

---

# Generalizing Consistent Multi-Class Classification with Rejection to be Compatible with Arbitrary Losses

---

Yuzhou Cao<sup>1\*</sup> Tianchi Cai<sup>2</sup> Lei Feng<sup>3,4†</sup> Lihong Gu<sup>2</sup>  
Jinjie Gu<sup>2</sup> Bo An<sup>1</sup> Gang Niu<sup>4</sup> Masashi Sugiyama<sup>4,5</sup>

<sup>1</sup>School of Computer Science and Engineering, Nanyang Technological University, Singapore

<sup>2</sup>Ant Group, China

<sup>3</sup>College of Computer Science, Chongqing University, China

<sup>4</sup>RIKEN Center for Advanced Intelligence Project, Japan

<sup>5</sup>The University of Tokyo, Japan

## Abstract

*Classification with rejection* (CwR) refrains from making a prediction to avoid critical misclassification when encountering test samples that are difficult to classify. Though previous methods for CwR have been provided with theoretical guarantees, they are only compatible with certain loss functions, making them not flexible enough when the loss needs to be changed with the dataset in practice. In this paper, we derive a novel formulation for CwR that can be equipped with arbitrary loss functions while maintaining the theoretical guarantees. First, we show that  $K$ -class CwR is equivalent to a  $(K + 1)$ -class classification problem on the original data distribution with an augmented class, and propose an empirical risk minimization formulation to solve this problem with an estimation error bound. Then, we find necessary and sufficient conditions for the learning *consistency* of the surrogates constructed on our proposed formulation equipped with any classification-calibrated multi-class losses, where consistency means the surrogate risk minimization implies the target risk minimization for CwR. Finally, experimental results validate the effectiveness of our proposed method.

## 1 Introduction

In risk-sensitive multi-class classification applications (e.g., medical diagnosis, healthcare, autonomous driving, and product inspections [13, 22, 44]), misclassification can cause serious or even fatal consequences. To alleviate this issue, many studies have been conducted on *classification with rejection* (CwR) [11, 6, 63, 13, 14, 16, 22, 52, 48, 44, 8], which can abstain from making an unsure prediction to prevent such critical misclassification.

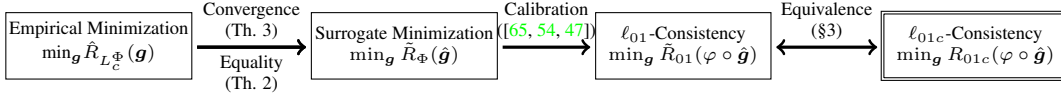
Most of the previous studies follow the framework that provides the reject option with a pre-defined cost  $c$  which is lower than the misclassification cost 1. Given cost  $c$ , the problem is further formulated as a risk minimization problem that aims to minimize the expectation of the zero-one- $c$  loss, i.e., the zero-one- $c$  risk. With the risk minimization process, the obtained classifier can balance the cost of rejection and prediction by choosing to incur a rejection cost  $c$  if the misclassification risk is high.

Due to the discontinuous nature of the zero-one- $c$  loss, recent works focused on finding its continuous surrogates to make the optimization problem tractable. A basic requirement for surrogate losses is the *consistency* [65, 7, 54, 47], i.e., the surrogate risk minimization implies the zero-one- $c$  risk minimization. Moreover, compared with the traditional  $K$ -class classification task where decisions

---

\*Work done when Yuzhou Cao was a research intern at Ant Group.

†Corresponding author: Lei Feng <lfeng@cqu.edu.cn>



**Figure 1:** Overview of the construction of consistent surrogates for classification with rejection in this work.

are normally made from the index of the maximum coordinate of a  $K$ -dimensional scoring function, the design of decision criteria in the CwR task is more elusive due to the existence of a reject option. By adopting different classification and rejection criteria, various surrogates of the zero-one- $c$  loss have been proposed with consistency analyses [6, 63, 13, 14, 48, 44, 8].

Classical studies focused on developing *confidence-based methods* [6, 63, 48, 44], which use the outputs of classifiers as confidence values and set a real-valued threshold as the rejection rule. Representative methods [63, 44] used surrogates that depend on *class-posterior probability estimation* (CPE) [50, 58], which is challenging when using deep models [25]. Though some of them [6, 48, 35, 24] could avoid CPE, most of them applied the modification of non-differentiable hinge/ramp-like surrogates, and their performance was only validated with linear models.

To avoid the use of the confidence threshold, Cortes et al. [13] provided an upper bound of the zero-one- $c$  loss as the surrogate that allows the use of a separated rejector and can be trained simultaneously with the classifier, which is regarded as *classifier-rejector methods*. Though these methods achieved state-of-the-art performance in binary classification scenarios, they only provided a consistency guarantee for hinge-like and exponential losses and cannot be directly generalized to the multi-class scenario as shown in Ni et al. [44]. Charoenphakdee et al. [8] showed that  $K$ -class CwR can be decomposed into  $K$  binary cost-sensitive classification problems [17, 51, 12] and proposed a family of surrogates are the ensembles of arbitrary binary classification losses, which can avoid CPE and the use of confidence threshold with properly chosen losses when the cost function is constant. Mozannar and Sontag [41] provided a modified version of the cross entropy loss as the surrogate for the task of learning to defer [41, 42] that can also be used in CwR, while its optimal solution still relies on CPE. In summary, previous works only took limited types of losses into consideration, and there lacks a theoretically grounded framework that can cover all the surrogates used in multi-class classification.

In this paper, we propose a novel framework for CwR that allows the use of arbitrary surrogate losses used in traditional multi-class classification as long as they are classification-calibrated, including but not limited to the well-known cross entropy loss, mean absolute error, focal loss [33, 9], and the pairwise/one-versus-all generalizations of binary margin losses [65]. Thanks to the flexible choices of losses, we are free of the restricted analyses on the consistency of certain surrogates. The access to a full range of surrogates also make it possible to custom the loss contingent on the characteristics and actual requirements of different datasets and tasks. An overview of our framework is shown in Figure 1. We summarize the main contributions of this work as follows:

- We disclose the equivalence between  $K$ -class CwR and a  $(K+1)$ -class classification problem on the original data distribution with an augmented class, by showing the equality between their classification risks.
- We propose a formulation of surrogates for  $\ell_{01c}$  that can recover the surrogate risk of a  $(K+1)$ -class classification task only with the  $K$ -class training distribution, and derive an estimation error bound for its empirical risk minimization.
- We find a sufficient condition for the consistency of the proposed family of surrogates *w.r.t.* the zero-one- $c$  loss that allows the use of any calibrated multi-class surrogates and show its necessity when considering a large family of loss functions.
- We for the first time provide an analysis on the calibration of the *generalized cross entropy* loss [66] that benefits from both the cross entropy loss and mean absolute error, and experimentally demonstrate that it is suitable for our proposed framework.

## 2 Preliminaries

In this section, we provide preliminary knowledge of CwR and calibrated surrogate losses, and discuss the consistency in CwR.

## 2.1 Classification with Rejection

The problem setting of CwR is based on the cost-based framework [11]. Let us denote by  $\mathcal{X}$  the feature space,  $\mathcal{Y} = \{1, 2, \dots, K\}$  the label space, and  $\mathcal{Y}^\circledast = \{1, 2, \dots, K, \circledast\}$  the label space with a reject option. We are given instance-label pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  independently and identically drawn from an underlying distribution with probability density  $p(\mathbf{x}, y)$ . The goal of CwR is to train a classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}^\circledast$  that can abstain from making a decision, where  $\circledast$  denotes the reject option. The evaluation metric of this task is the zero-one- $c$  loss  $\ell_{01c}$ , which can be expressed as a variant of the traditional zero-one loss  $\ell_{01}(f(\mathbf{x}), y) = \mathbb{I}[f(\mathbf{x}) \neq y]$ :

$$\ell_{01c}(f(\mathbf{x}), y) = \begin{cases} c, & f(\mathbf{x}) = \circledast, \\ \mathbb{I}[f(\mathbf{x}) \neq y], & f(\mathbf{x}) \in \{1, 2, \dots, k\}, \end{cases}$$

where  $\mathbb{I}[\cdot]$  is the Iverson bracket notation as suggested by Knuth [29] and the cost  $c$  can be further extended to an instance-dependent function  $c(\mathbf{x})$ . Our goal is to train a classifier that can minimize the expectation of  $\ell_{01c}$  over the data distribution:

$$R_{01c}(f) = \mathbb{E}_{p(\mathbf{x}, y)}[\ell_{01c}(f(\mathbf{x}), y)]. \quad (1)$$

Let us denote by  $f^* = \operatorname{argmin}_f R_{01c}(f)$  the Bayes optimal solution and  $\boldsymbol{\eta}(\mathbf{x}) = \{p(y|\mathbf{x})\}_{y=1}^K$  the posterior probabilities. When evaluated by  $\ell_{01c}$ , a classifier receives a standard classification error in  $\{0, 1\}$  if it makes a prediction and a cost of  $c$  if it does not make a prediction (i.e., chooses the reject option). Intuitively, an optimal solution  $f^*$  should balance the possibility of misclassification and the rejection cost  $c$ . This explanation is theoretically justified by Chow's rule [11]:

**Definition 1.** (Chow's Rule) A classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}^\circledast$  is the optimal solution of (1) if and only if it meets the following condition almost surely:

$$f(\mathbf{x}) = \begin{cases} \circledast, & \max_y \eta_y(\mathbf{x}) \leq 1 - c, \\ \operatorname{argmax}_y \eta_y(\mathbf{x}), & \text{else.} \end{cases}$$

Chow's rule shows that the optimal solution should refrain from making a decision if the most confident prediction of an example is still not confident enough given a rejection cost  $c$ .

## 2.2 Calibrated Surrogate Losses

Most classification problems can be formalized as the minimization of the target risk, which is the expectation of a target loss. Then *empirical risk minimization* (ERM) is conducted to obtain models with performance guarantees. However, most of the target losses are discontinuous, e.g., the zero-one loss in multi-class classification and the Hamming/ranking loss in multi-label classification [20]. Therefore, directly optimizing them is usually difficult and even NP-hard [18].

In order to optimize the target risk efficiently, surrogate risk minimization is preferred that minimizing the expectation of a continuous surrogate loss instead, e.g., the hinge loss in binary classification and the cross entropy loss in multi-class classification. For the statistical consistency of learning, *calibration* [53] is considered as a basic requirement for surrogate losses, which is a pointwise version of consistency and means that the minimization of the surrogate loss yields that of the target loss for each possible sample. A commonly adopted definition of the calibration of surrogates in multi-class classification is given as follows:

**Definition 2.** ( $\ell_{01}$ -Calibration [7, 54, 47]) For a  $K$ -class classification problem with target loss  $\ell_{01}$ , we say  $\Phi : \mathbb{R}^K \times \mathcal{Y} \rightarrow \mathbb{R}_+$  is  $\ell_{01}$ -calibrated if for any  $\mathbf{p} \in \Delta^K$ :

$$\inf_{\mathbf{u} \in \mathbb{R}^K, \mathbf{u} \notin \operatorname{argmin}_{\mathbf{u}} \mathbf{p}^T \mathbf{L}_{01}(\mathbf{u})} \mathbf{p}^T \Phi(\mathbf{u}) > \inf_{\mathbf{u} \in \mathbb{R}^K} \mathbf{p}^T \Phi(\mathbf{u}),$$

where  $\Phi(\mathbf{u}) = \{\Phi(\mathbf{u}, y)\}_{y=1}^K$ ,  $\mathbf{L}_{01}(\mathbf{u}) = \{\ell_{01}(\operatorname{argmax}_{y' \in \mathcal{Y}} \mathbf{u}_{y'}, y)\}_{y=1}^K$ .

The definition of  $\ell_{01}$ -calibration requires that a surrogate loss should be able to distinguish between optimal solutions and non-optimal ones *w.r.t.* any potential posterior distribution  $\mathbf{p}$ . This property is shown to be a necessary and sufficient condition for the statistical consistency of surrogate risk minimization, and fruitful research on the verification of  $\ell_{01}$ -calibrated surrogates has been conducted [7, 65, 54, 47, 46, 19].

**Table 1:** Comparisons between our proposed method and previous works of multi-class classification with rejection. Since our method is induced from a  $(K+1)$ -class classification problem, we can render a consistent learning guarantee with arbitrary surrogate losses that are calibrated *w.r.t.* the zero-one loss. Thanks to the abundant choices of losses, our proposed method can avoid CPE and the use of confidence thresholds.

Method	CPE-Free	Instance-Dependent Cost	Confidence Threshold-Free	Arbitrary Losses
[48]	✓	✓	✗	✗
[44]	✗	✓	✗	✗
[41]	✗	✓	✓	✗
[8]	✓	✗	✓	✗
Proposed	✓	✓	✓	✓

Besides multi-class classification, the calibration of surrogate losses also has been studied in various aspects of statistical learning, including but not limited to, multi-label classification [20, 64, 30, 59], AUC optimization [21, 39], general linear-fractional utility maximization [3], cost-sensitive learning [12, 51], top- $K$  classification [32, 62], and adversarially robust classification [4, 2, 1].

### 2.3 Consistency in Classification with Rejection

In the field of CwR, we are also interested in the consistency of surrogate losses. Let  $\mathcal{C} \subset \mathbb{R}^d$  where  $d \in \mathbb{N}$  and  $\Phi : \mathcal{C} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  is a surrogate loss, the consistency is defined as follows:

**Definition 3.** ( $\ell_{01c}$ -Consistency) A surrogate loss  $\Phi : \mathcal{C} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  is  $\ell_{01c}$ -consistent if there exists a function  $\varphi : \mathcal{C} \rightarrow \mathcal{Y}^{\otimes}$  for all probability distributions and all the sequences of functions  $\{\mathbf{g}_i\}_{i \in \mathbb{N}} : \mathcal{X} \rightarrow \mathcal{C}$ :

$$R_{\Phi}(\mathbf{g}_i) \rightarrow R_{\Phi}^* \Rightarrow R_{01c}(\varphi \circ \mathbf{g}_i) \rightarrow R_{01c}^*, \quad (2)$$

where  $R_{\Phi}(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)}[\Phi(\mathbf{g}(\mathbf{x}), y)]$ ,  $R_{\Phi}^* = \inf_{\mathcal{X} \rightarrow \mathcal{C}} R_{\Phi}(\mathbf{g})$ , and  $R_{01c}^* = \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}^{\otimes}} R_{01c}(f)$ .

This definition is inspired by the problem of general multi-class classification [47]. For an  $\ell_{01c}$ -consistent surrogate loss  $\Phi$ , we can safely minimize the surrogate risk  $R_{\Phi}$  instead while retaining the consistency guarantee of  $R_{01c}$ .

To ensure the consistency of  $\Phi$ , it is routine to discuss the calibration of surrogate losses. However, unlike the classical multi-class classification problem, where  $\varphi$  is usually an argmax operator, the design of  $\varphi$  for CwR can be quite complicated and hard to be unified, which makes it difficult to directly do calibration analysis on  $\Phi$ . The flexibility of  $\varphi$  also limits the discussions to specific types of surrogate losses. In Ramaswamy et al. [48], the authors considered the multi-class extensions of the hinge-loss with a confidence threshold. Ni et al. [44] indicated that the confidence-based method is indispensable and only focuses on class probability estimation via surrogate risk minimization. Both of Mozannar and Sontag [41] and Charoenphakdee et al. [8] gave surrogate losses for the zero-one- $c$  loss that does not depend on the accurate estimation of the class probability, while Mozannar and Sontag [41] focused on a variant of the cross entropy loss and Charoenphakdee et al. [8] constructed calibrated surrogate losses with the ensemble of  $K$  calibrated losses for binary classification.

In this paper, instead of directly discussing the calibration of surrogate  $\Phi$ , we show that there is an equivalence between classical multi-class classification and CwR. Based on this equivalence, we show that it is sufficient for  $\Phi$  to be  $\ell_{01c}$ -consistent by letting it be a simple variant of **any** calibrated surrogate loss *w.r.t.* the traditional zero-one loss  $\ell_{01}$ . The comparison of the proposed method and related works is shown in Table 1.

## 3 Equivalence between Classification with Rejection and Ordinary Classification

In this section, we first show that the risk  $R_{01c}(f)$  can be formalized as a  $(K+1)$ -class classification problem, and show that we can obtain  $\ell_{01c}$ -consistent surrogates with a variant of any calibrated surrogate *w.r.t.*  $\ell_{01}$ , which enables the use of  $\mathcal{C} \subset \mathbb{R}^{K+1}$  and  $\varphi(\cdot) = \text{argmax}(\cdot)$  as in the traditional multi-class classification tasks. We also show that such equivalence also holds when the cost  $c$

depends on sample  $\mathbf{x}$ . The proof of the conclusions in this section can be found in Appendix A. We start by considering the following distribution  $\mathcal{D}_c^\circledast$  over  $\mathcal{X} \times \mathcal{Y}^\circledast$  with probability density  $\tilde{p}(\mathbf{x}, \tilde{y})$ :

**Definition 4.** (Self-Augmented Distribution) A distribution  $\mathcal{D}_c^\circledast$  is called a  $c$ -self-augmented distribution *w.r.t.*  $\mathcal{D}$  if its probability density meets the following conditions:

$$\tilde{p}(\mathbf{x}, \tilde{y}) = \begin{cases} \frac{p(\mathbf{x}, y)}{2-c}, & \tilde{y} \in \{1, 2, \dots, K\}, \\ \frac{(1-c)p(\mathbf{x})}{2-c}, & \tilde{y} = \circledast. \end{cases}$$

It can be seen that distribution  $\mathcal{D}_c^\circledast$  shares the same marginal density of  $\mathbf{x}$  as the original distribution  $\mathcal{D}$  while  $\mathcal{D}_c^\circledast$  has an augmented class  $\circledast$  with class probability determined by the rejection cost  $c$ . Based on the connection between  $\mathcal{D}_c^\circledast$  and  $\mathcal{D}$ , we can further explore the relation between the two tasks: classification on  $\mathcal{D}_c^\circledast$  and CwR on  $\mathcal{D}$ .

**Theorem 1.** For any classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}^\circledast$ , the following equation holds:

$$R_{01c}(f) - R_{01c}^* = (2 - c) \left( \tilde{R}_{01}(f) - \tilde{R}_{01}^* \right),$$

where  $\tilde{R}_{01}(f) = \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})}[\ell_{01}(f(\mathbf{x}), \tilde{y})]$  and  $\tilde{R}_{01}^* = \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}^\circledast} \tilde{R}_{01}(f)$ .

Theorem 1 reveals the equivalence between the two tasks in a straightforward manner. Since the multiplication of the classification risk on  $\mathcal{D}_c^\circledast$  with a positive constant is equal to  $R_{01c}(f)$ , the minimization of  $\tilde{R}_{01}(f)$  immediately yields the minimization of  $R_{01c}(f)$  and vice versa. Furthermore, according to the linear correlation between  $\tilde{R}_{01}(f)$  and  $R_{01c}(f)$ , we can directly quantify the excess error  $R_{01c}(f) - R_{01c}^*$  by bounding  $\tilde{R}_{01}(f) - \tilde{R}_{01}^*$ , which is an easier work thanks to the existing research of multi-class classification. In conclusion, risk minimization with  $\tilde{R}_{01}(f)$  can also give a classifier with a rejection option with the optimality guarantee, and then we can consider a surrogate risk minimization problem for multi-class classification instead of CwR.

When the cost  $c(\mathbf{x})$  is an instance-dependent function, we show that such equivalence still holds with a minor modification. Considering the reweighted zero-one loss:  $\bar{\ell}_{01}(f(\mathbf{x}), y) = (2 - c(\mathbf{x}))\mathbb{1}[f(\mathbf{x}) \neq y]$  and its expectation  $\bar{R}_{01}(f)$  on  $\mathcal{D}_c^\circledast$ , we have the following conclusion:

**Corollary 1.** For any classifier  $f : \mathcal{X} \rightarrow \mathcal{Y}^\circledast$ , the following equation holds:

$$\bar{R}_{01}(f) - \bar{R}_{01}^* = R_{01c}(f) - R_{01c}^*.$$

It is obvious that Theorem 1 is a special case of Corollary 1 with constant cost functions. Though here we consider a reweighted classification task, the calibration result of multi-class surrogate losses can still be applied without any modification since the minimization of  $\bar{R}_{01}(f)$  can be seen as ordinary classification risk minimization with a slightly different marginal density  $p'(\mathbf{x})$ , which does not affect the calibration result since the class-posterior possibilities remain unchanged. All the conclusions in the rest of this paper can be extended to the scenario of instance-dependent cost (see Appendix G).

## 4 $\ell_{01c}$ -Consistent Surrogates with Arbitrary $\ell_{01}$ -Calibrated Losses

As discussed in Section 3, CwR can be safely replaced by multi-class classification on a special distribution  $\mathcal{D}_c^\circledast$ . Following the surrogate risk minimization in multi-class classification, we can replace the zero-one loss  $\ell_{01}$  with a surrogate risk  $\Phi : \mathbb{R}^{K+1} \times \mathcal{Y} \cup \{K+1\} \rightarrow \mathbb{R}_+$  and minimizing the surrogate risk with a score-based classifier  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^{K+1}$  instead, which is defined as follows:

$$\tilde{R}_\Phi(\mathbf{g}) = \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})}[\Phi(\mathbf{g}(\mathbf{x}), t(\tilde{y}))], \quad (3)$$

where  $t(\tilde{y}) = K + 1$  if  $\tilde{y} = \circledast$  and  $t(\tilde{y}) = \tilde{y}$  otherwise. (3) is a typical formulation of the multi-class classification risk and we can asymptotically minimize it following the ERM framework [55]. After the risk minimization process, the prediction is generated with the following link function  $\varphi : \mathbb{R}^{K+1} \rightarrow \mathcal{Y}^\circledast$ :

$$\varphi(\mathbf{u}) = \begin{cases} \circledast, & \operatorname{argmax}_{y \in \mathcal{Y} \cup \{K+1\}} \mathbf{u}_y(\mathbf{x}) = K + 1, \\ \operatorname{argmax}_{y \in \mathcal{Y} \cup \{K+1\}} \mathbf{u}_y(\mathbf{x}), & \text{else.} \end{cases}$$

With a classification-calibrated surrogate  $\Phi$ , it is sufficient to say that the minimization of  $\tilde{R}_\Phi(\mathbf{g})$  can lead to that of  $\tilde{R}_{01c}(\varphi(\mathbf{g}))$ , which indicates the minimization of  $R_{01c}(\varphi(\mathbf{g}))$  according to Theorem 1 and Corollary 1. The theory of how to find such surrogates has been thoroughly studied in the field of the classification-calibration of multi-class surrogates [65, 54, 47].

However, we do not have direct access toward  $\mathcal{D}_c^\circledast$  though it is closely related to the available data distribution  $\mathcal{D}$ . In this section, we propose a family of surrogate losses based on the conclusions in the previous section, which allows the use of any multi-class classification surrogates. With this formulation of surrogates, we can recover the classification risk of  $\tilde{R}_\Phi(\mathbf{g})$  without access to  $\mathcal{D}_c^\circledast$  by taking its expectation in  $\mathcal{D}$ . Based on the loss formulation, we also provide the estimation error bound to show the validity of ERM.

#### 4.1 Formulation of Surrogates

Here, we begin with the definition of a family of surrogates for the zero-one- $c$  loss, and then show how it can relate  $\mathcal{D}$  and  $\mathcal{D}_c^\circledast$ . With any multi-class classification loss  $\Phi$ , we have the following formulation of surrogate losses for  $\ell_{01c}$ :

**Definition 5.** Given a pre-defined rejection cost  $c$ , we have the following formulation of surrogate  $L_c^\Phi : \mathbb{R}^{K+1} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  for CwR:

$$L_c^\Phi(\mathbf{u}, y) = \Phi(\mathbf{u}, y) + (1 - c)\Phi(\mathbf{u}, K + 1), \quad (4)$$

where  $\Phi : \mathbb{R}^{K+1} \times \mathcal{Y} \cup \{K + 1\} \rightarrow \mathbb{R}_+$  and  $\mathbf{u} \in \mathbb{R}^{K+1}$ .

The proposed surrogate loss is the linear combination of a  $(K + 1)$  dimensional multi-class classification loss with coefficient determined by the predefined cost  $c$ . It can also be learned from Appendix A.2 of Charoenphakdee et al. [8] that when  $\Phi$  is the softmax cross entropy loss, (4) is equivalent to Mozannar and Sontag [41]. The following theorem reveals the connection between  $\tilde{R}_\Phi(\mathbf{g})$  and the expectation of  $L_c^\Phi$  on  $\mathcal{D}$ :

**Theorem 2.** For any  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^{K+1}$  and  $R_{L_c^\Phi}(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)}[L_c^\Phi(\mathbf{g}(\mathbf{x}), y)]$ :

$$R_{L_c^\Phi}(\mathbf{g}) = (2 - c)\tilde{R}_\Phi(\mathbf{g}).$$

The proof is provided in Appendix B. From Theorem 2, we can obtain the risk  $\tilde{R}_\Phi(\mathbf{g})$  without access to  $\mathcal{D}_c^\circledast$  with the use of the proposed surrogate (4). Following the common practice, we can finally conduct ERM [55] that minimizes the unbiased estimator of  $R_{L_c^\Phi}(\mathbf{g})$ , which is also that of  $\tilde{R}_\Phi(\mathbf{g})$  according to Theorem 2:

$$\hat{R}_{L_c^\Phi}(\mathbf{g}) = \frac{1}{n} \sum_{i=1}^n L_c^\Phi(\mathbf{g}(\mathbf{x}_i), y_i) \quad (5)$$

After minimizing  $\hat{R}_{L_c^\Phi}(\mathbf{g})$  and obtaining the empirically optimal  $\hat{\mathbf{g}}$ , we can use it for predicting with the link function  $\varphi \circ \hat{\mathbf{g}}$ , where  $\varphi \circ \hat{\mathbf{g}}(\mathbf{x}) = \varphi(\hat{\mathbf{g}}(\mathbf{x}))$ .

According to the unbiasedness of (3), it is promising that the induced prediction rule  $\varphi \circ \hat{\mathbf{g}}$  can approximate Chow's rule (Definition 1). To quantify such approximation, there remains two questions: what is the relation between the minimization of empirical risk  $\hat{R}_{L_c^\Phi}(\mathbf{g})$  and  $R_{L_c^\Phi}(\mathbf{g})$ , and whether the minimization of  $R_{L_c^\Phi}(\mathbf{g})$  yields that of  $R_{01c}(\varphi \circ \mathbf{g})$ . We will answer the two problems in Section 4.2 and Section 5, respectively.

#### 4.2 Estimation Error Bound

In Section 4.1, we proposed a family of surrogates that can recover the surrogate risk on  $\mathcal{D}_c^\circledast$  with only  $\mathcal{D}$  and provided an ERM framework to learn the empirically optimal  $\hat{\mathbf{g}} = \min_{\mathbf{g} \in \mathcal{G}} \hat{R}_{L_c^\Phi}(\mathbf{g})$ . Here we further justify the use of ERM by showing that the minimization of  $\hat{R}_{L_c^\Phi}$  can also result in that of  $R_{L_c^\Phi}$  with the following estimation error bound.

**Theorem 3.** For any  $\delta \in (0, 1)$ , suppose the model class of  $g_y$  is  $\mathcal{G}_y$  and  $\mathbf{g} \in \mathcal{G}$ , where  $\mathcal{G}_y \subset \mathcal{X} \rightarrow \mathbb{R}$  and  $\mathcal{G} \subset \mathcal{X} \rightarrow \mathbb{R}^{K+1}$  is composed of  $\{\mathcal{G}_y\}_{y=1}^{K+1}$ .  $\Phi(\cdot, y)$  is  $\rho$ -Lipschitz continuous and is bounded

by  $C_\Phi > 0$ . Assume that the identifiable condition holds, i.e.,  $\min_{\mathbf{g} \in \mathcal{G}} R_{L_c^\Phi}(\mathbf{g}) = R_{L_c^\Phi}^*$ , then the following inequality holds with probability at least  $1 - \delta$ :

$$R_{L_c^\Phi}(\hat{\mathbf{g}}) - R_{L_c^\Phi}^* \leq 4\sqrt{2}(2-c)\rho \sum_{y=1}^{K+1} \mathfrak{R}_n(\mathcal{G}_y) + (2-c)C_\Phi \sqrt{\frac{2 \log 2/\delta}{n}}, \quad (6)$$

where  $\mathfrak{R}_n(\mathcal{G}_y)$  is the Rademacher complexity [5] w.r.t.  $\mathcal{G}_y$  on the distribution with density  $p(\mathbf{x})$  that often decays in the rate of  $O(\frac{1}{\sqrt{n}})$ .

We prove this conclusion in Appendix C. From the theorem above, we can learn that with the identifiable condition which is a common assumption with the use of complex models [4, 27, 34],  $R_{L_c^\Phi}(\hat{\mathbf{g}})$  converges to  $R_{L_c^\Phi}^*$  in  $O_p(1/\sqrt{n})$ , which is the optimal parametric convergence rate without additional assumptions [38]. According to Theorem 2, it is straightforward that  $\tilde{R}_\Phi(\mathbf{g}) \xrightarrow{P} \tilde{R}_\Phi^*$  also holds. Nevertheless, the relation between the minimization of surrogate risk  $\tilde{R}_\Phi(\mathbf{g})$  and that of the target risk  $\tilde{R}_{01}(\varphi \circ \mathbf{g})$  is still unknown. According to Theorem 1, the minimization of  $\tilde{R}_{01}(\varphi \circ \mathbf{g})$  is equivalent to zero-one- $c$  risk minimization, which is the goal of CwR. We answer this question in the next section by giving a sufficient condition for the  $\ell_{01c}$ -consistency for  $L_c^\Phi$ .

## 5 Theoretical Analysis

In this section, we first point out the sufficient condition for  $L_c^\Phi$  to be  $\ell_{01c}$ -calibrated and show that it is also necessary if the candidates of  $\Phi$  are permutation-invariant, which consists of a large number of commonly used loss functions, including but not limited to cross-entropy loss, focal loss, and mean absolute error. Then we further specify the regret transfer bounds for a family of CPE-free surrogates [66], which has not been provided with theoretical analysis before.

### 5.1 Necessary and Sufficient Conditions for $\ell_{01c}$ -Consistency

Given the loss formulation (4), a natural idea is to construct surrogate  $L_c^\Phi$  with commonly used multi-class loss functions. Here, we justify this idea by showing that we can borrow the calibration analyses of multi-class surrogates and set  $\Phi$  to any  $(K+1)$ -class  $\ell_{01}$ -calibrated surrogates according to the following conditions:

**Theorem 4.**  $L_c^\Phi$  is  $\ell_{01c}$ -consistent for any  $c \in [0, 1]$  if  $\Phi$  is an  $\ell_{01}$ -calibrated surrogate loss. Furthermore, if  $\Phi$  is permutation-invariant, i.e.,  $\Phi(P\mathbf{g}) = P\Phi(\mathbf{g})$  for all permutation matrices  $P$ ,  $\Phi$ 's  $\ell_{01}$ -calibration is also necessary for the  $\ell_{01c}$ -consistency of  $L_c^\Phi$ .

The complete proof is shown in Appendix D and here we provide its sketch. The equivalence between CwR and multi-class classification on  $\mathcal{D}_c^\otimes$  shown in Theorem 1, 2, and Corollary 1 directly yields the sufficiency of this condition. Though the equivalent classification problem is limited on  $\mathcal{D}_c^\otimes$ , the permutation-invariance of  $\Phi$  and the arbitrariness of  $c$  require that the minimizers of the expectation of  $\Phi$  should be included in that of  $\ell_{01}$  for any potential distributions, which indicates the necessity of  $\Phi$ 's calibration.

As a result, we can use any  $\Phi$  in an off-the-shelf manner, i.e., to the consistency of different  $L_c^\Phi$ , we only have to check if  $\Phi$  is  $\ell_{01}$ -calibrated, which has been studied thoroughly [7, 54, 47], instead of tedious case-based discussions. It is noticeable that when considering multi-class surrogates that are not permutation-invariant, the calibration of  $\Phi$  may not be necessary, which indicates that surrogates that are not calibrated and not permutation-invariant may also be potential options. Nevertheless, a classification-calibrated surrogate  $\Phi$  can always be a safe choice given the sufficiency in Theorem 4 and fruitful calibration analyses.

### 5.2 Calibration Result for Generalized Cross Entropy Loss

Given the sufficient condition for  $\ell_{01c}$ -consistency, we can construct  $L_c^\Phi$  with any  $\ell_{01}$ -calibrated surrogates. However, it has been shown in Charoenphakdee et al. [8] that it can lead to a model that rejects more data than necessary if the cross entropy (CE) loss is used as  $\Phi$ , which is a popular choice as a surrogate. Another common surrogate is the mean absolute error (MAE). Though it can avoid CPE and only focus on the crucial class with the maximum posterior probability, it usually takes more training epochs before convergence [66], which can be costly in practical use.

Here, we consider the *generalized cross entropy* (GCE) loss [66] that can take the advantages of the CE loss and MAE, which is defined as below:

**Definition 6.** (Generalized cross entropy losses) For any  $\gamma \in (0, 1]$ , the GCE loss is defined as below:

$$\Phi^\gamma(\mathbf{g}(\mathbf{x}), y) = (1 - S(\mathbf{g})_y^\gamma) / \gamma,$$

where  $S(\cdot)$  is the softmax-transformation.

It can be seen that the loss formulation is equivalent to MAE if  $\gamma = 1$  and it is also reported in Zhang and Sabuncu [66] that the GCE loss can approximate the CE loss if  $\gamma \rightarrow 0$ . Though the GCE loss has proved to be effective in practical use, to the best of our knowledge, its calibration results remain unknown. To justify the combination of GCE loss and  $L_c^\Phi$ , we give the calibration analysis and further show its analytical solution.

**Theorem 5.** The GCE loss  $\Phi^\gamma$  is  $\ell_{01}$ -calibrated for any  $\gamma \in (0, 1]$ . For the optimal model  $\mathbf{g}^*$ ,  $S(\mathbf{g}^*)_y = \eta_y^{\frac{1}{1-\gamma}} / \sum_{y'=1}^K \eta_{y'}^{\frac{1}{1-\gamma}}$  for all the  $\mathbf{x} \in \mathcal{X}$  almost surely if  $\gamma \in (0, 1)$ . If  $\gamma = 1$ ,  $S(\mathbf{g}^*)_{\text{argmax}_y \eta_y} = 1$ .

The proof can be found in Appendix E. After verifying the calibration result of the GCE loss, we can combine it with the loss formulation  $L_c^\Phi$  and obtain an  $\ell_{01c}$ -consistent surrogate. We will experimentally demonstrate its effectiveness in the next section.

## 6 Experiments

In this section, we provide the experiment results of CwR with deep models, which are evaluated by the zero-one- $c$  loss following the common practice [44, 8]. We also show the misclassification rate of the accepted data and the ratio of the rejected data. Details of the setup and the experiments for instance-dependent cost can be found in Appendix F and G, respectively.

**Datasets and Models.** In the experiments, we evaluate the proposed methods and baselines on three widely-used benchmarks Fashion-MNIST [60], SVHN [43], CIFAR-10 [31] with cost  $c$  selected from  $\{0.05, 0.06, 0.07, 0.08, 0.09, 0.10\}$  for Fashion-MNIST and  $\{0.05, 0.10, 0.15, 0.20, 0.25, 0.30\}$  the other two. We conduct data augmentation for CIFAR-10 and use the original datasets of Fashion-MNIST and SVHN in the experiments. For Fashion-MNIST, we use a CNN defined in Charoenphakdee et al. [8], and ResNet-18 and ResNet-34 [26] are used for SVHN and CIFAR-10, respectively.

**Baselines.** We compare our method with state-of-the-art methods in CwR, including confidence-based cross entropy loss (CE) [44], learning to defer (DEFER) [41], and cost-sensitive learning-based method with sigmoid loss (CS) [8], in which DEFER is a special case of our method that use cross entropy loss as  $\Phi$ . For CE, we also conduct the temperature scaling [25] to alleviate overconfidence. For the proposed method, we use GCE with default parameter  $\gamma = 0.7$  as suggested in Zhang and Sabuncu [66] and pairwise-sigmoid (Sigmoid) loss [65] to construct the surrogate  $L_c^\Phi$ . We implemented all the methods by Pytorch [45], and conducted all the experiments on NVIDIA GeForce 3090 GPUs.

**Experimental Results.** As can be seen from the experimental results reported in Table 2, our proposed method (i.e., either GCE or Sigmoid) significantly outperforms other compared methods in most cases. Obviously, for all the datasets and cost  $c$ , our GCE method outperforms the baseline DEFER method, which indicates that CwR cannot be simply solved by the methods used for learning to defer. It can be also seen that confidence-based CE is only comparable to the proposed method on FMNIST with a simple CNN. When complex models are used, the effect of overconfidence is inevitable even with the use of temperature scaling, which can be induced from the fact that CE often rejects less data than GCE on SVHN and CIFAR-10. Though CS is comparable to GCE on CIFAR-10 when the rejection cost is high, its performance degrades drastically when the rejection cost decreases, which shows that it is not the best choice in highly error-critical tasks. When ResNet-18 and ResNet-34 are used on SVHN and CIFAR-10 respectively, our GCE method outperforms or is comparable to all the baselines, which shows that GCE is more stable on complex models. Our proposed Sigmoid method performs better than most baselines and is comparable to CE with the use of a simple CNN model, which aligns with the existing observations that pairwise losses are



**Table 2:** The mean and standard error of the zero-one- $c$  losses (**01c**, rescaled to 0-100), rejection ratio (**Rej**), and missclassification rates (**01**) of the accepted data for 5 trails. The best and comparable methods based on the paired t-test at the significance level 5% are highlighted in boldface.

Method	Cost	CE			CS			DEFER			GCE			Sigmoid		
		01c	Rej	01	01c	Rej	01	01c	Rej	01	01c	Rej	01	01c	Rej	01
FMNIST	0.05	2.30 (0.07)	25.17 (3.17)	1.39 (0.11)	2.93 (0.25)	34.95 (1.94)	1.81 (0.48)	3.79 (0.28)	50.461 (2.51)	2.58 (0.46)	3.22 (0.07)	50.47 (2.49)	1.39 (0.30)	<b>2.23</b> ( <b>0.01</b> )	30.98 (0.62)	0.99 (0.05)
	0.06	<b>2.58</b> ( <b>0.07</b> )	22.92 (1.45)	1.56 (0.09)	3.37 (0.15)	33.13 (1.27)	2.07 (0.28)	4.63 (0.10)	56.45 (3.69)	2.84 (0.42)	3.78 (0.17)	50.46 (1.24)	1.53 (0.30)	<b>2.62</b> ( <b>0.08</b> )	26.76 (3.37)	1.37 (0.21)
	0.07	<b>2.73</b> ( <b>0.14</b> )	21.17 (2.23)	1.58 (0.31)	3.45 (0.17)	35.77 (2.62)	1.47 (0.04)	5.18 (0.47)	56.46 (6.85)	2.86 (0.41)	4.23 (0.21)	48.05 (5.24)	1.66 (0.25)	2.94 (0.07)	29.87 (0.85)	1.21 (0.17)
	0.08	<b>3.12</b> ( <b>0.11</b> )	20.71 (1.68)	1.85 (0.07)	4.13 (0.36)	33.68 (0.32)	2.17 (0.52)	5.86 (0.30)	54.08 (3.47)	3.36 (0.29)	4.50 (0.06)	45.66 (2.36)	1.55 (0.25)	<b>3.14</b> ( <b>0.17</b> )	26.10 (0.23)	1.43 (0.25)
	0.09	<b>3.55</b> ( <b>0.21</b> )	23.64 (1.82)	1.86 (0.18)	4.20 (0.21)	31.90 (1.74)	1.96 (0.15)	6.31 (0.40)	54.62 (4.33)	3.09 (0.49)	4.95 (0.06)	44.05 (1.74)	1.77 (0.23)	<b>3.50</b> ( <b>0.05</b> )	23.71 (2.28)	1.79 (0.18)
	0.10	<b>3.59</b> ( <b>0.16</b> )	18.32 (1.56)	2.15 (0.32)	4.45 (0.20)	28.96 (0.13)	2.18 (0.41)	6.72 (0.07)	52.69 (0.74)	3.08 (0.18)	5.06 (0.23)	39.01 (4.87)	1.89 (0.26)	<b>3.73</b> ( <b>0.05</b> )	23.96 (1.90)	1.76 (0.20)
SVHN	0.05	3.33 (0.14)	14.37 (0.94)	3.05 (0.13)	4.42 (0.13)	12.81 (0.14)	4.33 (0.12)	4.19 (0.29)	33.05 (1.59)	3.80 (0.37)	<b>2.68</b> ( <b>0.17</b> )	19.79 (0.72)	2.10 (0.24)	<b>2.70</b> ( <b>0.14</b> )	29.56 (1.16)	1.73 (0.17)
	0.10	4.66 (0.20)	10.91 (0.57)	4.01 (0.21)	4.48 (0.14)	12.85 (0.42)	3.67 (0.11)	5.55 (0.56)	30.72 (2.64)	3.58 (0.52)	<b>4.13</b> ( <b>0.11</b> )	14.83 (0.54)	3.10 (0.10)	<b>4.13</b> ( <b>0.39</b> )	19.16 (1.94)	2.74 (0.43)
	0.15	5.40 (0.09)	8.52 (0.15)	4.50 (0.07)	5.14 (0.10)	13.21 (0.62)	3.64 (0.19)	6.37 (0.21)	21.19 (0.94)	4.05 (0.25)	<b>4.66</b> ( <b>0.06</b> )	11.47 (0.41)	3.31 (0.07)	<b>4.83</b> ( <b>0.44</b> )	18.38 (1.37)	2.54 (0.61)
	0.20	6.16 (0.13)	7.74 (0.26)	4.99 (0.09)	<b>5.51</b> ( <b>0.20</b> )	12.78 (1.03)	3.19 (0.24)	5.99 (0.17)	12.33 (0.51)	4.02 (0.16)	<b>5.44</b> ( <b>0.04</b> )	10.02 (0.25)	3.82 (0.03)	6.39 (0.45)	15.86 (0.72)	3.82 (0.48)
	0.25	7.08 (0.32)	6.51 (1.06)	5.83 (0.36)	6.77 (0.16)	12.96 (0.97)	4.06 (0.18)	6.69 (0.16)	9.18 (0.35)	4.33 (0.16)	<b>5.75</b> ( <b>0.14</b> )	8.64 (0.20)	3.93 (0.12)	6.74 (0.13)	13.79 (0.33)	3.82 (0.14)
	0.30	7.12 (0.16)	5.31 (0.36)	5.83 (0.18)	7.26 (0.33)	13.21 (1.20)	3.80 (0.41)	7.07 (0.31)	12.35 (2.31)	4.55 (0.34)	<b>6.30</b> ( <b>0.09</b> )	8.72 (0.11)	4.04 (0.09)	7.69 (0.22)	10.79 (0.76)	5.00 (0.13)
CIFAR-10	0.05	4.43 (0.23)	29.93 (1.85)	4.18 (0.33)	6.59 (0.27)	20.20 (0.51)	7.00 (0.35)	4.62 (0.47)	44.97 (5.24)	4.30 (0.88)	3.80 (0.20)	34.52 (2.77)	3.16 (0.35)	<b>3.67</b> ( <b>0.03</b> )	42.69 (8.74)	2.63 (0.49)
	0.10	7.13 (0.11)	21.13 (0.81)	6.35 (0.18)	7.68 (0.32)	20.31 (0.66)	7.08 (0.42)	6.56 (0.26)	26.21 (1.12)	5.34 (0.39)	<b>5.84</b> ( <b>0.12</b> )	25.47 (0.98)	4.41 (0.15)	<b>6.11</b> ( <b>0.13</b> )	31.66 (2.17)	4.30 (0.30)
	0.15	9.03 (0.32)	7.76 (0.39)	7.74 (0.37)	8.35 (0.29)	21.83 (0.92)	6.49 (0.45)	8.39 (0.19)	20.39 (1.59)	6.69 (0.35)	<b>7.56</b> ( <b>0.14</b> )	20.43 (0.60)	5.65 (0.23)	8.18 (0.10)	23.39 (0.82)	6.10 (0.18)
	0.20	10.45 (0.29)	14.53 (0.47)	8.82 (0.38)	<b>9.32</b> ( <b>0.21</b> )	21.86 (0.46)	6.33 (0.33)	9.65 (0.14)	17.16 (1.04)	7.50 (0.11)	<b>9.09</b> ( <b>0.14</b> )	18.45 (1.93)	6.62 (0.42)	9.69 (0.15)	19.54 (1.55)	7.20 (0.07)
	0.25	11.64 (0.26)	11.20 (0.30)	9.96 (0.32)	<b>10.46</b> ( <b>0.24</b> )	22.02 (0.40)	6.35 (0.35)	10.85 (0.08)	14.22 (1.35)	8.50 (0.30)	<b>10.31</b> ( <b>0.23</b> )	15.39 (1.47)	7.64 (0.38)	10.96 (0.11)	14.99 (1.71)	8.48 (0.40)
	0.30	12.20 (0.18)	10.02 (0.53)	10.89 (0.15)	<b>11.43</b> ( <b>0.23</b> )	22.23 (0.81)	6.13 (0.24)	11.90 (0.17)	11.48 (0.75)	9.55 (0.31)	<b>11.23</b> ( <b>0.16</b> )	12.52 (0.122)	8.55 (0.14)	12.14 (0.12)	11.08 (0.60)	9.91 (0.25)

often effective with simple models [57, 15]. These results show that our method can benefit from the flexibility of the choices of loss functions.

## 7 Conclusion

In this paper, we studied the problem of classification with rejection, which can refrain from making a prediction to avoid critical misclassification. We derived a novel formulation for CwR that can be equipped with arbitrary loss functions while maintaining the theoretical guarantees, making them highly adaptive to the dataset in practical use. First, we showed the equivalence between  $K$ -class CwR and a  $(K+1)$ -class classification problem, and proposed an empirical risk minimization formulation to solve this problem with an estimation error bound. Then, we pointed out necessary and sufficient conditions for the learning consistency of the surrogates constructed on our proposed formulation equipped with any classification-calibrated multi-class losses. Finally, experimental results demonstrated the effectiveness of our proposed method.

## Acknowledgement

This research is supported by the National Research Foundation, Singapore under its Industry Alignment Fund – Pre-positioning (IAF-PP) Funding Initiative. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore. Yuzhou Cao was also supported by Ant Group through Ant Research Intern Program. Lei Feng was supported by the National Natural Science Foundation of China (Grant No. 62106028), Chongqing Overseas Chinese Entrepreneurship and Innovation Support Program, and CAAI-Huawei MindSpore Open Fund. MS was supported by JST CREST Grant Number JPMJCR18A2.

## References

- [1] Pranjal Awasthi, Natalie Frank, Anqi Mao, Mehryar Mohri, and Yutao Zhong. Calibration and consistency of adversarial surrogate losses. *CoRR*, abs/2104.09658, 2021. URL <https://arxiv.org/abs/2104.09658>.
- [2] Pranjal Awasthi, Anqi Mao, Mehryar Mohri, and Yutao Zhong. A finer calibration analysis for adversarial robustness. *CoRR*, abs/2105.01550, 2021. URL <https://arxiv.org/abs/2105.01550>.
- [3] Han Bao and Masashi Sugiyama. Calibrated surrogate maximization of linear-fractional utility in binary classification. In *AISTATS*, volume 108 of *Proceedings of Machine Learning Research*, pages 2337–2347. PMLR, 2020.
- [4] Han Bao, Clayton Scott, and Masashi Sugiyama. Calibrated surrogate losses for adversarially robust classification. In *COLT*, volume 125 of *Proceedings of Machine Learning Research*, pages 408–451. PMLR, 2020.
- [5] Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *J. Mach. Learn. Res.*, 3:463–482, 2002. URL <http://jmlr.org/papers/v3/bartlett02a.html>.
- [6] Peter L. Bartlett and Marten H. Wegkamp. Classification with a reject option using a hinge loss. *J. Mach. Learn. Res.*, 9:1823–1840, 2008.
- [7] Peter L. Bartlett, Michael I. Jordan, and Jon D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- [8] Nontawat Charoenphakdee, Zhenghang Cui, Yivan Zhang, and Masashi Sugiyama. Classification with rejection based on cost-sensitive classification. In *ICML*, volume 139 of *Proceedings of Machine Learning Research*, pages 1507–1517. PMLR, 2021.
- [9] Nontawat Charoenphakdee, Jayakorn Vongkulbhisal, Nuttapong Chairatanakul, and Masashi Sugiyama. On focal loss for class-posterior probability estimation: A theoretical perspective. In *CVPR*, pages 5202–5211, 2021.
- [10] Mohammad-Amin Charusaie, Hussein Mozannar, David A. Sontag, and Samira Samadi. Sample efficient learning of predictors that complement humans. In *ICML*, volume 162, pages 2972–3005, 2022.
- [11] C. Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*, 16(1):41–46, 1970. doi: 10.1109/TIT.1970.1054406.
- [12] Scott Clayton. Calibrated asymmetric surrogate losses. *Electronic Journal of Stats*, 6:238–238, 2012.
- [13] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Boosting with abstention. In *NeurIPS*, pages 1660–1668, 2016.
- [14] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In *ALT*, volume 9925, pages 67–82, 2016.
- [15] Ürün Dogan, Tobias Glasmachers, and Christian Igel. A unified view on multi-class support vector classification. *J. Mach. Learn. Res.*, 17:45:1–45:32, 2016.
- [16] Ran El-Yaniv and Yair Wiener. On the foundations of noise-free selective classification. *J. Mach. Learn. Res.*, 11:1605–1641, 2010.
- [17] Charles Elkan. The foundations of cost-sensitive learning. In *IJCAI*, pages 973–978. Morgan Kaufmann, 2001.
- [18] Vitaly Feldman, Venkatesan Guruswami, Prasad Raghavendra, and Yi Wu. Agnostic learning of monomials by halfspaces is hard. *SIAM J. Comput.*, 41(6):1558–1590, 2012.

- [19] Jessica Finocchiaro, Rafael M. Frongillo, and Bo Waggoner. An embedding framework for consistent polyhedral surrogates. In *NeurIPS*, pages 10780–10790, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/9ec51f6eb240fb631a35864e13737bca-Abstract.html>.
- [20] Wei Gao and Zhi-Hua Zhou. On the consistency of multi-label learning. *Artif. Intell.*, 199–200: 22–44, 2013.
- [21] Wei Gao and Zhi-Hua Zhou. On the consistency of AUC pairwise optimization. In *IJCAI*, pages 939–945. AAAI Press, 2015.
- [22] Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In *NeurIPS*, pages 4878–4887, 2017.
- [23] Gustavo L. Gilardoni. On pinsker’s and vajda’s type inequalities for csiszár’s f-divergences. *IEEE Trans. Inf. Theory*, 56(11):5377–5386, 2010. doi: 10.1109/TIT.2010.2068710.
- [24] Yves Grandvalet, Alain Rakotomamonjy, Joseph Keshet, and Stéphane Canu. Support vector machines with a reject option. In *NeurIPS*, pages 537–544, 2008.
- [25] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 1321–1330. PMLR, 2017.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 463–469, 2016.
- [27] Yen-Chang Hsu, Zhaoyang Lv, Joel Schlosser, Phillip Odom, and Zsolt Kira. Multi-class classification without multi-class labels. In *ICLR*, 2019.
- [28] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, volume 37, pages 448–456. JMLR.org, 2015.
- [29] D. E. Knuth. Two notes on notation. *American Mathematical Monthly*, 99(5):403–422, 1992.
- [30] Oluwasanmi Koyejo, Nagarajan Natarajan, Pradeep Ravikumar, and Inderjit S. Dhillon. Consistent multilabel classification. In *NeurIPS*, pages 3321–3329, 2015.
- [31] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [32] Maksim Lapin, Matthias Hein, and Bernt Schiele. Analysis and optimization of loss functions for multiclass, top-k, and multilabel classification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 40(7):1533–1554, 2018.
- [33] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, pages 2999–3007, 2017.
- [34] Nan Lu, Shida Lei, Gang Niu, Issei Sato, and Masashi Sugiyama. Binary classification from multiple unlabeled datasets via surrogate set classification. In *ICML*, 2021.
- [35] Naresh Manwani, Kalpit Desai, Sanand Sasidharan, and Ramasubramanian Sundararajan. Double ramp loss based reject option classifier. In *PAKDD*, volume 9077, pages 151–163, 2015.
- [36] Andreas Maurer. A vector-contraction inequality for rademacher complexities. In Ronald Ortner, Hans Ulrich Simon, and Sandra Zilles, editors, *ALT*, pages 3–17, 2016.
- [37] C. McDiarmid. *On the method of bounded differences*. Surveys in Combinatorics, 1989.
- [38] Shahar Mendelson. Lower bounds for the empirical minimization algorithm. *IEEE Trans. Inf. Theory*, 54(8):3797–3803, 2008. doi: 10.1109/TIT.2008.926323. URL <https://doi.org/10.1109/TIT.2008.926323>.
- [39] Aditya Krishna Menon and Robert C. Williamson. In *COLT*, volume 35 of *JMLR Workshop and Conference Proceedings*, pages 68–106. JMLR.org, 2014. URL <http://proceedings.mlr.press/v35/menon14.html>.

- [40] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. Adaptive computation and machine learning. MIT Press, 2012. ISBN 978-0-262-01825-8.
- [41] Hussein Mozannar and David A. Sontag. Consistent estimators for learning to defer to an expert. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 7076–7087. PMLR, 2020.
- [42] Hussein Mozannar, Arvind Satyanarayan, and David A. Sontag. Teaching humans when to defer to a classifier via exemplars. *CoRR*, abs/2111.11297, 2021. URL <https://arxiv.org/abs/2111.11297>.
- [43] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Ng. Reading digits in natural images with unsupervised feature learning. *NeurIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 1 2011.
- [44] Chenri Ni, Nontawat Charoenphakdee, Junya Honda, and Masashi Sugiyama. On the calibration of multiclass classification with rejection. In *NeurIPS*, pages 2582–2592, 2019.
- [45] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8024–8035, 2019.
- [46] Bernardo Ávila Pires and Csaba Szepesvári. Multiclass classification calibration functions. *CoRR*, abs/1609.06385, 2016. URL <http://arxiv.org/abs/1609.06385>.
- [47] Harish G. Ramaswamy and Shivani Agarwal. Convex calibration dimension for multiclass loss matrices. *J. Mach. Learn. Res.*, 17:14:1–14:45, 2016. URL <http://jmlr.org/papers/v17/14-316.html>.
- [48] Harish G Ramaswamy, Ambuj Tewari, and Shivani Agarwal. Consistent algorithms for multiclass classification with an abstain option. *Electronic Journal of Statistics*, 12(1):530–554, 2018.
- [49] Alexey E. Rastegin. Bounds of the pinsker and fannes types on the tsallis relative entropy. *Mathematical Physics Analysis & Geometry*, 16(3):213–228, 2013.
- [50] Mark D. Reid and Robert C. Williamson. Composite binary losses. *J. Mach. Learn. Res.*, 11: 2387–2422, 2010.
- [51] Clayton Scott. Surrogate losses and regret bounds for cost-sensitive classification with example-dependent costs. In Lise Getoor and Tobias Scheffer, editors, *ICML*, pages 153–160. Omnipress, 2011. URL [https://icml.cc/2011/papers/138\\_icmlpaper.pdf](https://icml.cc/2011/papers/138_icmlpaper.pdf).
- [52] Song-Qing Shen, Bin-Bin Yang, and Wei Gao. AUC optimization with a reject option. In *AAAI*, pages 5684–5691, 2020.
- [53] Ingo Steinwart. How to compare different loss functions and their risks. *Constructive Approximation*, 26:225–287, 2007.
- [54] Ambuj Tewari and Peter L. Bartlett. On the consistency of multiclass classification methods. *J. Mach. Learn. Res.*, 8:1007–1025, 2007.
- [55] Vladimir Vapnik. *Statistical learning theory*. Wiley, 1998. ISBN 978-0-471-03003-4.
- [56] Rajeev Verma and Eric T. Nalisnick. Calibrated learning to defer with one-vs-all classifiers. In *ICML*, volume 162, pages 22184–22202, 2022.
- [57] Yutong Wang and Clayton Scott. Weston-watkins hinge loss and ordered partitions. In *NeurIPS*, 2020.
- [58] Robert C. Williamson, Elodie Vernet, and Mark D. Reid. Composite multiclass losses. *J. Mach. Learn. Res.*, 17:223:1–223:52, 2016.

- [59] Guoqiang Wu, Chongxuan Li, Kun Xu, and Jun Zhu. Rethinking and reweighting the univariate losses for multi-label ranking: Consistency and generalization. *CoRR*, abs/2105.05026, 2021. URL <https://arxiv.org/abs/2105.05026>.
- [60] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747.
- [61] F. Y. Yang. Maximum lq-likelihood estimation. *Annals of Statistics*, 2010.
- [62] Forest Yang and Sanmi Koyejo. On the consistency of top-k surrogate losses. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 10727–10735. PMLR, 2020.
- [63] Ming Yuan and Marten H. Wegkamp. Classification methods with reject option based on convex risk minimization. *J. Mach. Learn. Res.*, 11:111–130, 2010.
- [64] Mingyuan Zhang, Harish Guruprasad Ramaswamy, and Shivani Agarwal. Convex calibrated surrogates for the multi-label f-measure. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 11246–11255. PMLR, 2020.
- [65] Tong Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5(Oct):1225–1251, 2004.
- [66] Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. In *NeurIPS*, pages 8792–8802, 2018.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] See Appendix H.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Appendix H.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
  - (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See the supplemental material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Section 6 and Appendix.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] See Table 2.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Section 6.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 6.
  - (b) Did you mention the license of the assets? [N/A] The used datasets are open benchmarks.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] Please refer to the supplemental materials.
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

## A Proof of Theorem 1 and Corollary 1

We begin with the proof of Corollary 1 and show that Theorem 1 is its special case.

*Proof.* First of all, we prove that the Bayes optimal solution on  $\tilde{p}(\mathbf{x}, \tilde{y})$  coincide with the Chow's rule of  $p(\mathbf{x}, y)$  with cost  $c$ . According to the optimality condition of multi-class classification, the optimal classifier  $f^*(\mathbf{x})$  on  $\tilde{p}(\mathbf{x}, \tilde{y})$  should fulfill the following condition almost surely:

$$f^*(\mathbf{x}) = \operatorname{argmax}_{\tilde{y}} \tilde{p}(\tilde{y}|\mathbf{x}), \tilde{y} \in \{1, \dots, K, \textcircled{\ast}\}.$$

According to the definition of  $\tilde{p}$ , we can further rewrite it as:

$$f^*(\mathbf{x}) = \begin{cases} \textcircled{\ast}, & \max_{\tilde{y} \in \{1, \dots, K\}} \frac{p(\tilde{y}|\mathbf{x})}{2-c(\mathbf{x})} \leq \frac{1-c(\mathbf{x})}{2-c(\mathbf{x})}, \\ \operatorname{argmax}_{\tilde{y} \in \{1, \dots, K\}} \frac{p(\tilde{y}|\mathbf{x})}{2-c(\mathbf{x})}, & \text{else,} \end{cases}$$

which coincides with the Chow's rule. Then we have the following conclusions:

$$\begin{aligned} \bar{R}_{01}(f) &= \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})}[(2-c(\mathbf{x}))\ell_{01}(f(\mathbf{x}), \tilde{y})] \\ &= \int_{\mathbf{x}} \sum_{\tilde{y}=1}^K (2-c(\mathbf{x}))\ell_{01}(f(\mathbf{x}), \tilde{y}) \frac{p(\mathbf{x}, \tilde{y})}{2-c(\mathbf{x})} d\mathbf{x} + \int_{\mathbf{x}} (2-c(\mathbf{x}))\ell_{01}(f(\mathbf{x}), K+1) \frac{p(\mathbf{x})}{2-c(\mathbf{x})} d\mathbf{x} \\ &= \int_{\mathbf{x}} \sum_{\tilde{y}=1}^K (2-c(\mathbf{x}))\ell_{01}(f(\mathbf{x}), \tilde{y}) \frac{p(\mathbf{x}, \tilde{y})}{2-c(\mathbf{x})} d\mathbf{x} + \int_{\mathbf{x}} (2-c(\mathbf{x}))\ell_{01}(f(\mathbf{x}), \textcircled{\ast}) \frac{(1-c(\mathbf{x}))p(\mathbf{x})}{2-c(\mathbf{x})} d\mathbf{x} \\ &= \int_{\mathbf{x}} \sum_{\tilde{y}=1}^K \ell_{01}(f(\mathbf{x}), \tilde{y}) p(\mathbf{x}, \tilde{y}) d\mathbf{x} + \int_{\mathbf{x}} \ell_{01}(f(\mathbf{x}), \textcircled{\ast}) (1-c(\mathbf{x})) p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

Suppose  $C(f(\mathbf{x})) = \sum_{\tilde{y}=1}^K \ell_{01}(f(\mathbf{x}), \tilde{y}) p(\tilde{y}|\mathbf{x}) + \ell_{01}(f(\mathbf{x}), \textcircled{\ast}) (1-c(\mathbf{x}))$  is the inner risk and  $f^*$  is the Chow's rule, we have that

- If  $f^*(\mathbf{x}) = \textcircled{\ast}$  and  $f(\mathbf{x}) \neq f^*(\mathbf{x})$ :

$$C(f(\mathbf{x})) - C(f^*(\mathbf{x})) = 1 - c(\mathbf{x}) - p(f(\mathbf{x})|\mathbf{x}).$$

- If  $f^*(\mathbf{x}) \in \{1, \dots, K\}$  and  $f(\mathbf{x}) = \textcircled{\ast}$ :

$$C(f(\mathbf{x})) - C(f^*(\mathbf{x})) = c(\mathbf{x}) - 1 + p(f^*(\mathbf{x})|\mathbf{x}).$$

- If  $f^*(\mathbf{x}), f(\mathbf{x}) \in \{1, \dots, K\}$  and  $f(\mathbf{x}) \neq f^*(\mathbf{x})$ :

$$C(f(\mathbf{x})) - C(f^*(\mathbf{x})) = p(f^*(\mathbf{x})|\mathbf{x}) - p(f(\mathbf{x})|\mathbf{x}).$$

These conclusion shows that

$$C(f(\mathbf{x})) - C(f^*(\mathbf{x})) = \mathbb{E}_{p(y|\mathbf{x})}[\ell_{01c}(f(\mathbf{x}), y)] - \mathbb{E}_{p(y|\mathbf{x})}[\ell_{01c}(f^*(\mathbf{x}), y)].$$

We can conclude the proof by taking the expectation over  $p(\mathbf{x})$  on both sides of the equation.  $\square$

It can be seen that when  $c(\mathbf{x})$  is constant, we can divide each side of Corollary 1 to get the proof of Theorem 1.

## B Proof of Theorem 2

*Proof.*

$$\begin{aligned} R_{L_c^\Phi}(\mathbf{g}) &= \mathbb{E}_{p(\mathbf{x}, y)}[L_c^\Phi(\mathbf{g}(\mathbf{x}), y)] \\ &= \mathbb{E}_{p(\mathbf{x}, y)}[\Phi(\mathbf{g}(\mathbf{x}), y)] + (1-c)\mathbb{E}_{p(\mathbf{x})}[\Phi(\mathbf{g}(\mathbf{x}), K+1)] \\ &= (2-c)\tilde{R}_\Phi(\mathbf{g}) \end{aligned}$$

$\square$

### C Proof of Theorem 3

We first give the definition of Rademacher complexity:

**Definition 7.** (*Rademacher complexity [5]*) Let  $Z_1, \dots, Z_n$  be *n i.i.d.* random variables drawn from a probability distribution  $\mu$  and  $\mathcal{F} = \{f : Z \rightarrow \mathbb{R}\}$  be a class of measurable functions. Then the expected Rademacher complexity of function class  $\mathcal{F}$  is given as follow:

$$\mathfrak{R}_n(\mathcal{F}) = \mathbb{E}_{Z_1, \dots, Z_n \sim \mu} \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(Z_i) \right], \quad (7)$$

where  $\sigma_1, \dots, \sigma_n$  are the Rademacher variables that take the value from  $\{-1, +1\}$  uniformly.

Then we can begin proving Theorem 3.

*Proof.* According to the conditions in Theorem 3, we can learn that  $L_c^\Phi$  is  $(2-c)\rho$ -Lipschitz continuous and is bounded by  $(2-c)C_\Phi$ . By applying the McDiarmid's inequality [37], it is routine [40] to show that the following inequalities holds with probability at least  $1 - \frac{\delta}{2}$ , respectively:

$$\begin{aligned} \sup_{g \in \mathcal{G}} \left( R_{L_c^\Phi}(g) - \hat{R}_{L_c^\Phi}(g) \right) &\leq \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_n} \left[ \sup_{g \in \mathcal{G}} \left( R_{L_c^\Phi}(g) - \hat{R}_{L_c^\Phi}(g) \right) \right] + (2-c)C_\Phi \sqrt{\frac{\log \frac{2}{\delta}}{2n}} \\ \sup_{g \in \mathcal{G}} \left( \hat{R}_{L_c^\Phi}(g) - R_{L_c^\Phi}(g) \right) &\leq \mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_n} \left[ \sup_{g \in \mathcal{G}} \left( \hat{R}_{L_c^\Phi}(g) - R_{L_c^\Phi}(g) \right) \right] + (2-c)C_\Phi \sqrt{\frac{\log \frac{2}{\delta}}{2n}} \end{aligned}$$

By applying Talagrand's contraction lemma [36], we can learn that:

$$\mathbb{E}_{\mathbf{x}_1, \dots, \mathbf{x}_n} \left[ \sup_{g \in \mathcal{G}} \left( R_{L_c^\Phi}(g) - \hat{R}_{L_c^\Phi}(g) \right) \right] \leq \sqrt{2}(2-c)\rho \sum_{y=1}^{K+1} \mathfrak{R}_n(\mathcal{G}_y)$$

and this conclusion also holds for another direction. Plugging this conclusion into the former inequalities and using the union bound, we can learn:

$$\sup_{g \in \mathcal{G}} \left| R_{L_c^\Phi}(g) - \hat{R}_{L_c^\Phi}(g) \right| \leq \sqrt{2}(2-c)\rho \sum_{y=1}^{K+1} \mathfrak{R}_n(\mathcal{G}_y) + (2-c)C_\Phi \sqrt{\frac{\log \frac{2}{\delta}}{2n}}$$

According to the definition of empirical risk minimization and identifiable condition, we can get the following conclusion, where  $\mathbf{g}^*$  is the optimal solution among all the measurable functions:

$$\begin{aligned} R_{L_c^\Phi}(\hat{\mathbf{g}}) - \min_{g \in \mathcal{G}} R_{L_c^\Phi} &= R_{L_c^\Phi}(\hat{\mathbf{g}}) - R_{L_c^\Phi}^* \\ &= \left( R_{L_c^\Phi}(\hat{\mathbf{g}}) - \hat{R}_{L_c^\Phi}(\hat{\mathbf{g}}) \right) + \left( \hat{R}_{L_c^\Phi}(\hat{\mathbf{g}}) - \hat{R}_{L_c^\Phi}(\mathbf{g}^*) \right) + \left( \hat{R}_{L_c^\Phi}(\mathbf{g}^*) - R_{L_c^\Phi}^* \right) \\ &\leq \left( R_{L_c^\Phi}(\hat{\mathbf{g}}) - \hat{R}_{L_c^\Phi}(\hat{\mathbf{g}}) \right) + \left( \hat{R}_{L_c^\Phi}(\mathbf{g}^*) - R_{L_c^\Phi}^* \right) \\ &\leq 2 \sup_{g \in \mathcal{G}} \left| R_{L_c^\Phi}(g) - \hat{R}_{L_c^\Phi}(g) \right| \end{aligned}$$

which concludes the proof.  $\square$

### D Proof of Theorem 4

*Proof.* According to Theorem 1, Theorem 2, and Theorem 3 in Ramaswamy and Agarwal [47], we can immediately learn the sufficiency of this condition. We begin the proof of the necessity of the calibration of  $\Phi$ .

First, we give a useful property for permutation-invariant surrogates. For any  $\mathbf{p} \in \Delta^{K+1}$ ,  $\mathcal{U}(\mathbf{p})$  denotes all the  $\mathbf{u} \in \mathbb{R}^{K+1} : \mathbf{u} \notin \operatorname{argmin}_{\mathbf{u}'} \mathbf{p}^T \mathbf{L}_{01}(\mathbf{u}')$  that meet the following condition:

$$\mathbf{p}^T \Phi(\mathbf{u}) = \inf_{\mathbf{u}' \in \mathbb{R}^{K+1}} \mathbf{p}^T \Phi(\mathbf{u}'),$$



and  $\mathcal{P}$  denotes the collection of  $\mathbf{p} \in \Delta^{K+1}$  that  $\mathcal{U}(\mathbf{p})$  is non-empty. Given the assumption of permutation-invariance, we can conclude that if  $\mathbf{p} \in \mathcal{P}$ , the set  $\mathcal{P}$  contains all of its re-permutation, since given any permutation matrix  $P$  and  $\mathbf{u} \in \mathcal{U}(\mathbf{p})$ :

$$\begin{aligned}
(P\mathbf{p})^T \Phi(P\mathbf{u}) &= \mathbf{p}^T P^T P \Phi(\mathbf{u}) \\
&= \mathbf{p}^T \Phi(\mathbf{u}) \\
&= \inf_{\mathbf{u} \in \mathbb{R}^{K+1}} \mathbf{p}^T \Phi(\mathbf{u}) \\
&= \inf_{\mathbf{u} \in \mathbb{R}^{K+1}} \mathbf{p}^T P^T P \Phi(\mathbf{u}) \\
&= \inf_{\mathbf{u} \in \mathbb{R}^{K+1}} (P\mathbf{p})^T \Phi(P\mathbf{u}) \\
&= \inf_{\mathbf{u}' \in \mathbb{R}^{K+1}} (P\mathbf{p})^T \Phi(\mathbf{u}')
\end{aligned}$$

which indicates  $P\mathbf{u} \in \mathcal{U}(P\mathbf{p})$  and thus  $P\mathbf{p} \in \mathcal{P}$ .

Then we continue to prove the necessity of  $\Phi$ 's calibration by contradiction. Suppose  $\Phi$  is not classification-calibrated and thus  $\mathcal{P}$  is non-empty. Denote with  $\tilde{\mathcal{P}}$  the collection of all the potential confidence vectors  $\{\tilde{p}(1|\mathbf{x}), \dots, \tilde{p}(K+1|\mathbf{x})\}$  for all the  $c \in [0, 1]$ . We can learn that  $L_c^\Phi$  cannot be  $\ell_{01c}$ -consistent if  $\mathcal{P} \cap \tilde{\mathcal{P}}$  is non-empty or we can construct a distribution with density  $\hat{p}$  whose support only contains a single point  $\mathbf{x}^3$  with posterior probability  $\mathbf{p}' \in \mathcal{P} \cap \tilde{\mathcal{P}}$  for contradiction.

Then we show that  $\mathcal{P} \cup \tilde{\mathcal{P}}$  must be non-empty if  $\Phi$  is permutation-invariant and not classification-calibrated. For any  $\mathbf{p} \in \mathcal{P}$ , if  $p_{K+1} \leq \frac{1}{2}$ , we can learn that  $\mathbf{p} \in \tilde{\mathcal{P}}$  constructed by  $\bar{\mathbf{p}} \in \Delta^K$  and cost  $c$  following the definition of self-augmented distribution:

$$c = \frac{1 - 2p_{K+1}}{1 - p_{K+1}}, \bar{p}_y = (2 - c)p_y = \frac{p_y}{1 - p_{K+1}}, \forall y \in \{1, \dots, K\},$$

and thus  $\mathcal{P} \cap \tilde{\mathcal{P}}$  is non-empty. Can we get a non-empty  $\mathcal{P}$  that all its elements have a  $(K+1)$ th element that is larger than  $\frac{1}{2}$ ? The answer is no: for any  $\mathbf{p} \in \Delta^{K+1}$ ,  $\min_{\bar{y} \in \{1, \dots, K+1\}} p_{\bar{y}} < 1/2$  and the fact that  $\mathcal{P}$  is closed w.r.t. the operation of permutation indicates that we can exchange the minimal value in  $\mathbf{p}$  and its  $(K+1)$ th element to get a  $\mathbf{p}' \in \mathcal{P}$  and  $p'_{K+1} < 1/2$ .

In conclusion,  $\mathcal{P} \cup \tilde{\mathcal{P}}$  must be non-empty if  $\Phi$  is not classification-calibrated and permutation-invariant, which concludes the proof of necessity.  $\square$

## E Proof of Theorem 5

*Proof.* According to [61], we can directly get the formulation of the optimal solution of GCE. Based on this formulation, we prove the classification-calibration of GCE constructively by giving an regret transfer bound.

First of all, we show that the excess error of GCE loss for any  $\mathbf{x}$  is a reweighted version of the Tsallis relative entropy [23, 49] in actual. Denote by  $S(\mathbf{g}^*)_y = \mathbf{q}_y^*$ ,  $S(\mathbf{g})_y = \mathbf{q}_y$  for any  $\mathbf{g}$ , and  $p(y|\mathbf{x}) = \eta_y$ . We substitute  $\gamma$  with  $r$  in the proof for simplicity:

$$\begin{aligned}
Ex(\mathbf{q}, \mathbf{x}) &= \sum_{y=1}^y \eta_y \frac{(1 - \mathbf{q}_y^r)}{r} - \sum_{y=1}^y \eta_y \frac{(1 - \mathbf{q}_y^{*r})}{r} \\
&= \frac{\sum_{y=1}^K \eta_y (\mathbf{q}_y^{*r} - \mathbf{q}_y^r)}{r} \\
&= \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{1-r} \frac{\left( 1 - \sum_{y=1}^K \mathbf{q}_y^{*(1-r)} \mathbf{q}_y^r \right)}{r}
\end{aligned}$$

It can be seen that the second term of the last equation is the Tsallis relative entropy between discrete possibilities  $\mathbf{q}^*$  and  $\mathbf{q}$ . According to the Corollary 9 of [23] and (4.13) of [49], we can lower bound

<sup>3</sup>An equivalent description is that  $\hat{p}(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{x}')$ , where  $\delta(\cdot)$  is the Dirac Delta function.

the excess error with the total variation distance between  $\mathbf{q}^*$  and  $\mathbf{q}$  and get a Pinsker's type inequality:

$$Ex(\mathbf{q}, \mathbf{x}) \geq \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{1-r} \frac{1-r}{2} \|\mathbf{q}^* - \mathbf{q}\|_1^2$$

Then we have to connect the *r.h.s.* of the inequality to the excess error w.r.t. 0-1 loss. When  $\operatorname{argmax}_y q_y(\mathbf{x}) \neq \operatorname{argmax}_y \eta_y$ , denote by  $\operatorname{argmax}_y q_y(\mathbf{x}) = \text{pred}$  and  $\operatorname{argmax}_y \eta_y = \text{max}$ :

$$\begin{aligned} \|\mathbf{q}^* - \mathbf{q}\|_1 &= \sum_{y=1}^K |q_y^* - q_y| \\ &\geq |q_{\text{max}}^* - q_{\text{max}}| + |q_{\text{pred}}^* - q_{\text{pred}}| \\ &\geq |q_{\text{max}}^* - q_{\text{pred}}^* + q_{\text{pred}} - q_{\text{max}}| \end{aligned}$$

According to the formulation of the optimal solution of GCE, we can learn that  $q_{\text{max}}^* \geq q_{\text{pred}}^*$ . Since  $\operatorname{argmax}_y q_y(\mathbf{x}) \neq \operatorname{argmax}_y \eta_y$ , we can learn that  $q_{\text{pred}} \geq q_{\text{max}}$ . Then we can further learn that:

$$\begin{aligned} \|\mathbf{q}^* - \mathbf{q}\|_1 &\geq |q_{\text{max}}^* - q_{\text{pred}}^*| \\ &= \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{-1} |\eta_{\text{max}}^{\frac{1}{1-r}} - \eta_{\text{pred}}^{\frac{1}{1-r}}| \\ &= \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{-1} (\eta_{\text{max}}^{\frac{1}{1-r}} - \eta_{\text{pred}}^{\frac{1}{1-r}}) \\ &= \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{-1} (\eta_{\text{max}} * \eta_{\text{max}}^{\frac{r}{1-r}} - \eta_{\text{pred}} * \eta_{\text{pred}}^{\frac{r}{1-r}}) \\ &\geq \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{-1} (\eta_{\text{max}} * \eta_{\text{max}}^{\frac{r}{1-r}} - \eta_{\text{pred}} * \eta_{\text{max}}^{\frac{r}{1-r}}) \\ &= \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{-1} \eta_{\text{max}}^{\frac{r}{1-r}} (\eta_{\text{max}} - \eta_{\text{pred}}) \end{aligned}$$

Then we can learn that:

$$\begin{aligned} Ex(\mathbf{q}, \mathbf{x}) &\geq \left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{-1-r} \eta_{\text{max}}^{\frac{2r}{1-r}} * \frac{1-r}{2} (\eta_{\text{max}} - \eta_{\text{pred}})^2 \\ &= \frac{1}{\left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{1+r} \eta_{\text{max}}^{\frac{2r}{1-r}}} * \frac{1-r}{2} (\eta_{\text{max}} - \eta_{\text{pred}})^2 \\ &\geq \frac{1}{\left( \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} \right)^{1+r}} \frac{1-r}{2K^{\frac{2r}{1-r}}} (\eta_{\text{max}} - \eta_{\text{pred}})^2 \tag{8} \end{aligned}$$

$$\geq \frac{1-r}{2K^{\frac{2r}{1-r}}} (\eta_{\text{max}} - \eta_{\text{pred}})^2 \tag{9}$$

The derivation of (8) to (9) is shown in the end of this proof. Then we have the following regret transfer bound:

$$\begin{aligned} R_{01}(\operatorname{argmax}_y \mathbf{g}_y) - R_{01}^* &\leq \mathbb{E}_{p(\mathbf{x})} \left[ \sqrt{C Ex(\mathbf{q}, \mathbf{x})} \right] \\ &\leq \sqrt{C \mathbb{E}_{p(\mathbf{x})} [Ex(\mathbf{q}, \mathbf{x})]} \quad (\text{Jensen's inequality}) \\ &= \sqrt{C(R_G(\mathbf{g}) - R_G^*)} \end{aligned}$$

where  $C = \frac{2K^{\frac{2r}{1-r}}}{1-r}$ ,  $R_G$  is the expected version of GCE loss, and  $R_G^*$  and  $R_{01}^*$  are the optimal value of the expected version of GCE loss and 0-1 loss, respectively. From this bound, we constructively prove the classification-calibration of GCE loss with  $r \in (0, 1)$ .

**Proof of (8) to (9):** To complete the step, we only have to show that the term  $\frac{1}{\left(\sum_{y=1}^K \eta_y^{\frac{1}{1-r}}\right)^{1+r}}$

achieves its minima at any  $e_y$ , where  $e_y = [0, \dots, 1, \dots, 0]$  that has the only non-zero value at dimension  $y$ . Notice that the following conclusion about  $p$ -norm holds for any real-valued  $K$ -dimensional vector and  $r \in (0, 1)$ :

$$\|\mathbf{x}\|_{\frac{1}{1-r}} \leq \|\mathbf{x}\|_1.$$

Suppose  $\boldsymbol{\eta}$  is the vector consists of  $\{\eta_y\}_{y=1}^K$ . Then:

$$\begin{aligned} \sum_{y=1}^K \eta_y^{\frac{1}{1-r}} &= \|\boldsymbol{\eta}\|_{\frac{1}{1-r}} \\ &\leq \|\boldsymbol{\eta}\|_1^{\frac{1}{1-r}} \\ &= 1 \end{aligned}$$

Notice that when  $\boldsymbol{\eta} = e_y$  for some  $y$ ,  $\sum_{y=1}^K \eta_y^{\frac{1}{1-r}} = 1$ , which indicates that  $\sum_{y=1}^K \eta_y^{\frac{1}{1-r}}$ 's maximum value is 1 when confined in probability simplex  $\Delta^K$ . Then we can further learn that  $\frac{1}{\left(\sum_{y=1}^K \eta_y^{\frac{1}{1-r}}\right)^{1+r}} \geq 1$ , which concludes the proof.  $\square$

It is noticeable that the bound does not hold for  $r = 1$ , e.g., the case of MAE loss, and the regret transfer bound becomes less compact when  $r$  increases. We prove the classification-calibration of MAE loss by showing its regret transfer bound.

**Corollary 2.** Suppose the expected version of MAE loss is  $R_M(\mathbf{g})$  and its minimal value is  $R_M^*$ . Then we have:

$$R_{01}(\operatorname{argmax}_y \mathbf{g}_y) - R_{01}^* \leq K(R_M(\mathbf{g}) - R_M^*).$$

*Proof.* Given the formulation of the optimal solution  $\mathbf{q}^*$  of expected MAE loss in Theorem 5, for any  $\mathbf{x}$ , the excess error can be written as:

$$\begin{aligned} Ex(\mathbf{q}, \mathbf{x}) &= \sum_{y=1}^K \eta_y(1 - q_y) - \sum_{y=1}^K \eta_y(1 - q_y^*) \\ &= \sum_{y=1}^K \eta_y(q_y^* - q_y) \\ &= \eta_{max} - \sum_{y=1}^K \eta_y q_y \end{aligned}$$

When  $\operatorname{argmax}_y q_y(\mathbf{x}) \neq \operatorname{argmax}_y \eta_y$ :

$$\begin{aligned} \eta_{max} - \sum_{y=1}^K \eta_y q_y &= \eta_{max} - \eta_{pred} q_{pred} - \sum_{y \neq pred}^K \eta_y q_y \\ &\geq \eta_{max} - \eta_{pred} q_{pred} - \eta_{max}(1 - q_{pred}) \\ &= q_{pred}(\eta_{max} - \eta_{pred}) \\ &\geq \frac{1}{K}(\eta_{max} - \eta_{pred}), \end{aligned}$$

which concludes the proof by taking the expectation on both sides.  $\square$

Combine the conclusions above and we can conclude the proof. Though the bound for GCE becomes less tight when  $r$  increases, the MAE loss has a better regret transfer bound, which indicates that the regret transfer bound of GCE for  $r \in (0, 1)$  may not be good enough. A potential reason is that [23, 49] considered the general case of Tsallis relative entropy while we only need the case that  $q$  is a probability distribution. It is promising to further tighten this bound by modifying the conclusions in [23, 49] and limiting  $q$  to a  $K - 1$ -dimensional probability simplex.

## F Details of the Experiment Setup

### F.1 Detailed Information of Benchmark Datasets

In the experiments, we used 3 widely-used benchmark datasets. Here, we report the sources of these datasets and the way we split them.

- Fashion-MNIST [60]. It is a 10-class dataset of fashion items. Each instance is a 28\*28 grayscale image. Source: <https://github.com/zalandoresearch/fashion-mnist>.
- SVHN [43] It is a 10-class dataset for 10 different digits and each instance is a 32\*32\*3 colored image in RGB format. Source: <http://ufldl.stanford.edu/housenumbers/>.
- CIFAR-10 [31]. It is a 10-class dataset for 10 different objects and each instance is a 32\*32\*3 colored image in RGB format. Source: <https://www.cs.toronto.edu/~kriz/cifar.html>.

For Fashion-MNIST and SVHN, we trained models on the whole training dataset. For CIFAR-10, we splitted 10% of the training dataset as the validation set and conducted random crop and flips for data augmentation. The cost  $c$  is less than 0.5 as suggested in [48] and further decreased on Fashion-MNIST since it is a less difficult dataset.

### F.2 Detailed Information of the Models and Optimization Algorithm

For Fashion-MNIST, we used the model defined in [8] for the experiments. For SVHN and CIFAR-10, ResNet-18 and ResNet-34 is used, respectively. For the cost-sensitive method [8], we use batch normalization [28] at the output layer as suggested in [8] since it fails to work without this modification.

Adam with default momentum was used for optimization in this paper. For Fashion-MNIST, the epoch number, batch size, learning rate, and weight decay are set to 20, 256, 1e-3, and 1e-4. For SVHN, the epoch number, batch size, learning rate, and weight decay are set to 20, 1024, 1e-3, and 1e-4. For CIFAR-10, the epoch number, batch size, learning rate, and weight decay are set to or selected from 200, 1024, {1e-3, 2e-3, 3e-3}, and 1e-4. For Fashion-MNIST and SVHN, we use the model after the 20th epoch for performance evaluation. For CIFAR-10, we report the performance of the model with the best performance on the validation dataset. Temperature scaling is further conducted for CE on CIFAR-10.

## G Details of Instance-dependent Rejection Cost

In practical applications, it can be beneficial letting the rejection cost  $c(x)$  vary among different samples. For example, when constructing a system to automatically prescribe for users, a wrong prescription can be fatal for users of advanced ages or with underlying diseases. To prevent such wrong prescriptions, the cost for this type of users can be decreased to encourage rejection. However, it is not suitable encouraging rejection for all the users, which makes the system meaningless. An acceptable choice is to increase the cost for rejection instead for users of low risk.

In this appendix, we expand the Theorem 2 and propose a surrogate for instance dependent cost based on Corollary 1, whose estimation error bound and calibration analysis can be derived almost symmetrically thanks to the equivalence shown in Corollary 3. Then we further evaluate its performance on SVHN dataset.

## G.1 Expansion of Theorem 2

Theorem 2 tells the equivalence between surrogate risk minimization of  $L_c^\Phi$  on  $p(\mathbf{x}, y)$  and surrogate risk minimization of  $\Phi$  on  $\tilde{p}(\mathbf{x}, \tilde{y})$ . Here we expand it to the case of instance-dependent cost.

Given the cost function  $c(\mathbf{x})$  and any function  $\Phi(\cdot) : \mathbb{R}^{K+1} \times \{1, \dots, K+1\} \rightarrow \mathbb{R}^+$ :

$$L_{c(\mathbf{x})}^\Phi(\mathbf{u}, y) = (\Phi(\mathbf{u}, y) + (1 - c(\mathbf{x}))\Phi(\mathbf{u}, K+1))/(2 - c(\mathbf{x})).$$

Then we have the following conclusion:

**Corollary 3.** For any  $\mathbf{g} : \mathcal{X} \rightarrow \mathbb{R}^{K+1}$  and  $R_{L_{c(\mathbf{x})}^\Phi}(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)}[L_{c(\mathbf{x})}^\Phi(\mathbf{g}(\mathbf{x}), y)]$ :

$$R_{L_{c(\mathbf{x})}^\Phi}(\mathbf{g}) = \tilde{R}_\Phi(\mathbf{g})$$

*Proof.*

$$\begin{aligned} R_{L_{c(\mathbf{x})}^\Phi}(\mathbf{g}) &= \mathbb{E}_{p(\mathbf{x}, y)}[L_{c(\mathbf{x})}^\Phi(\mathbf{g}(\mathbf{x}), y)] \\ &= \mathbb{E}_{p(\mathbf{x}, y)}[(\Phi(\mathbf{g}(\mathbf{x}), y) + (1 - c(\mathbf{x}))\Phi(\mathbf{g}(\mathbf{x}), K+1))/(2 - c(\mathbf{x}))] \\ &= \int_{\mathbf{x}} \sum_{y=1}^K \Phi(\mathbf{g}(\mathbf{x}), y) \frac{p(\mathbf{x}, y)}{2 - c(\mathbf{x})} d\mathbf{x} + \int_{\mathbf{x}} \frac{(1 - c(\mathbf{x}))p(\mathbf{x})}{2 - c(\mathbf{x})} \Phi(\mathbf{g}(\mathbf{x}), K+1) d\mathbf{x} \\ &= \tilde{R}_\Phi(\mathbf{g}) \end{aligned}$$

□

The derivation of its estimation error bound is similar to that of Theorem 3 by modifying the upper bound and Lipschitz constant, and the necessity and sufficiency of the  $\ell_{01}$ -calibration of  $\Phi$  can also be proved by utilizing the arbitrariness of  $\tilde{p}(\mathbf{x}, y)$  as in Appendix D.

## G.2 Experiments on SVHN

In this section, we compare our proposed surrogate  $L_{c(\mathbf{x})}^\Phi$  with CE and DEFER on SVHN. The cost-sensitive learning-based method [8] is not compared since it cannot tackle the case of instance-dependent cost.

In the experiments, we use SVHN [43] to demonstrate the effectiveness of  $L_{c(\mathbf{x})}^\Phi$ . To generate instance-dependent costs, we split 10% of the training dataset and manually corrupt it into a binary dataset by aggregating the 10 classes into ['0', '2', '3', '5', '6', '8', '9'] and ['1', '4', '7']. We train a binary classifier with on the corrupted dataset with 10 epochs. Then we further use the obtained classifier on training and testing set to split them into 2 parts. For any  $\mathbf{x}$  that is classified as ['0', '2', '3', '5', '6', '8', '9'], we set  $c(\mathbf{x}) = c_1$  and  $c_2$  otherwise. In the experiments, Adam with default momentum is used with learning rate, batch size and weight decay set to 1e-3, 1024, and 1e-4, respectively. The model used is ResNet-18.

**Table 3:** The mean and standard error of the zero-one- $c$  losses (**01c**, rescaled to 0-100), rejection ratio (**Rej**), and missclassification rates (**01**) of the accepted data for 5 trails. The best and comparable methods based on the paired t-test at the significance level 5% are highlighted in boldface.

Method	$(c_1, c_2)$	CE			DEFER			GCE		
		01c	Rej	01	01c	Rej	01	01c	Rej	01
SVHN	(0.50, 0.10)	8.03 (0.16)	4.46 (0.54)	7.60 (0.01)	8.00 (0.30)	9.20 (0.72)	5.07 (0.25)	<b>7.20</b> <b>(0.17)</b>	6.73 (6.73)	5.13 (5.13)
	(0.45, 0.15)	7.80 (0.26)	4.36 (0.31)	7.03 (0.23)	9.10 (0.46)	9.93 (0.41)	<b>5.07</b> (0.42)	6.93 <b>(0.31)</b>	7.00 (0.35)	4.70 (0.26)
	(0.40, 0.20)	7.70 (0.10)	4.50 (0.50)	6.83 (0.25)	7.80 (0.26)	11.13 (1.27)	5.00 (0.44)	<b>7.03</b> <b>(0.21)</b>	7.93 (0.55)	4.80 (0.17)
	(0.35, 0.25)	7.76 (0.12)	4.90 (0.20)	6.67 (0.15)	7.70 (0.26)	11.93 (0.45)	4.80 (0.10)	<b>6.83</b> <b>(0.20)</b>	8.43 (0.31)	4.63 (0.15)

The experimental results are reported in the table above. It can be seen that in the scenario of instance-dependent cost, the proposed surrogate with GCE loss still outperforms baseline methods, which aligns with the observations in Section 6.

## H Limitations and Potential Negative Social Impacts

**Limitations:** This framework is used for multi-class classification with rejection, while there are also other scenarios for learning with rejection, e.g., AUC optimization with rejection [52]. We believe that extensions to CwR with complex evaluation is a promising future direction.

**Potential Negative Social Impacts:** Though classification with rejection can be useful in risk-critical missions, it can lead to inefficient services once abused, i.e., used in risk-insensitive missions. This is also the potential negative social impact of all the methods for CwR.

## I Synthetic Experiments

In this section, we evaluate the performance of our proposed method GCE with  $q = 0.7$  on a synthetic dataset and compare it with the bayes optimal model. The synthetic dataset has a class number of 3 and its class-conditional probabilities *w.r.t.* different classes are generated by three different 2-dimensional Gaussian distributions whose means and covariance matrices are shown below:

$$\boldsymbol{\mu}_1 = [-0.3, 0.3]^\top, \boldsymbol{\mu}_2 = [0.3, 0.3]^\top, \boldsymbol{\mu}_3 = [0.3, -0.3]^\top, \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}_3 = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}.$$

The class-prior probabilities are 1/3 for each class. We generate 6000 data points for training a linear model with Adam and 30000 data points for testing. Since we have access to the *pdf*  $p(\boldsymbol{x})$  and  $p(\boldsymbol{x}, y)$ , we can calculate the class-posterior probability for each data point in the testing set and approximate the performance of the bayes optimal model efficiently by directly applying the Chow’s rule. We report the performances in the following table, and the notations are further described in the caption of this table. It can be seen that the proposed method can well approximate the optimal model even only linear model is considered.

**Table 4:** The mean values of the zero-one- $c$  losses (**01c**, rescaled to 0-100), rejection ratio (**Rej**), and missclassification rates (**01**) of the accepted data for 5 trails. The performance of the Bayes Optimal solution is calculated following the Chow’s rule on the test dataset.

$c$	GCE			Bayes Optimal		
	01c	Rej	01	01c	Rej	01
0.10	8.54	75.51	4.04	8.41	73.33	4.44
0.15	12.25	67.96	6.42	11.72	61.46	6.48
0.20	15.31	59.87	8.31	14.57	52.60	8.53
0.25	17.90	46.01	11.84	16.93	44.41	10.50

## J Connections with Learning to Defer

There are two concurrent papers about learning to defer [10, 56], and one of them [10] is closely related to this paper that our proposed surrogate formulation for instance-independent cost can be seen as their special case with a prior knowledge of expert prediction’s accuracy. The framework in Charusaie et al. [10] allows the use of any classification-calibrated surrogates for consistent learning to defer. They further compared the cases of joint-learning and two-stage learning by analyzing function spaces of classifiers and rejectors with fixed VC-dimension, and they also provided an active-learning extension of their framework. However, due to the stochastic nature of expert prediction, their regret transfer bound relies on an assumption on the calibration function, while our bound can be applied on any calibrated surrogates.

In fact, our proposed surrogate for instance-dependent rejection cost can provide a potential problem reduction for learning to defer if an extra step is added. Firstly, we can see the prediction accuracy

of the expert on a certain instance  $\mathbf{x}$ ,  $\mathbb{E}_{\mu_{M,Y|X=\mathbf{x}}}[\mathbb{I}(M \neq Y)]$ , as the instance-dependent cost  $c(\mathbf{x})$ . Given the fact that we only have  $(M_i, Y_i)$  that is drawn i.i.d. from distribution  $\mu_{M,Y|X}$ , we can split some samples out to train a regression model for predicting  $c(\mathbf{x})$  with consistency guarantee, and then plug the predicted  $c(\mathbf{x})$  into the loss formulation in Appendix G.1 to convert the problem of learning to defer into classification with rejection. This baseline can also be considered in the future study of learning to defer.