$  when the parameters satisfy the conditions of the penalty.

*Proof.* According to the definition of  $\hat{\mathcal{L}}_{COST}$ , we have:

$$\min_{\tau} c \cdot \pi \cdot \left( 1 - \frac{1}{n_+} \sum_{i=1}^{n_+} \mathbb{I}_{[s(\mathbf{w}, \mathbf{x}_i^+) - \tau]} \right) + (1-c) \cdot (1-\pi) \cdot \left( \frac{1}{n_-} \sum_{j=1}^{n_-} \mathbb{I}_{[s(\mathbf{w}, \mathbf{x}_j^-) - \tau]} \right) \quad (29)$$

$s.t. \quad 0 \leq \tau \leq 1.$

Then we have the equivalent formulation:

$$\begin{aligned} \min_{\tau, \mathbf{P}, \mathbf{N}} c \cdot \pi \cdot \left( 1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i \right) + (1-c) \cdot (1-\pi) \cdot \left( \frac{1}{n_-} \sum_{j=1}^{n_-} N_j \right) \\ s.t. \max(s(\mathbf{w}, \mathbf{x}_i^+) - \tau, 0) = P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) \\ \max(s(\mathbf{w}, \mathbf{x}_j^-) - \tau, 0) = N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) \\ 0 \leq \tau, P_i, N_j \leq 1, \end{aligned} \quad (30)$$

where  $\mathbf{P} \in \mathbb{R}^{n_+}$ ,  $\mathbf{N} \in \mathbb{R}^{n_-}$  and  $P_i \in \mathbf{P}$ ,  $N_j \in \mathbf{N}$ . Since the equality constraint is hard to process, we convert it to inequality constraint (e.g.,  $a = b \Leftrightarrow a \leq b, a \geq b$ ).

$$\begin{aligned} \min_{\tau, \mathbf{P}, \mathbf{N}} c \cdot \pi \cdot \left( 1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i \right) + (1-c) \cdot (1-\pi) \cdot \left( \frac{1}{n_-} \sum_{j=1}^{n_-} N_j \right) \\ s.t. \forall i \max(s(\mathbf{w}, \mathbf{x}_i^+) - \tau, 0) \geq P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) \\ \forall j \max(s(\mathbf{w}, \mathbf{x}_j^-) - \tau, 0) \geq N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) \\ \forall i \max(s(\mathbf{w}, \mathbf{x}_i^+) - \tau, 0) \leq P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) \\ \forall j \max(s(\mathbf{w}, \mathbf{x}_j^-) - \tau, 0) \leq N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) \\ \forall i, j \quad 0 \leq \tau, P_i, N_j \leq 1 \end{aligned} \quad (31)$$

Due to the fact that  $\max(s(\mathbf{w}, \mathbf{x}_i^+) - \tau, 0) \geq P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau)$  and  $\max(s(\mathbf{w}, \mathbf{x}_j^-) - \tau, 0) \geq N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau)$  is ground truth all the time. Hence, we have

$$\begin{aligned} \min_{\tau, \mathbf{P}, \mathbf{N}} c \cdot \pi \cdot \left( 1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i \right) + (1-c) \cdot (1-\pi) \cdot \left( \frac{1}{n_-} \sum_{j=1}^{n_-} N_j \right) \\ s.t. \forall i \max(s(\mathbf{w}, \mathbf{x}_i^+) - \tau, 0) \leq P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) \\ \forall j \max(s(\mathbf{w}, \mathbf{x}_j^-) - \tau, 0) \leq N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) \\ \forall i, j \quad 0 \leq \tau, P_i, N_j \leq 1 \end{aligned} \quad (32)$$

Then we apply the penalty function method to convert the constraint optimization into approximated unconstrained optimization:

$$\begin{aligned} \min_{\tau, \mathbf{P}, \mathbf{N}} c \cdot \pi \cdot \left( 1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i \right) + (1-c) \cdot (1-\pi) \cdot \left( \frac{1}{n_-} \sum_{j=1}^{n_-} N_j \right) + M\psi(\tau - 1) \\ + \frac{1}{n_+} \sum_{i=1}^{n_+} M'(\psi(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) - P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) + M\psi(P_i - 1)) \\ + \frac{1}{n_-} \sum_{j=1}^{n_-} M'(\psi(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) - N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) + M\psi(N_j - 1)) \\ \forall i, j \quad 0 \leq \tau, P_i, N_j \end{aligned} \quad (33)$$

where  $\psi(x) = \frac{\log(1+\exp(\kappa x))}{\kappa}$  is penalty function ( $\psi(x) \xrightarrow{\kappa \rightarrow \infty} \max(x, 0)$ ),  $M$  and  $M'$  denote positive number which are large enough. It's noticed that when  $\kappa, M, M' \rightarrow \infty$ , then Eq.(33) is equivalent to Eq.(32). Next, we will prove the strong convexity of  $\tau$  in Eq.(33). Firstly, we give the hessian matrix of Eq.(33):

$$H = M \begin{bmatrix} \frac{\kappa e^{\kappa(\tau-1)}}{(e^{\kappa(\tau-1)}+1)^2} & & \\ + \frac{1}{n_+} \sum_{i=1}^{n_+} \frac{\kappa e^{\kappa(P_i-\tau)}}{(e^{\kappa(P_i-\tau)}+1)^2} & M'/M & M'/M \\ + \frac{1}{n_-} \sum_{j=1}^{n_-} \frac{\kappa e^{\kappa(N_j-\tau)}}{(e^{\kappa(N_j-\tau)}+1)^2} & & \\ M'/M & \frac{1}{n_+} \sum_{i=1}^{n_+} \frac{\kappa e^{\kappa(P_i-1)}}{(e^{\kappa(P_i-1)}+1)^2} & 0 \\ M'/M & 0 & \frac{1}{n_-} \sum_{j=1}^{n_-} \frac{\kappa e^{\kappa(N_j-1)}}{(e^{\kappa(N_j-1)}+1)^2} \end{bmatrix} \quad (34)$$

For computational simplicity, we define

$$\begin{aligned} x &= \frac{\kappa e^{\kappa(\tau-1)}}{(e^{\kappa(\tau-1)}+1)^2} + \frac{1}{n_+} \sum_{i=1}^{n_+} \frac{\kappa e^{\kappa(P_i-\tau)}}{(e^{\kappa(P_i-\tau)}+1)^2} + \frac{1}{n_-} \sum_{j=1}^{n_-} \frac{\kappa e^{\kappa(N_j-\tau)}}{(e^{\kappa(N_j-\tau)}+1)^2} \\ y &= \frac{1}{n_+} \sum_{i=1}^{n_+} \frac{\kappa e^{\kappa(P_i-1)}}{(e^{\kappa(P_i-1)}+1)^2} \\ z &= \frac{1}{n_-} \sum_{j=1}^{n_-} \frac{\kappa e^{\kappa(N_j-1)}}{(e^{\kappa(N_j-1)}+1)^2} \end{aligned} \quad (35)$$

Then we reformulate the hessian matrix

$$H = M \begin{bmatrix} x & M'/M & M'/M \\ M'/M & y & 0 \\ M'/M & 0 & z \end{bmatrix} \quad (36)$$

where

$$x \in \left[ \frac{3\kappa e^\kappa}{(e^\kappa + 1)^2}, \frac{3\kappa}{4} \right], y, z \in \left[ \frac{\kappa e^\kappa}{(e^\kappa + 1)^2}, \frac{\kappa}{4} \right] \quad (37)$$

We calculate the principal minor of the hessian matrix

$$D_1 = x > 0 \quad (38)$$

$$D_2 = xy - \frac{M'^2}{M^2} > 0 \Rightarrow M'^2 \leq M^2 \frac{3\kappa^2 e^{2\kappa}}{(e^\kappa + 1)^4} < M^2 xy \quad (39)$$

$$D_3 = xyz - (y+z) \frac{M'^2}{M^2} > 0 \Rightarrow M'^2 \leq M^2 \frac{6\kappa^2 e^{3\kappa}}{(e^\kappa + 1)^6} < \frac{M^2 xyz}{y+z} \quad (40)$$

Hence, we find that if we have  $M'^2 < \min \left( M^2 \frac{3\kappa^2 e^{2\kappa}}{(e^\kappa + 1)^4}, M^2 \frac{6\kappa^2 e^{3\kappa}}{(e^\kappa + 1)^6} \right)$ , then we can ensure  $\tau$  is strongly convex for Eq.(33). We define the approximated equivalent formulation

$$\begin{aligned} \min_{\tau, \mathbf{P} \in \mathbb{R}^{n_+}, \mathbf{N} \in \mathbb{R}^{n_-}} \hat{\mathcal{L}}_{eq}(\mathbf{w}, \tau, c) &:= c \cdot \pi \cdot \left(1 - \frac{1}{n_+} \sum_{i=1}^{n_+} P_i\right) + (1-c) \cdot (1-\pi) \cdot \left(\frac{1}{n_-} \sum_{j=1}^{n_-} N_j\right) \\ &+ \frac{1}{n_+} \sum_{i=1}^{n_+} M' \psi(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) - P_i(s(\mathbf{w}, \mathbf{x}_i^+) - \tau) + M\psi(P_i - 1) + M\psi(\tau - 1) \\ &+ \frac{1}{n_-} \sum_{j=1}^{n_-} M' \psi(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) - N_j(s(\mathbf{w}, \mathbf{x}_j^-) - \tau) + M\psi(N_j - 1) \quad \forall i, j \quad 0 \leq \tau, P_i, N_j \end{aligned} \quad (41)$$

Then we can calculate the strong convexity of  $\tau$ . According to the definition of strong convex

$$\exists \mu > 0, \forall \tau \in [0, 1], \mathbf{P} \in [0, 1]^{n_+}, \mathbf{N} \in [0, 1]^{n_-} \quad \nabla^2 \hat{\mathcal{L}}_{eq} \succeq \mu_g I \quad (42)$$

Assuming that  $\mu_g > 0$ , in order to satisfy the strong convexity, we need to ensure the positive definiteness of the Hessian matrix

$$H = M \begin{bmatrix} x - \mu_g & M'/M & M'/M \\ M'/M & y - \mu_g & 0 \\ M'/M & 0 & z - \mu_g \end{bmatrix} \quad (43)$$



where

$$x \in \left[ \frac{3\kappa e^\kappa}{(e^\kappa + 1)^2}, \frac{3\kappa}{4} \right], y, z \in \left[ \frac{\kappa e^\kappa}{(e^\kappa + 1)^2}, \frac{\kappa}{4} \right] \quad (44)$$

We calculate the principal minor of the Hessian matrix

$$D_1 = x - \mu_g > 0 \Rightarrow \mu_g < \frac{3\kappa e^\kappa}{(e^\kappa + 1)^2} \quad (45)$$

$$D_2 = (x - \mu_g)(y - \mu_g) - \frac{M'^2}{M^2} > 0 \Rightarrow \mu_g > \frac{\kappa \sqrt{\kappa^2 + \frac{3\kappa^2}{16} M'^2 / M^2}}{2} \geq \frac{x + y + \sqrt{(x + y)^2 + 4xy M'^2 / M^2}}{2} \quad (46)$$

$$D_3 = (x - \mu_g)(y - \mu_g)(z - \mu_g) - (y + z - 2\mu_g) \frac{M'^2}{M^2} > 0 \Rightarrow \mu_g > \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} + \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} \quad (47)$$

where

$$p = -xy - xz - yz - \frac{(x + y + z)^2}{3} - \frac{M'(y + z)}{M} \quad (48)$$

$$q = xyz + \frac{2(x + y + z)^3}{27} - \frac{(9x + 9y + 9z) \left( -xy - xz - yz - \frac{M'(y + z)}{M} \right)}{27}$$

When  $p = -\frac{23\kappa^2}{24} - \frac{M'\kappa}{2M}$  and  $q = \frac{331\kappa^3}{1728} - \frac{5\kappa \left( -\frac{7\kappa^2}{16} - \frac{M'\kappa}{2M} \right)}{12}$ ,  $\mu$  has a lower bound. Hence, we find that if we have

$$\max \left( \sqrt[3]{-\frac{q}{2} + \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}} + \sqrt[3]{-\frac{q}{2} - \sqrt{\frac{q^2}{4} + \frac{p^3}{27}}}, \frac{\kappa \sqrt{\kappa^2 + \frac{3\kappa^2}{16} M'^2 / M^2}}{2} \right) < \mu_g < \min \left( \frac{\kappa e^\kappa}{(e^\kappa + 1)^2}, \frac{3\kappa e^\kappa}{(e^\kappa + 1)^2} \right) \quad (49)$$

then we can ensure  $\tau$  is  $\mu_g$ -strongly convex for Eq.(33). For computational simplicity, we use the upper bound of  $\mu = \frac{3M\kappa e^\kappa}{(e^\kappa + 1)^2} \geq M \cdot \mu_g$ .  $\square$