# Optimal strategies for reject option classifiers

**Vojtech Franc**                                                    XFRANCV@FEL.CVUT.CZ

**Daniel Prusa**                                                        PRUSA@FEL.CVUT.CZ

**Vaclav Voracek**                                                 VORACVA1@FEL.CVUT.CZ
*Department of Cybernetics, Faculty of Electrical Engineering*
*Czech Technical University in Prague, Czech Republic*

**Editor:** Editors

## Abstract

In classification with a reject option, the classifier is allowed in uncertain cases to abstain from prediction. The classical cost-based model of a reject option classifier requires the cost of rejection to be defined explicitly. An alternative bounded-improvement model, avoiding the notion of the reject cost, seeks for a classifier with a guaranteed selective risk and maximal cover. We coin a symmetric definition, the bounded-coverage model, which seeks for a classifier with minimal selective risk and guaranteed coverage. We prove that despite their different formulations the three rejection models lead to the same prediction strategy: a Bayes classifier endowed with a randomized Bayes selection function. We define a notion of a proper uncertainty score as a scalar summary of prediction uncertainty sufficient to construct the randomized Bayes selection function. We propose two algorithms to learn the proper uncertainty score from examples for an arbitrary black-box classifier. We prove that both algorithms provide Fisher consistent estimates of the proper uncertainty score and we demonstrate their efficiency on different prediction problems including classification, ordinal regression and structured output classification.

**Keywords:** Reject option classification, prediction uncertainty, selective classifiers

## 1. Introduction

In safety critical applications of classification models, prediction errors may lead to serious losses. In such cases estimating when the model makes an error can be as important as its average performance. These two objectives are taken into account in classification with a reject option when the classifier is allowed in uncertain cases to abstain from prediction.

The cost-based model of a classification strategy with the reject option was proposed by Chow in his pioneering work Chow (1970). The goal is to minimize the expected loss equal to the cost of misclassification, when the classifier predicts, and to the reject cost, when the classifier abstains from prediction. An optimal strategy leads to the Bayes classifier abstaining from prediction when the conditional expected risk exceeds the reject cost. The known form of the optimal strategy allows to construct the classifier by plugging in an estimate of the class posterior probabilities to the formula for the conditional risk. Besides the plug-in rule, the reject option classifiers can be learned by empirical risk minimization based approaches like e.g. modifications of Support Vector Machines (Grandvalet et al.,

2008), Boosting (Cortes et al., 2016), or Prototype-based classifiers (Villman et al., 2016) to name a few.

The cost-based model requires the reject cost to be defined explicitly which is difficult in some applications e.g. when the misclassification costs have different physical units than the reject cost. An alternative bounded-improvement model coined in Pietraszek (2005) avoids explicit definition of the reject cost. The rejection strategy is evaluated by two antagonistic quantities: i) a selective risk defined as the expected misclassification cost on accepted predictions and ii) a coverage which corresponds to the probability that the prediction is accepted. An optimal strategy for the bounded-improvement model is the one which maximizes the coverage under the condition that the selective risk does not exceed a target value. In contrast to the cost-based model, it has not been formally shown what is the optimal prediction strategies when the underlying model is known. A solution has been proposed only for special instances of the task. Pietraszek (2005) coined a method based on ROC analysis applicable when a score proportional to posterior probabilities is known and the task is to find only the optimal thresholds. El-Yaniv and Wiener (2010) proposed an algorithm learning the optimal strategy in the noise-free setting, i.e. when a perfect strategy with zero selective risk exists. Geifman and El-Yaniv (2017) shows how to equip a trained classifier with a reject option provided an uncertainty measure is known and the task is to find only a rejection threshold optimal under the bounded-improvement model.

A large number of other works address the problem of uncertainty prediction, including recent papers related to deep learning like e.g. Lakshminarayanan et al. (2017); Jiang et al. (2018); Corbiere et al. (2019). They seek for an uncertainty score [1] defined informally as a real valued summary of an input observation that is predictive of the classification error. These works do not formulate the problem to be solved explicitly as a rejection model. However, most of these works asses performance of their uncertainty scores using evaluation metrics for the rejection models, namely, using the Risk-Coverage (RC) curve and the Area under the RC curve (AuRC).

This article unifies and extends existing formulations of an optimal reject option classifier and proposes theoretically grounded algorithms to learn the classifiers from examples. The main contributions are as follows:

1. We provide necessary and sufficient conditions for an optimal prediction strategy of the bounded-improvement model when the underlying distribution is known. We show that an optimal solution is the Bayes classifier endowed with a rejection strategy, which we call *randomized Bayes selection function*. The randomized Bayes selection function is constructed from the conditional expected risk and two parameters: a decision threshold and an acceptance probability. The strategy rejects prediction when the conditional risk is above the threshold, accepts prediction when it is below the threshold and randomizes with the acceptance probability otherwise. We provide an explicit relation between the decision threshold, the acceptance probability and the target risk.

2. We formulate a bounded-coverage model whose definition is symmetric to the bounded-improvement model. The optimal prediction strategy of the bounded-coverage model

---

1. Some works use term confidence score which is inverse to the uncertainty score utilized in this paper.

minimizes the selective risk under the condition that the coverage is not below a target value. We provide necessary and sufficient conditions for an optimal strategy and we show that the conditions are satisfied by the Bayes classifier endowed with the randomized Bayes selection function. We provide an explicit relation between the decision threshold, the acceptance probability and the target coverage.

3. We define a notion of a *proper uncertainty score* as a function which preserves ordering of the inputs induced by the conditional expected risk. A proper uncertainty score is sufficient for construction of the randomized Bayes selection function. We propose two generic algorithms to learn the proper uncertainty score from examples for an arbitrary black-box classifier. The first is based on regression of the classifier loss. The second is based on minimization of a newly proposed loss function which we call SELEctive classifier learning (SELE) loss. We show that SELE loss is a tight approximation of the AuRC and at the same time amenable to optimization. We prove that both proposed algorithms provide Fisher consistent estimate of the proper uncertainty score. As a proof of concept we apply the proposed algorithms to learn proper uncertainty scores for different prediction problems including classification, ordinal regression and structured output classification. We demonstrate that the algorithm based on the SELE loss minimization learns uncertainty scores which consistently outperform common baselines and work on par with the state-of-the-art methods that are, unlike our algorithm, applicable only to particular prediction models.

Besides the proposed algorithms applicable to learning uncertainty score for an arbitrary classification model, our contributions may have the following uses. Firstly, our analysis shows that despite their different objectives the cost-based, the bounded-improvement and the bounded-coverage rejection models are equivalent in the sense that they lead to the same prediction strategy. Secondly, the explicit characterization of optimal strategies provides a recipe how to construct plug-in rules which has been so far possible only for the cost-based model. That is, any method estimating the class posterior distribution can be turned into an algorithm learning the reject option classifier that solves the bounded-improvement and the bounded-coverage model, respectively. Thirdly, there is a tight connection between the proposed bounded-coverage model and the RC curve. The RC curve represents quality of all solutions of the bounded-coverage model that can be constructed from a pair of a classifier and an uncertainty score. The AuRC is then an expected quality of the reject option classifier constructed from the pair when the target coverage is selected uniformly at random. This connection sheds light on many published methods which do not explicitly define the target objective but use the RC curve and the AuRC as evaluation metrics.

This article is an extension of our previous work published in Franc and Prusa (2019). The major extensions involve introduction of the bounded-coverage model and its analysis, analysis of the learning algorithms including the proof of Fisher consistency, and most of the experiments.

The paper is organized as follows. Section 2 introduces the three rejection models and provides characterization of their optimal solutions. Algorithms to learn a proper uncertainty score from examples are discussed in Section 3. Survey of related literature is given in Section 4. Experimental evaluation of the proposed learning algorithms is provided in Section 5. Section 6 concludes the paper. Proofs of all theorems are deferred to Appendix.

## 2. Reject Option Models and Their Optimal Strategies

Let $\mathcal{X}$ be a set of input observations and $\mathcal{Y}$ a finite set of labels. Let us assume that inputs and labels are generated by a random process with p.d.f. $p(x, y)$ defined over $\mathcal{X} \times \mathcal{Y}$. A goal in the non-reject setting is to find a *classifier* $h\colon \mathcal{X} \to \mathcal{Y}$ with a small *expected risk*

$$R(h) = \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \ell(y, h(x)) \, dx \, ,$$

where $\ell\colon \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$ is a *loss* penalizing the predictions.

The expected risk can be reduced by abstaining from prediction in uncertain cases. To this end, we use a *selective classifier* [2] $(h, c)$ composed of a classifier $h\colon \mathcal{X} \to \mathcal{Y}$ and a *selection function* $c\colon \mathcal{X} \to [0, 1]$. When applying the selective classifier to input $x \in \mathcal{X}$ it outputs

$$(h, c)(x) = \begin{cases} h(x) & \text{with probability} \quad c(x) \, , \\ \text{reject} & \text{with probability} \quad 1 - c(x) \, . \end{cases}$$

In the sequel we introduce three models of an optimal selective classifier: the cost-based, the bounded-improvement and the bounded-coverage model. We characterize optimal strategies of the three models provided the underlying distribution $p(x, y)$ is known.

### 2.1 Cost-based model

Besides the label loss $\ell\colon \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$, let us define a reject loss $\varepsilon \in \mathbb{R}_+$ incurred when a classifiers rejects to predict. The selective classifier $(h, c)$ is evaluated in terms of the expected risk

$$R_B(h, c) = \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \left( \ell(y, h(x)) c(x) + (1 - c(x)) \varepsilon \right) dx \, .$$

**Problem 1 (Cost-based model)** *The optimal selective classifier* $(h_B, c_B)$ *is a solution to the minimization problem*

$$\min_{h,c} R_B(h, c) \, , \tag{1}$$

*where we assume that both minimizers exist.*

The well-known optimal strategy $(h_B, c_B)$ solving Problem 1 reads

$$h_B(x) \in \operatorname*{argmin}_{\hat{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} p(y \mid x) \, \ell(y, \hat{y}) \, , \tag{2}$$

$$c_B(x) = \begin{cases} 1 & \text{if} \quad r^*(x) < \varepsilon \, , \\ \tau & \text{if} \quad r^*(x) = \varepsilon \, , \\ 0 & \text{if} \quad r^*(x) > \varepsilon \, , \end{cases} \tag{3}$$

---

2. The *classifier with a reject option* is usually represented by a single function $h'\colon \mathcal{X} \to \mathcal{Y} \cup \{\text{reject}\}$. We use the decomposition $h'(x) = (h, c)(x)$, and the terminology *selective classifier* from El-Yaniv and Wiener (2010) because we analyze $h$ and $c$ separately.

where

$$r^*(x) = \min_{\hat{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} p(y \mid x) \, \ell(y, \hat{y})$$

is the minimal conditional expected risk associated to input $x$, and $\tau$ is any real number from the interval $[0, 1]$. In the boundary cases when $r^*(x) = \varepsilon$ one can arbitrarily reject or return the best label $h_B(x)$ without affecting the value of the risk $R_B(h, c)$. In turn there always exist a deterministic optimal strategy solving the cost-based model. In the sequel we denote the classifier $h_B$ alone as the *Bayes classifier*.

## 2.2 Bounded-improvement model

One can characterize the selective classifier by two antagonistic quantities: i) the *coverage*

$$\phi(c) = \int_{\mathcal{X}} p(x) \, c(x) \, dx$$

corresponding to the probability that the prediction is accepted [3] and ii) the *selective risk*

$$R_S(h, c) = \frac{\int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \, \ell(y, h(x)) \, c(x) \, dx}{\phi(c)} \, ,$$

defined for non-zero $\phi(c)$ as the expected classification loss on the accepted predictions.

**Problem 2 (Bounded-improvement model)** *Given a* target risk $\lambda > 0$, *the optimal selective classifier* $(h_I, c_I)$ *is a solution to the problem*

$$\max_{h,c} \phi(c) \quad s.t. \quad R_S(h, c) \leq \lambda \, , \tag{4}$$

*where we assume that both maximizers exist.*

**Theorem 1** *Let* $(h, c)$ *be an optimal solution to* (4). *Then,* $(h_B, c)$, *where* $h_B$ *is the Bayes classifier* (2), *is also optimal to* (4).

According to Theorem 1 the Bayes classifier $h_B$ is also optimal for the task (4) defining the bounded-improvement model which is not surprising. Note however that the Bayes classifier is not a unique solution to (4) because the predictions on the reject region $\mathcal{X}_{c(x)=0}$ do not count to the selective risk and hence they can be arbitrary.

Theorem 1 allows to solve the bounded-improvement task (4) in two consecutive steps: First, set $h_I$ to be the Bayes classifier $h_B$. Second, when $h_I$ is fixed, the optimal selection function $c_I$ is obtained by solving the task (4) only with respect to $c$ which boils down to

**Problem 3 (Bounded-improvement model for known $h(x)$)** *Given a classifier* $h(x)$, *the optimal selection function* $c^*(x)$ *is a solution to*

$$\max_{c \in [0,1]^{\mathcal{X}}} \phi(c) \quad s.t. \quad R_S(h, c) \leq \lambda \, . \tag{5}$$

---

3. For a function $f \colon \mathcal{X} \to \mathbb{R}$ and $a \in \mathbb{R} \cup \{\infty\}$, we define $\mathcal{X}_{f(x) \leq a} = \{x \in \mathcal{X} \mid f(x) \leq a\}$, $\mathcal{X}_{f(x) < a} = \{x \in \mathcal{X} \mid f(x) < a\}$, $\mathcal{X}_{f(x)=a} = \{x \in \mathcal{X} \mid f(x) = a\}$, $\mathcal{X}_{f(x) > a} = \{x \in \mathcal{X} \mid f(x) > a\}$, $\mathcal{X}_{f(x) \geq a} = \{x \in \mathcal{X} \mid f(x) \geq a\}$.

Note that Problem 3 is meaningful even if $h$ is not the Bayes classifier $h_B$. In practice we can seek for an optimal selection function $c^*$ for any fixed $h$ which is usually our best approximation of $h_B$ learned from data. We will show that the key concept to characterize an optimal selection function of Problem 3 is the conditional expected risk of $h$ defined as

$$r(x) = \sum_{y \in \mathcal{Y}} p(y \mid x)\, \ell(y, h(x))\,. \tag{6}$$

**Theorem 2** *A selection function $c^* : \mathcal{X} \to [0,1]$ is an optimal solution to Problem 3 if and only if it holds*

$$\int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)dx = \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)dx, \tag{7}$$

$$\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx = \begin{cases} -\frac{\rho(\mathcal{X}_{\overline{r}(x)<b})}{b} & if \quad b > 0\,, \\ \int_{\mathcal{X}_{\overline{r}(x)=0}} p(x)dx & if \quad b = 0\,, \end{cases} \tag{8}$$

$$\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c^*(x)dx = 0\,, \tag{9}$$

*where $\overline{r}(x) = r(x) - \lambda$ measures how much the conditional risk $r(x)$ of the classifier $h(x)$ exceeds the target $\lambda$,*

$$\rho(\mathcal{X}') = \int_{\mathcal{X}'} p(x)\overline{r}(x)\, dx \tag{10}$$

*is the expectation of $\overline{r}(x)$ restricted to inputs in $\mathcal{X}'$, and*

$$b = \sup\left\{ a \mid \rho(\mathcal{X}_{\overline{r}(x)\leq a}) \leq 0 \right\} \geq 0\,. \tag{11}$$

Theorem 2 defines behaviour of an optimal selection function $c^*(x)$ on a partition of the input space $\mathcal{X}$ into three regions $\mathcal{X}_{\overline{r}(x)<b}$, $\mathcal{X}_{\overline{r}(x)=b}$ and $\mathcal{X}_{\overline{r}(x)>b}$. In each region the expected value of $c^*(x)$ is constrained to a particular constant the value of which depends on parameters of the problem. A particular selection function satisfying the optimality condition is given by the following theorem.

**Theorem 3** *Let $r \colon \mathcal{X} \to \mathbb{R}$ be the conditional risk (6) of a classifier $h \colon \mathcal{X} \to \mathcal{Y}$, $\gamma = b + \lambda$ the rejection threshold given by the target risk $\lambda$ and a constant $b$ computed by (11). Then the selection function*

$$c^*(x) = \begin{cases} 1 & if \quad r(x) < \gamma\,, \\ \tau & if \quad r(x) = \gamma\,, \\ 0 & if \quad r(x) < \gamma\,, \end{cases} \tag{12}$$

*where $\tau$ is the acceptance probability given by*

$$\tau = \begin{cases} 1 & if \quad \rho(\mathcal{X}_{r(x)=\gamma}) = 0\,, \\ -\frac{\rho(\mathcal{X}_{r(x)<\gamma})}{\rho(\mathcal{X}_{r(x)=\gamma})} & if \quad \rho(\mathcal{X}_{r(x)=\gamma}) > 0\,, \end{cases} \tag{13}$$

*satisfies the optimality condition of Theorem 2, and hence it is a solution to Problem 3.*

6

The selection function (12) is defined by the conditional risk $r(x)$, the decision threshold $\gamma$ and the acceptance probability $\tau$. The prediction is always accepted when $r(x) < \gamma$ and always rejected when $r(x) > \gamma$. In the boundary cases, when $r(x) = \gamma$, the strategy randomizes and the prediction is accepted with probability $\tau$. The decision threshold is given by $\gamma = b + \lambda$ where $\lambda$ is the target risk in the definition of Problem 3 and $b$ is given by (11). Solving (11) is hard and it requires knowledge of $p(x, y)$. When the probability mass of the set of boundary cases $\mathcal{X}_{r(x)=\gamma}$ is zero, which usually happens in case of continuous $p(x)$, the acceptance probability is $\tau = 1$ and the boundary cases are always accepted, i.e. no randomization is needed.

## 2.3 Bounded-coverage model

In this section we introduce *bounded-coverage model* the definition of which is symmetric to the definition of the bounded-improvement model. Although the problem seems equally useful in practice we are unaware of its formal definition in literature.

**Problem 4 (Bounded-coverage model)** *Given a* target coverage $\omega > 0$, *the optimal selective classifier* $(h_C, c_C)$ *is a solution to the problem*

$$\min_{h,c} R_S(h, c) \quad s.t. \quad \phi(c) \geq \omega \,, \tag{14}$$

*where we assume that both minimizers exist.*

**Theorem 4** *Let* $(h, c)$ *be an optimal solution to (14). Then,* $(h_B, c)$, *where* $h_B$ *is the optimal Bayes classifier (2), is also optimal to (14).*

Theorem 4 ensures that the Bayes classifier $h_B$ is an optimal solution to (14) defining the bounded-coverage model. Note that the solution is not unique as the predictions on $\mathcal{X}_{c(x)=0}$ do not count to the selective risk hence they can be arbitrary. After fixing the classifier $h = h_B$ the search for an optimal selection function leads to:

**Problem 5 (Bounded-coverage model for known** $h(x)$**)** *Given a classifier* $h(x)$ *and a target coverage* $0 < \omega \leq 1$, *the optimal selection function* $c^*(x)$ *is a solution to the problem*

$$\min_{c\in[0,1]^{\mathcal{X}}} R_S(h, c) \quad s.t. \quad \phi(c) \geq \omega \,, \tag{15}$$

*where we assume that the minimizer exists.*

**Theorem 5** *A selection function* $c^* : \mathcal{X} \to [0, 1]$ *is an optimal solution to Problem 5 if and only if it holds*

$$\int_{\mathcal{X}_{r(x)<\beta}} p(x)c^*(x)dx = \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx, \tag{16}$$

$$\int_{\mathcal{X}_{r(x)=\beta}} p(x)c^*(x)dx = \omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx, \tag{17}$$

$$\int_{\mathcal{X}_{r(x)>\beta}} p(x)c^*(x)dx = 0 \,, \tag{18}$$

7

*where*

$$\beta = \inf \left\{ a \mid \int_{\mathcal{X}_{r(x)<a}} p(x)dx \geq \omega \right\} . \tag{19}$$

Theorem 5 defines necessary and sufficient conditions on an optimal solution to Problem 5. A particular selection function satisfying the optimality conditions is given by the following theorem.

**Theorem 6** *Let* $r\colon \mathcal{X} \to \mathbb{R}$ *be the conditional risk (6) of a classifier* $h\colon \mathcal{X} \to \mathcal{Y}$, $1 \geq \omega > 0$ *be a target coverage and* $\beta$ *be the constant computed by (19). Then the selection function*

$$c^*(x) = \begin{cases} 1 & if \quad r(x) < \beta \,, \\ \kappa & if \quad r(x) = \beta \,, \\ 0 & if \quad r(x) > \beta \,, \end{cases} \tag{20}$$

*where* $\kappa$ *is the acceptance probability given by*

$$\kappa = \begin{cases} 0 & if \ \int_{\mathcal{X}_{r(x)=\beta}} p(x)dx = 0 \,, \\ \frac{\omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx}{\int_{\mathcal{X}_{r(x)=\beta}} p(x)dx} & otherwise \,, \end{cases} \tag{21}$$

*satisfies the optimality condition of Theorem 6, and hence it is a solution of Problem 5.*

The selection function (20) is determined by the conditional risk $r(x)$, the decision threshold $\beta$ and the acceptance probability $\kappa$. Both computations of the decision threshold $\beta$, defined by (19), and the acceptance probability $\kappa$, defined by (21), involve integration of $p(x)$. When the probability mass of the set of boundary cases $\mathcal{X}_{r(x)=\beta}$ is zero, the acceptance probability is $\kappa = 0$ and the boundary cases are always rejected without any randomization.

## 2.4 Summary

We have shown that the three rejection models, namely, the cost-based model (c.f. Problem 1), the bounded-improvement model (c.f. Problem 2) and the bounded-coverage model (c.f. Problem 4), share the same class of optimal prediction strategies. An optimal selective classifier $(h, c)$ can be always constructed from the Bayes classifier $h = h_B$ given by (3) and a selection function

$$c_R(x) = \begin{cases} 1 & if \quad r(x) < \alpha \,, \\ \nu & if \quad r(x) = \alpha \,, \\ 0 & if \quad r(x) > \alpha \,, \end{cases} \tag{22}$$

where $r(x)$ is the conditional expected risk of $h(x)$ given by (6), $\alpha \in \mathbb{R}$ is a decision threshold and $\nu \in [0, 1]$ is an acceptance probability. We denote $c_R$ defined by (22) as the *randomized Bayes selection function*. Note that the randomized Bayes selection function $c_R$ is also an optimal solution of the rejection models defined for an arbitrary (i.e. non-Bayes) classifier $h$, that is, an optimal solution of Problem 3 and Problem 5.

The constants $(\nu, \alpha)$ are defined for each rejection model differently and their value depends on parameters of the model (i.e. reject cost $\varepsilon$, target risk $\lambda$ or target coverage $\omega$), the conditional risk $r(x)$ and the underlying distribution $p(x, y)$. For example, in

case of the cost-based model the acceptance threshold $\alpha$ equals to the reject cost $\varepsilon$ and the acceptance probability $\tau$ can be arbitrary. In case of the bounded-improvement and the bounded-coverage model the constants $(\nu, \alpha)$ are defined implicitly via optimization problems and integral equations. In practice $(\nu, \alpha)$ can be tuned on data. For example, Geifman and El-Yaniv (2017) show how to find $\alpha$ from a finite sample such that it is optimal for the bounded-improvement model in PAC sense.

The key component of randomized Bayes selection function $c_R$ is ranking of the inputs $\mathcal{X}$ according to $r(x)$. We introduce notion of a *proper uncertainty score* which is less informative than the conditional risk $r(x)$, yet it is sufficient to construct $c_R$.

**Definition 1** *Let* $h\colon \mathcal{X} \to \mathcal{Y}$ *be a classifier and* $r(x) = \sum_{y \in \mathcal{Y}} p(y \mid x)\, \ell(y, h(x))$ *its conditional expected risk. We say that function* $s\colon \mathcal{X} \to \mathbb{R}$ *is a* proper uncertainty score *of h iff* $\forall (x, x') \in \mathcal{X} \times \mathcal{X}\colon r(x) < r(x') \Rightarrow s(x) < s(x').$

By definition the proper uncertainty score $s(x)$ preserves ordering of the inputs $\mathcal{X}$ induced by the conditional risk $r(x)$. Therefore replacing $r(x)$ by $s(x)$ in function (22), and changing the decision threshold $\alpha$ appropriately, leads to the same optimal selection function.

## 3. Learning the uncertainty function

Assume we want to construct a selective classifier $(h, c)$ solving any of the three rejection models described in Section 2. We have shown that regardless the rejection model, an optimal $h$ is the Bayes classifier $h_B$ given by (2) and an optimal $c$ is the randomized Bayes selection function $c_R$ given by (22). In this section we consider the scenario when $h\colon \mathcal{X} \to \mathcal{Y}$ has been already trained and we want to endow it with $c_R$. The key component of $c_R$ is a proper uncertainty score $s\colon \mathcal{X} \to \mathbb{R}$ satisfying Definition 1. In this section we address problem of learning a proper uncertainty score from examples $\mathcal{T}_n = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \mid i = 1, \ldots, n\}$ assumed to be generated from $n$ i.i.d. random variables with distribution $p(x, y)$. Before describing the algorithms, in Section 3.1 we introduce the notion of Risk-Coverage (RC) curve and Area under Risk-Coverage curve (AuRC). In line with the literature we use the AuRC as a metric to evaluate performance of the learned uncertainty scores. We also show a connection between the RC curve, the AuRC and the bounded-coverage model. In Section 3.2 we describe a *plug-in condition risk rule* and point out that a frequently used Maximum Class Probability rule (MCP) is its special instance. In Section 3.3 we outline a learning approach based on a *loss regression*. In Section 3.4 we introduce a proxy of the AuRC which we call a loss *for SELEctive classifier learning* (SELE). We prove that both proposed methods learn the Fisher consistent estimator of the proper uncertainty score.

### 3.1 Area under Risk Coverage curve

Majority of existing methods that learn selective classifiers output a classifier $h\colon \mathcal{X} \to \mathcal{Y}$ and a deterministic selection function $c\colon \mathcal{X} \to [0, 1]$ defined as [4]

$$c(x) = [\![\, s(x) \le \theta \,]\!]\,, \tag{23}$$

---

4. Note that the deterministic (23) and the randomized Bayes selection function $c_R$ coincide if the acceptance probability is $\nu = 1$, which is a usual case when $p(x)$ is continuous.

where $s\colon \mathcal{X} \to \mathbb{R}$ is an uncertainty score and $\theta \in \mathbb{R}$ a decision threshold. Performance of the pair $(h, s)$ is evaluated by the RC curve obtained after computing the empirical selective risk and the coverage for all settings of the threshold $\theta$. Namely, the computation is as follows. Let us order the examples $\mathcal{T}_n = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \mid i = 1, \ldots, n\}$ according to $s(x)$ so that $s(x_{\pi(1)}) \leq s(x_{\pi(2)}) \leq \cdots \leq s(x_{\pi(n)})$, where $\pi\colon \{1, \ldots, n\} \to \{1, \ldots, n\}$ is a permutation defining the order [5]. Let $L(i, s) = \sum_{j=1}^{i} \ell(y_{\pi_j}, h(x_{\pi_j}))$ be a sum of losses incurred by the classifier $h(x)$ on the examples with uncertainty not higher than the $i$-th highest uncertainty on the examples $\mathcal{T}_n$. The Risk-Coverage curve $\mathcal{C} = \{(\frac{1}{i}L(i, s), \frac{i}{n}) \mid i = 1, \ldots, n\}$ is a set of 2-dimensional points, where the pair $(\frac{1}{i}L(i, s), \frac{i}{n})$ corresponds to the empirical estimate of the selective risk $R_S(h, c)$ and the coverage $\phi(c)$ of a selective classifier $(h, c)$ with the deterministic selective function (23) and decision threshold $\theta = s(x_{\pi_i})$. The area under the RC curve $\mathcal{C}$ is then

$$\mathrm{AuRC}(s, \mathcal{T}_n) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{i} L(i, s) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{i} \sum_{j=1}^{i} \ell(y_{\pi_j}, h(x_{\pi_j})). \tag{24}$$

The value of $\mathrm{AuRC}(s, \mathcal{T}_n)$ can be interpreted as an arithmetic mean of the empirical selective risks corresponding to coverage equidistantly spread over the interval $[0, 1]$ with step $\frac{1}{n}$.

There is a tight connection between RC curve, AuRC and the bounded-coverage model. The RC curve $\mathcal{C}$ represents quality of all admissible solutions of the bounded-coverage model that can be constructed from the pair $(h, s)$ when using the sample $\mathcal{T}_n$ for evaluation. The value of $\mathrm{AuRC}(s, \mathcal{T}_n)$ is an estimate of the expected quality of the selective classifier constructed from the pair $(h, s)$ when the target coverage is selected uniformly at random.

### 3.2 Plug-in conditional risk rule

Prediction models, like e.g. Logistic Regression or Neural Networks learned by cross-entropy loss, use the training set $\mathcal{T}_n$ to learn an estimate $\hat{p}(y \mid x)$ of the class posterior distribution $p(y \mid x)$. The estimate is then used to construct a plug-in Bayes classifier $\hat{h}(x) \in \mathrm{argmin}_{\hat{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \hat{p}(y \mid x) \ell(y, \hat{y})$. Similarly, using $\hat{p}(y \mid x)$ instead of $p(y \mid x)$ in (6) yields the plug-in rule for the conditional risk of classifier $h$ defined as

$$\hat{r}(x) = \sum_{y \in \mathcal{Y}} \hat{p}(y \mid x)\, \ell(y, h(x)) .$$

Provided $p(y \mid x) = \hat{p}(y \mid x), \forall x \in \mathcal{X}, \forall y \in \mathcal{Y}$, the plug-in conditional risk $\hat{r}(x)$ is by definition a proper uncertainty score and it can be used to construct the randomized Bayes selection function $c_R$ which is an optimal rejection strategy for all the three rejection models.

**Example 1 (Maximum Class Probability rule)** *In case of 0/1-loss $\ell(y, y') = [\![y \neq y']\!]$ the plug-in Bayes classifier decides based on the maximum posterior probability $\hat{h}(x) \in \mathrm{argmax}_{y \in \mathcal{Y}} \hat{p}(y \mid x)$ and the plug-in conditional risk rule is*

$$\hat{r}(x) = \sum_{y \in \mathcal{Y}} \hat{p}(y \mid x)\, \ell(y, \hat{h}(x)) = 1 - \max_{y \in \mathcal{Y}} \hat{p}(y \mid x) .$$

---

5. To break ties we use the index of the input in case the scores are the same.

### 3.3 Loss regression

A straightforward approach to learn the uncertainty score is to pose it as a regression problem. The regression function gets an input $x \in \mathcal{X}$ and outputs an estimate of the classification loss $\ell(y, h(x))$. Formally, given a hypothesis space $\mathcal{F} \subset \{s \colon \mathcal{X} \to \mathbb{R}\}$, classifier $h(x)$ and training set $\mathcal{T}_n$, the *loss regression score* $s \colon \mathcal{X} \to \mathbb{R}$ is a solution to $\min_{s \in \mathcal{F}} F_{\mathrm{reg}}(s)$ where

$$F_{\mathrm{reg}}(s) = \frac{1}{n} \sum_{i=1}^{n} \Big( \ell(y_i, h(x_i)) - s(x_i) \Big)^2 \,.$$

It is easy to show that the loss regression score is Fisher consistent estimate of the proper uncertainty score. This amounts to defining the expectation $F_{\mathrm{reg}}(x)$ with respect to i.i.d. generated training set $\mathcal{T}_n$, i.e.,

$$
\begin{aligned}
E_{\mathrm{reg}}(s) &= \mathbb{E}_{\mathcal{T}_n \sim p(x,y)} F_{\mathrm{reg}}(s) \\
&= \int_{\mathcal{X}^n} \sum_{\boldsymbol{y} \in \mathcal{Y}^n} \prod_{i=1}^{n} p(x_i, y_i) \left[ \frac{1}{n} \sum_{i=1}^{n} \Big( \ell(y_i, h(x_i)) - s(x_i) \Big)^2 \right] dx_1 \cdots dx_n \quad (25) \\
&= \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \Big( \ell(y, h(x)) - s(x) \Big)^2 dx \,,
\end{aligned}
$$

and showing that its minimizer is the conditional risk $r(x)$ which is by definition a proper uncertainty score. This is ensured by the following theorem.

**Theorem 7** *The conditional risk $r(x)$ defined by (6) is an optimal solution to $\min\limits_{s \colon \mathcal{X} \to \mathbb{R}} E_{\mathrm{reg}}(s)$.*

### 3.4 Minimization of SELE loss

In this section we define a computationally manageable proxy of AuRC which we call SElective classifier LEarning (SELE) loss. The SELE loss $\Delta_{\mathrm{sele}} \colon \mathbb{R}^n \times \mathcal{X}^n \times \mathcal{Y}^n \to \mathbb{R}_+$ is defined as [6]

$$\Delta_{\mathrm{sele}}(s, \mathcal{T}_n) = \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \ell(y_i, h(x_i)) [\![ s(x_i) \leq s(x_j) ]\!] \,. \quad (26)$$

In contrast to AuRC the computation of $\Delta_{\mathrm{sele}}$ does not require sorting the examples, i.e. we got rid of the permutations that make the evaluation hard. $\Delta_{\mathrm{sele}}$ is still hard to optimize directly due to the step function in its definition. After replacing the step function $[\![ \cdot ]\!]$ by a logistic function we obtain its proxy $\psi_{\mathrm{sele}} \colon \mathbb{R}^n \times \mathcal{X}^n \times \mathcal{Y}^n \to \mathbb{R}_+$ defined as

$$\psi_{\mathrm{sele}}(s, \mathcal{T}_n) = \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \ell(y_i, h(x_i)) \log \big( 1 + \exp(s(x_j) - s(x_i)) \big) \,. \quad (27)$$

The function $\psi_{\mathrm{sele}}(s, \mathcal{T}_n)$ is smooth and convex w.r.t. the argument $s$ and hence it is amenable to optimization. Minimization of $\psi_{\mathrm{sele}}$ is in the core of the proposed learning algorithm which works as follows.

---

6. We assume that the training set $\mathcal{T}_n$ has at least two examples, i.e. $n \geq 2$.

Given a hypothesis space $\mathcal{F} \subset \{s\colon \mathcal{X} \to \mathbb{R}\}$, a classifier $h(x)$ and training set $\mathcal{T}_n$, the *SELE score* $s\colon \mathcal{X} \to \mathbb{R}$ is a solution to the problem $\min_{s \in \mathcal{F}} \psi_{\mathrm{sele}}(s, \mathcal{T}_n)$. We justify the proposed algorithm empirically in Section 5. The theoretical justification is based on the following three arguments:

1. The value of $\Delta_{\mathrm{sele}}$ is a tight approximation of AuRC. In case when value of $s(x)$ is different for each input in $\mathcal{T}_n$, i.e. $s(x_i) \neq s(x_j)$, $\forall i \neq j$, then it holds that

$$\Delta_{\mathrm{sele}}(s, \mathcal{T}_n) \leq \mathrm{AuRC}(s, \mathcal{T}_n) \leq 2 \cdot \Delta_{\mathrm{sele}}(s, \mathcal{T}) .$$

   The first inequality follows from

$$\Delta_{\mathrm{sele}}(s, \mathcal{T}_n) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{n} \hat{L}(i, s) \leq \frac{1}{n} \sum_{i=1}^{n} \frac{1}{i} \hat{L}(i, s) = \mathrm{AuRC}(s, \mathcal{T}_n) .$$

   The second inequality is ensured by Theorem 8.

2. The uncertainty score estimator defined as $\hat{s} \in \mathrm{argmin}_{s \in [0,1]^{\mathcal{X}}} \Delta_{\mathrm{sele}}(s, \mathcal{T}_n)$ is *Fisher consistent*. Namely, Theorem 9 ensures that a population minimizer of $\Delta_{\mathrm{sele}}$ is a proper uncertainty score.

3. The Fisher consistency is preserved even for the smooth proxy $\psi_{\mathrm{sele}}$ the minimization of which is in the core of the proposed algorithm. Namely, Theorem 11 ensures the population minimizer of $\psi_{\mathrm{sele}}$ is a proper uncertainty score.

**Theorem 8** *The inequality* $\mathrm{AuRC}(s, \mathcal{T}_n) \leq 2 \cdot \Delta_{\mathrm{sele}}(s, \mathcal{T}_n)$ *holds true for any* $s\colon \mathcal{X} \to \mathbb{R}$ *and* $\mathcal{T}_n = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \mid i = 1, \ldots, n\}$.

To show the Fisher consistency of $\Delta_{\mathrm{sele}}$ we define its expectation with respect to i.i.d. generated examples $\mathcal{T}_n$, i.e.,

$$
\begin{aligned}
E_{\mathrm{sele}}(s) &= \int_{\mathcal{X}^n} \sum_{\boldsymbol{y} \in \mathcal{Y}^n} \prod_{i=1}^{n} p(x_i, y_i) \Delta_{\mathrm{sele}}(s, \mathcal{T}_n) dx_1 \cdots dx_n \\
&= \frac{1}{n^2} \int_{\mathcal{X}^n} \prod_{i=1}^{n} p(x_i) \sum_{i=1}^{n} \sum_{j=1}^{n} r(x_i) \, [\![ s(x_i) \leq s(x_j) ]\!] \, dx_1 \cdots dx_n \\
&= \frac{1}{n^2} \sum_{i \neq j} \int_{\mathcal{X}} \int_{\mathcal{X}} p(x_i) p(x_j) r(x_i) [\![ s(x_i) \leq s(x_j) ]\!] dx_i \, dx_j \\
&\quad + \frac{1}{n^2} \sum_{i=1}^{n} \int_{\mathcal{X}} \int_{\mathcal{X}} p(x_i) p(x_i) r(x_i) [\![ s(x_i) \leq s(x_i) ]\!] dx_i \, dx_i \\
&= \frac{n^2 - n}{n^2} \int_{\mathcal{X}} p(x) \, r(x) \left( \int_{\mathcal{X}} p(z) \, [\![ s(x) \leq s(z) ]\!] dz \right) dx + \frac{1}{n} \int_{\mathcal{X}} \int_{\mathcal{X}} p^2(x) r(x) dx \, dx .
\end{aligned}
\tag{28}
$$

Minimizers of $E_{\mathrm{sele}}$ are characterized by the following theorem [7].

---

7. $\int_{\mathcal{X}} \int_{\substack{z \neq x \\ s^*(z) = s^*(x)}} f(x,z) dz \, dx$ stands for $\int_{\mathcal{X}} \int_{\mathcal{X}'} f(x,z) dz \, dx$ where $\mathcal{X}' = \{z \in \mathcal{X} \mid z \neq x \wedge s^*(z) = s^*(x)\}$, etc.

**Theorem 9** *A function $s^* : \mathcal{X} \to \mathbb{R}$ is an optimal solution to $\min_{s:\mathcal{X}\to\mathbb{R}} E_{\mathrm{sele}}(s)$ iff*

$$\int_{\mathcal{X}} \int_{\substack{z \neq x \\ s^*(z)=s^*(x)}} \max\{r(x), r(z)\} p(x) p(z) dz\, dx = 0\,, \text{ and} \tag{29}$$

$$\int_{\mathcal{X}} \int_{\substack{r(z)<r(x) \\ s^*(z)>s^*(x)}} (r(x) - r(z))\, p(x) p(z) dz\, dx = 0\,. \tag{30}$$

The conditions (29) and (30) imply that the conditional expectations $\mathbb{E}_{x,z\sim p(x)}[\max\{r(x), r(z)\} \mid z \neq x \land s^*(x) = s^*(z)]$ and $\mathbb{E}_{x,z\sim p(x)}[r(x) - r(z) \mid r(z) < r(x) \land s^*(z) > s^*(x)]$ are both zero. If combined it further implies that a subset of input space $\mathcal{X}' = \{(x, z) \in \mathcal{X} \times \mathcal{X} \mid r(z) < r(x) \land s^*(z) > s^*(x)\}$, on which the order is violated, has probability measure zero. In other words the optimal $s^*(x)$ preserves the ordering induced by $r(x)$ *almost surely* [8].

**Corollary 10** *Any function $s : \mathcal{X} \to \mathbb{R}$ fulfilling*

$$\forall (x, x') \in \mathcal{X} \times \mathcal{X} : x \neq x' \Rightarrow s(x) \neq s(x'), \text{ and} \tag{31}$$

$$\forall (x, x') \in \mathcal{X} \times \mathcal{X} : r(x) < r(x') \Rightarrow s(x) < s(x') \tag{32}$$

*satisfies the optimality conditions of Theorem 9.*

Note that (31) requires the minimizer of $\Delta_{\mathrm{sele}}$ to assign a unique value to each input $x \in \mathcal{X}$ which is not necessary for the score to be proper. Hence the minimizers of $\Delta_{\mathrm{sele}}$ form a subset of all proper uncertainty scores.

To show the Fisher consistency of the smooth proxy $\psi_{\mathrm{sele}}$ we define its expectation with respect to i.i.d. generated examples $\mathcal{T}_n$, i.e.,

$$
\begin{aligned}
E_{\mathrm{proxy}}(s) &= \int_{\mathcal{X}^n} \sum_{\boldsymbol{y}\in\mathcal{Y}^n} \prod_{i=1}^{n} p(x_i, y_i)\, \psi_{\mathrm{sele}}(s, \mathcal{T}_n)\, dx_1 \cdots dx_n \\
&= \frac{n^2 - n}{n^2} \int_{\mathcal{X}} p(x)\, r(x) \left( \int_{\mathcal{X}} p(z)\, \log\left(1 + \exp(s(z) - s(x))\right) dz \right) dx \\
&\quad + \frac{\log(2)}{n} \int_{\mathcal{X}} p(x)\, r(x) dx\,.
\end{aligned}
\tag{33}
$$

We omitted the derivation as it is similar to (28). The key property of the minimizers of $E_{\mathrm{proxy}}$ is stated in the following theorem.

**Theorem 11** *Let $s^* : \mathcal{X} \to \mathbb{R}$ be an optimal solution to $\min_{s:\,\mathcal{X}\to\mathbb{R}} E_{\mathrm{proxy}}(s)$. Then, the condition*

$$\forall (x, x') \in \mathcal{X} \times \mathcal{X} : r(x) < r(x') \Rightarrow s^*(x) < s^*(x')$$

*is satisfied almost surely.*

---

8. This means that the condition can be violated at most on a subset of $\mathcal{X}$ with probability measure zero.

## 4. Related Works

The cost-based rejection model was proposed in Chow (1970) who also provides the optimal strategy in case the distribution $p(x, y)$ is known, analyzes the error-reject trade-off, and proves basic properties of the error-rate and the reject-rate, e.g. that both functions are monotone with respect to the reject cost. The original paper considers the risk with 0/1-loss only. The model with arbitrary classification costs was analyzed e.g. in Tortorella (2000); Schlesinger and Hlaváč (2002).

The bounded-improvement model was coined in Pietraszek (2005). He assumes that a classifier score proportional to the posterior probabilities is known and the task is to find only an optimal decision threshold, which is done numerically based on ROC analysis. The original formulation assumes two classes and 0/1-loss. In this article we consider a straightforward generalization of the bounded-improvement model, see Problem 2, which allows arbitrary number of classes, arbitrary loss and puts no constraint on the class of optimal strategies. We have proved necessary and sufficient conditions on an optimal strategy in case $p(x, y)$ is known. We showed that a particular optimal solution is composed of the Bayes classifier and the randomized Bayes selection function. In addition, we coined the bounded-coverage model, see Problem 4, the definition of which is symmetric to the bounded-improvement model. Although the bounded-coverage model seems equally useful in practice we are unaware of its formal definition in a literature.

There exist another formulations of rejection models for two-class classifiers. For example, in Hanczar and Dougherty (2008) the objective is to maximize the coverage under the constraints that each class has an error rate below a specific threshold. Hence it can be seen as a generalization of the bounded-improvement model. The objective of a rejection model proposed in Lei (2014) is to maximize the total coverage under the constraint that each class has coverage above a specific threshold. None of the two papers analyzes optimal strategies of the corresponding models.

A common approach to construct the selective classifiers for the cost-based model is based on the plug-in rule, which involves learning the class posterior distribution from examples and plugging the distribution to the formula defining the Bayes-optimal strategy (2). In the case of 0/1-loss the plug-in rule rejects based on the maximal class posterior which is denoted as *Maximal Class Probability* (MCP) rule (see Example 1). The MCP rule is probably most frequently used uncertainty score in the literature. The statistically consistency of the plug-in rejection rule is discussed in Herbei and Wegkamp (2006). Fumera et al. (2000) investigate how errors in estimation of the posterior distribution affect effectiveness of the plug-in rule and they also try to improve its performance by using class-specific thresholds. Other methods trying to improve the plug-in rule by tuning multiple thresholds were proposed in Kummert et al. (2016); Fischer et al. (2016). In our work we have derived the optimal strategies for the bounded-improvement model and the newly proposed bounded-coverage model. Our results thus provides a recipe to construct the plug-in rules also for these two rejection models.

There exist many modifications of standard prediction models to learn reject option classifiers for the cost-based model. For example, extensions of the Support Vector Machines to learn a reject option classifier have been studied extensively in Grandvalet et al. (2008); Bartlett and Wegkamp (2008); Yuan and Wegkamp (2010). These works are limited

to two-class problems and 0/1-loss. Learning leads to minimization of a convex surrogate of the cost-based model's objective function. Under some conditions the algorithms are statistically consistent. A boosting algorithm for learning a two-class classifier with reject option is proposed in Cortes et al. (2016). The algorithm minimizes a convex surrogate for the cost-based model and show that the surrogate is calibrated with Bayes solution. Learning prototype-based classifier with rejection option has been addressed in Villman et al. (2016). All these methods require the reject cost to be fixed at the time of learning, and hence changing the cost requires re-training. In contrast, we propose algorithms to learn the proper uncertainty score on top of a pre-trained classifier so that the risk-coverage trade-off can be set by tuning the reject threshold without re-training.

For many prediction models it is easy to devise an ordinal uncertainty score from outputs of the learned (non-reject) classifier. Such strategies are often heuristically based but work reasonably in practice. For example, LeCun et al. (1990) proposed a reject strategy for a Neural Network classifier based on thresholding either the output of the maximally activated unit of the last layer or a difference between the maximal and runner upper output units. Other heuristically based strategies for neural networks were evaluated in Zaragoza and d'Alche Buc (1998); Fisher et al. (2015). In case of Support Vector Machine classifiers (Vapnik, 1998) the trained linear score, proportional to the distance between the input and the decision hyper-plane, is directly used as the uncertainty score (Fumera and Roli, 2002). We denote this approach as the *margin score* and use it as a baseline in our experiments.

Learning of a selective classifier optimal for the bounded-improvement model was discussed in El-Yaniv and Wiener (2010). Their method requires a noisy-free scenario, i.e. they maximize the coverage under the constraint that the selective risk is zero. They provide a characterization of the lower and upper bound of the risk-coverage curves in PAC setting. Geifman and El-Yaniv (2017) assume a selective classifier based on thresholding an uncertainty score and show how to find a decision threshold for the bounded-coverage model which is optimal in PAC sense. They do not address the problem of learning the uncertainty score. We complement their work by showing that the thresholding based selective classifier is an optimal solution when the uncertainty function is proper, and we propose algorithms to learn the proper uncertainty score from examples.

Recent works address uncertainty prediction in context of deep learning (Lakshminarayanan et al., 2017; Jiang et al., 2018; Corbiere et al., 2019). These works do not formulate the problem to be solved explicitly as a rejection model. However, they evaluate their uncertainty scores empirically in terms of the Risk-Coverage curve and the Area under the RC curve which we have shown to be connected with the bounded-coverage model (see Section 3.1). Lakshminarayanan et al. (2017) construct the MCP rule from a posterior distribution modeled as an ensemble of neural networks trained from multiple random initialization. They use adversarial examples to smooth the posterior estimate. Jiang et al. (2018) propose a Trust Score as the ratio between the distance from the test sample to the samples of the nearest class with a different label and the distance to the samples with the same labels as the predicted class. Corbiere et al. (2019) propose a True Class Probability (TCP) score as a measure of prediction confidence. The TCP predicts the value of $p(y^* \mid x)$ where $y^*$ is the ground truth label. They learn a NN, so called ConfNet, by minimizing L2-loss between the ConfNet output and $\hat{p}(y_i \mid x)$ on training examples, where $\hat{p}(y \mid x)$ is the soft-max

distribution trained by standard cross-entropy loss. They show that the TCP empirically outperforms the MCP score and the Trust Score (Jiang et al., 2018) in terms of the AuRC metric. Both Jiang et al. (2018); Corbiere et al. (2019) consider the two-stage approach to learn the uncertainty score similarly to our paper. We show empirically that our proposed SELE score, besides having a theoretical backing and being applicable for a generic classification problem, outperforms the the state-of-the-art TCP score.

## 5. Experiments

In Section 3, we outlined two risk minimization based methods to learn the uncertainty score $s(x)$ for a pre-trained predictor $h(x)$, namely, the algorithm based on i) *loss regression* (Section 3.3) and ii) *minimization of SELE loss* (Section 3.4). We have shown that both methods are Fisher consistent, i.e. they are guaranteed to find the proper score in the idealized setting when the distribution $p(x, y)$ is known (estimation error is zero), the hypothesis space $\mathcal{F}$ contains the proper score (approximation error is zero) and the loss minimizer can be found exactly (optimization error is zero). In this section we evaluate these methods experimentally on real data when all assumptions are presumably violated. We design the experiments so that the optimization and the estimation error are small by using a large number of training examples and linear rules making the loss minimization a convex problem. We compare against the recently proposed True Class Probability (TCP) score (Corbiere et al., 2019) which is learned from examples like the proposed methods. Unlike the proposed methods, the TCP requires the prediction model $h(x)$ to provide an estimate of the posterior $p(y \mid x)$, hence it is not applicable to fully discriminative models like e.g. SVMs. We emphasize that the experiments are meant to be a proof of concept rather than an exhaustive comparison to all existing methods. On the other hand, we are not aware of any other generic method (i.e. being not connected to a particular prediction model) we could compare against.

To demonstrate that the proposed methods are generic we consider three different categories of prediction problems: classification, ordinal regression and structured output classification. For each prediction problem we use several benchmark datasets and frequently used prediction models like the Logistic Regression (LR), three variants of Support Vector Machines (SVMs) and Gradient Boosted Trees. For each prediction model there exists an uncertainty score that is commonly used in practice, like e.g. Maximal Class Probability (MCP) for logistic regression or distance to the decision hyper-plane (a.k.a. margin score) for SVMs. We use these uncertainty scores as additional baselines in our experiments.

### 5.1 Compared methods for uncertainty score learning

In this section we describe three algorithms that use a training set $\mathcal{T}_n = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \mid i = 1, \ldots, n\}$ to learn an uncertainty score $s(x)$ for a pre-trained classifier $h(x)$. We consider linear scores $s_{\boldsymbol{\theta}}(x) = \langle \boldsymbol{\theta}, \boldsymbol{\psi}(x) \rangle$, where $\boldsymbol{\theta} \in \mathbb{R}^m$ are parameters to be learned and $\boldsymbol{\psi} \colon \mathcal{X} \to \mathbb{R}^m$ is a fixed mapping that will be defined for each prediction model separately in the following sections. All evaluated methods are instances of regularized risk minimization framework. In all cases learning leads to an unconstrained minimization of a convex objective $F(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \hat{R}(\boldsymbol{\theta}, \mathcal{T}_n)$, where $C > 0$ is a regularization constant and $\hat{R}(\boldsymbol{\theta}, \mathcal{T}_n)$ is an

empirical risk defined by each method differently. The optimal value of $C$ is selected from $\{0, 1, 10, 100, 1000\}$ based on the minimal value of the AuRC evaluated on a validation set.

### 5.1.1 Regression score

The parameters $\boldsymbol{\theta} \in \mathbb{R}^m$ are learned by minimizing a convex function

$$F_{\mathrm{REG}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{n}\sum_{i=1}^{n} \big(\ell(y_i, h(x_i)) - s_{\boldsymbol{\theta}}(x_i)\big)^2 \,.$$

Minimization of $F_{\mathrm{REG}}(\boldsymbol{\theta})$ is an instance of the ridge regression that can be solved efficiently e.g. by Singular Value Decomposition (SVD).

### 5.1.2 SELE score

Evaluation of the proposed SELE loss $\psi_{\mathrm{sele}}(s, \mathcal{T}_n)$ as defined by (27) requires $\mathcal{O}(n^2)$ operations. To decrease the complexity we approximate its value by splitting the examples into chunks and computing the average loss over the chunks. Namely, the parameters $\boldsymbol{\theta} \in \mathbb{R}^m$ are learned by minimizing a convex function

$$F_{\mathrm{SELE}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{P}\sum_{i=1}^{P} \psi_{\mathrm{sele}}(s, \mathcal{T}_n^k) \,,$$

where $C > 0$ is a regularization constant and $\mathcal{T}_n^1 \cup \mathcal{T}_n^2 \cup \cdots \cup \mathcal{T}_n^P$ is a randomly generated partition of the training set $\mathcal{T}_n$ into $P$ approximately equally sized batches. In all experiments we used $P = \mathrm{round}(n/500)$, i.e. the chunks contain around 500 examples. We minimize $F_{\mathrm{SELE}}(\boldsymbol{\theta})$ by the Bundle Method for Risk Minimization (BMRM) algorithm (Teo et al., 2010) which is set to find a solution whose objective is at most 1% off the optimum [9]. The total computation time of the BMRM algorithm is in order of units of minutes for all dataset using a contemporary PC.

### 5.1.3 True Class Probability score

The TCP (Corbiere et al., 2019) was originally designed for getting uncertainty score on top of a Convolution Neural Network (CNN) trained with cross-entropy loss. The setting we consider here can be seen as the original method applied to a single layer CNN. Namely, let $\hat{p}(y \mid x)$ be an estimate of the posterior distribution, in our experiments provided by the Logistic Regression (a.k.a. single layer NN). The parameters $\boldsymbol{\theta} \in \mathbb{R}^m$ are learned by minimizing a convex function

$$F_{\mathrm{TCP}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{n}\sum_{i=1}^{n} (\hat{p}(y_i \mid x_i) - s_{\boldsymbol{\theta}}(x_i))^2 \,.$$

Minimization of $F_{\mathrm{REG}}(\boldsymbol{\theta})$ is an instance of the ridge regression which we solve by SVD.

---

9. We use $(F_{\mathrm{primal}} - F_{\mathrm{dual}})/F_{\mathrm{primal}} \leq 0.01$ as the stopping condition of the BMRM algorithm.

## 5.2 Benchmark problems

### 5.2.1 CLASSIFICATION

Given real valued features $\boldsymbol{x} \in \mathcal{X} \subseteq \mathbb{R}^d$, the task is to predict a hidden state $y \in \mathcal{Y} = \{1, \ldots, Y\}$ so that the expectation of 0/1-loss $\ell(y, y') = 100 \, [\![y \neq y']\!]$ is as small as possible [10]. We selected 11 classification problems from the UCI repository (Dua and Taniskidou, 2017) and libSVM datasets (Chang and C.J.Lin, 2011). The datasets are summarized in Table 1. We chose the datasets with sufficiently large number of examples relative to the number of features, as we need to learn both the classifier and the uncertainty score and simultaneously keep the estimation error low. Each dataset was randomly split 5 times into 5 subsets, Trn1/Val1/Trn2/Val2/Tst, in ratio 30/10/30/10/20 (up to CODRNA with ratio 25/5/20/20/30 and COVTYPE with ratio 28/20/2/20/30). The subsets Trn1/Val1 were used for learning and tuning the best regularization constant of the classifier $h(x)$. The subsets Trn2/Val2 were used for learning and tuning the regularization constant of the uncertainty score $s(x)$ as described in Section 5.1. All features were normalized to have zero mean and unit variance. The normalization coefficients were estimated using only the Trn1 and Trn2 subsets, respectively. The Tst subset was used solely to evaluate the test performance.

We used two prediction models: Logistic Regression (LR) (Hastie et al., 2009) and Support Vector Machines (SVM) (Vapnik, 1998).

**Logistic Regression** learns parameters $\boldsymbol{\theta}_{\mathrm{LR}} = ((\boldsymbol{w}_y, b_y) \in \mathbb{R}^d \times \mathbb{R} \mid y \in \mathcal{Y})$ of the posterior probabilities $\hat{p}_{\boldsymbol{\theta}}(y \mid \boldsymbol{x}) \approx \exp(\langle \boldsymbol{w}_y, \boldsymbol{x} \rangle + b_y)$ by maximizing the regularized log-likelihood $F_{\mathrm{LR}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{m}\sum_{i=1}^{m} \log\big(\hat{p}_{\boldsymbol{\theta}}(y_i \mid \boldsymbol{x}_i)\big)$. The optimal $C$ was selected from $\{1, 10, 100, 1000\}$ based on the validation classification error. After learning $\boldsymbol{\theta}_{\mathrm{LR}}$ we used the plug-in Bayes classifier $h(\boldsymbol{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} \hat{p}_{\boldsymbol{\theta}_{\mathrm{LR}}}(y \mid \boldsymbol{x})$. As a baseline uncertainty score we use the *plug-in* class conditional risk $\hat{r}(\boldsymbol{x}) = 1 - \hat{p}_{\boldsymbol{\theta}_{\mathrm{LR}}}(h(\boldsymbol{x}) \mid \boldsymbol{x})$. In accordance with the literature we refer to this baseline as the *Maximal Class Probability (MCP) rule*. As shown in Section 3.2, the MCP score is the proper uncertainty score provided the estimate $\hat{p}(y \mid \boldsymbol{x})$ matches the true posterior $p(y \mid \boldsymbol{x})$.

**Support Vector Machines** learn parameters $\boldsymbol{\theta}_{\mathrm{SVM}} = ((\boldsymbol{w}_y, b_y) \in \mathbb{R}^d \times \mathbb{R} \mid y \in \mathcal{Y})$ of the linear classifier $h(\boldsymbol{x}) = \operatorname{argmax}_{y \in \mathcal{Y}}(\langle \boldsymbol{w}_y, \boldsymbol{x} \rangle + b_y)$ by minimizing $F_{\mathrm{SVM}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{m}\sum_{i=1}^{m} \max_{y \in \mathcal{Y}} \big([\![y \neq y_i]\!] + \langle \boldsymbol{w}_y - \boldsymbol{w}_{y_i}, \boldsymbol{x}_i \rangle\big)$. The optimal $C$ was selected in the same way as in case of the LR. As the baseline uncertainty measure we use $s(\boldsymbol{x}) = \max_{y \in \mathcal{Y}}\langle \boldsymbol{w}_y, \boldsymbol{x} \rangle + b_y$. In the binary case $|\mathcal{Y}| = 2$, the setting was $\boldsymbol{\theta}_{\mathrm{SVM}} = (\boldsymbol{w}, b)$, $h(\boldsymbol{x}) = \operatorname{sgn}(\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b)$, $F_{\mathrm{SVM}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{m}\sum_{i=1}^{m} \max\{0, 1 - y_i(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b)\}$ and $s(\boldsymbol{x}) = |\langle \boldsymbol{w}, \boldsymbol{x} \rangle + b|$. In both cases, the value of $s(\boldsymbol{x})$ is proportional to a distance between the input $\boldsymbol{x}$ and the decision boundary. We denote this baseline as the *margin score*.

Given the pre-trained LR or SVM classifier $h(x)$, we apply the methods from Section 5.1 to learn the uncertainty score

$$s_{\boldsymbol{\theta}}(\boldsymbol{x}) = \langle \boldsymbol{w}_{h(\boldsymbol{x})}, \boldsymbol{x} \rangle + b_y \,, \tag{34}$$

where $\boldsymbol{\theta} = ((\boldsymbol{w}_y, b_y) \in \mathbb{R}^d \times \mathbb{R} \mid y \in \mathcal{Y})$ are the parameters to be learned. It is seen that the rule (34) can be re-written as $s_{\boldsymbol{\theta}}(x) = \langle \boldsymbol{\theta}, \boldsymbol{\psi}(x) \rangle$, where $\boldsymbol{\psi} \colon \mathcal{X} \to \mathbb{R}^d$ is an appropriately

---

10. Due to the factor 100 the reported errors correspond to the percentage of misclassified examples.

defined feature map. The form of the score (34) can be justified by noting that its special instance is the margin score which is obtained after substituting $\boldsymbol{\theta}_{\text{SVM}}$ for $\boldsymbol{\theta}$.

### 5.2.2 ORDINAL REGRESSION

The task is to predict a hidden state from $\mathcal{Y} = \{1, \ldots, Y\}$ based on real valued features $\mathcal{X} \subseteq \mathbb{R}^d$. Unlike the classification problem, the hidden states $\mathcal{Y}$ are assumed to be ordered and the goal is to minimize the expectation of the Mean Absolute Error $\ell(y, y') = |y - y'|$. We selected 11 regression problems from UCI repository (Dua and Taniskidou, 2017). The datasets are summarized in Table 1. The real-valued hidden states were discretized into $Y$ bins which are constructed to get uniform class prior. Each dataset was randomly split 5 times into 5 subsets, Trn1/Val1/Trn2/Val2/Tst, in ratio 30/10/30/10/20. We used the same normalization and evaluation protocol as described for the classification benchmarks.

As a prediction model we used a variant of the Support Vector Machine algorithm developed for ordinal regression (Chu and Keerthi, 2005) [11].

**Support Vector Ordinal Regression (SVOR)** learns parameters $\boldsymbol{\theta}_{\text{SVOR}} = (\boldsymbol{w} \in \mathbb{R}^d, (b_1, \ldots, b_{Y-1}) \in \mathbb{R}^{Y-1})$ of the ordinal linear classifier $h(\boldsymbol{x}) = 1 + \sum_{y=1}^{Y-1} [\![\langle \boldsymbol{x}, \boldsymbol{w} \rangle > b_y]\!]$ by minimizing $F_{\text{SVOR}}(\boldsymbol{\theta}) = \frac{C}{2}\|\boldsymbol{\theta}\|^2 + \frac{1}{m}\sum_{i=1}^{m} \big( \sum_{y=1}^{y_i-1} \max(0, 1 - \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b_y) + \sum_{y=y_i}^{Y-1} \max(0, 1 + \langle \boldsymbol{x}_i, \boldsymbol{w} \rangle - b_y))$. The optimal C was selected from $\{1, 10, 100, 1000\}$ based on the validation MAE. The ordinal classifier can be thought of as a standard linear classifier composed of parallel decision hyper-planes. Similarly to the standard SVM, we use $s(\boldsymbol{x}) = \min_{y \in \{1, \ldots, Y-1\}} |\langle \boldsymbol{x}, \boldsymbol{w} \rangle - b_y|$ as a baseline uncertainty score. The value of $s(\boldsymbol{x})$ is proportional to the distance of $\boldsymbol{x}$ to the closest hyper-plane hence we also denote it as the *margin score*. When learning the uncertainty from examples we use the parametrization (34).

### 5.2.3 STRUCTURED OUTPUT CLASSIFICATION

Given an RGB image $\boldsymbol{x} \in \mathcal{X} = \{0, \ldots, 255\}^{W \times H \times 3}$ capturing a human face, the task is to predict a pixel positions of 68 landmarks $\boldsymbol{y} = (\boldsymbol{l}_1, \ldots, \boldsymbol{l}_{68}) \in \mathcal{Y} = (\{1, \ldots, W\} \times \{1, \ldots, H\})^{68}$ corresponding to contours of eyes, mouth, nose, etc. We use the 300-W dataset and the associated evaluation protocol which was created by organizers of landmark detection challenge (Sagonas et al., 2016). The 300-W dataset contains 5,807 faces each annotated with 68 landmarks. The faces are split into 3,484 training, 1,161 validation and 1,162 test examples. The prediction accuracy is measured in terms of normalized average localization error $\ell(\boldsymbol{y}, \hat{\boldsymbol{y}}) = \frac{100}{\text{iod}(\boldsymbol{y})} \sum_{i=1}^{68} \|\boldsymbol{l}_i - \hat{\boldsymbol{l}}_i\|$, where $\text{iod}(\boldsymbol{y})$ is the inter-ocular distance computed from the ground-truth landmark positions $\boldsymbol{y}$.

As the structured classifier $h(x)$ we use the landmark detector from DLIB package (King, 2009). The detector predicts the landmark positions based on HOG descriptors (Dalal and Triggs, 2005) of the input image using an ensemble of regression trees that are trained by gradient boosting (Kazemi and Sullivan, 2014). The DLIB landmark detector has been widely used by developers due to its robustness and exceptional speed even on a low-end hardware. The detector does not provide any measure of prediction uncertainty and it is unclear how to derive it from outputs of the regression trees. A commonly used uncertainty score for

---

11. (Chu and Keerthi, 2005) introduced two variants of SVOR algorithm. We use so called SVOR with *implicit constraints* which is designed for minimization of the MAE loss.

face recognition related prediction problems is a score function of a face detector. The face detector score is an output of a binary classifier trained to distinguish face from non-face images. The value of the score is high for well looking prototypical faces and low for corrupted or "difficult" faces. As a baseline we use the score of the DLIB face detector which is a linear SVM classifier on top of HOG descriptors extracted from the image.

When learning the uncertainty score from examples the score is a linear regressor $s_{\boldsymbol{\theta}}(x) = \langle \boldsymbol{\theta}, \boldsymbol{\psi}(x) \rangle$ on top of a feature vector $\boldsymbol{\psi}(x) \in \mathbb{R}^{2,448}$ which is a concatenation of HOG descriptors extracted from the input facial image $x$ along landmark positions predicted by the landmark detector $h(x)$. The DLIB detector uses the same features to predict the landmark positions hence the extra computational time required to evaluate the uncertainty score is neglectable.

| Classification problems | | | | Ordinal regression problems | | | |
|---|---|---|---|---|---|---|---|
| dataset | examples | feat | cls | dataset | examples | feat | cls |
| AVILA | 20,867 | 10 | 12 | ABALONE | 4,177 | 10 | 19 |
| CODRNA | 331,152 | 8 | 2 | BANK | 8,192 | 32 | 10 |
| COVTYPE | 581,012 | 54 | 7 | BIKESHARE | 17,379 | 11 | 10 |
| IJCNN | 49,990 | 22 | 2 | CALIFORNIA | 20,640 | 8 | 10 |
| LETTER | 20,000 | 16 | 26 | CCPP | 9568 | 4 | 10 |
| MARKETING | 45,211 | 51 | 2 | CPU | 8192 | 21 | 10 |
| PENDIGIT | 10,992 | 16 | 10 | FACEBOOK | 50,993 | 53 | 10 |
| PHISHING | 11,055 | 68 | 2 | GPU | 24,1600 | 14 | 10 |
| SATTELITE | 6,435 | 36 | 6 | METRO | 48,204 | 30 | 10 |
| SENSORLESS | 58,509 | 48 | 11 | MSD | 499,671 | 90 | 41 |
| SHUTTLE | 58,000 | 9 | 7 | SUPERCOND | 21,263 | 81 | 10 |

Table 1: Summary of 11 classification problems (left) and 11 ordinal regression problems (right) selected from UCI repository (Dua and Taniskidou, 2017) and libSVM datasets (Chang and C.J.Lin, 2011). The table shows the total number of examples, the number of features and the number of classes.

## 5.3 Results

### 5.3.1 Classification problems

For both classification models, LR and SVM, we recorded the test risk of the classifier $h(x)$ and the AuRC computed from $h(x)$ and uncertainty score $s(x)$ produced by the corresponding method under evaluation. In case of LR, we compare the baseline MCP score (sec 5.2.1) and scores learned from examples including the state-of-the-art TCP score (sec 5.1.3) and the two proposed SELE score (sec 5.1.2) and Regression (REG) score (sec 5.1.1). In case of SVM, we compare the baseline Margin score (sec 5.2.1) against the proposed SELE and REG scores. Note that TCP score is not applicable for SVM classifier as it does not provide an estimate of the posterior probability $p(y \mid x)$. The results are summarized in Table 2.

For each dataset we rank the compared methods according to the AuRC [12]. Following the methodology of Demšar (2006) we summarize performance of each method by its average rank and use the Friedman test and the post-hoc Nemenyi test to analyze significance of the results.

- We used the Friedman test to check whether the measured average ranks are significantly different from the mean rank. The null hypothesis states that the compared scores are equivalent so that their average ranks should be equal. In both cases the null hypothesis is rejected for p-value 0.05, i.e., the *performance of compared methods is significantly different.*

- We used post-hoc Nemenyi test for pair-wise comparison. For each pair of methods it checks whether their average ranks are significantly different. In case of LR, when we compare $K = 4$ methods using $N = 11$ datasets, the critical difference for p-value 0.10 is $CD = 1.26$. By comparing the average ranks we conclude that *SELE score performs significantly better than MCP score and REG score.* In case of SVM, when we compare $K = 3$ methods using $N = 11$ dataset, the critical difference for p-value 0.10 is $CD = 0.98$. We conclude that *SELE performs significantly better than REG score and Margin score.* The data is not sufficient to reach any conclusion about other pairwise comparisons. The result of the Nemenyi test is visualized in Figure 1.

We further computed relative improvement gained by using the scores SELE, TCP and REG, that are all learned from examples, with respect to the baseline scores derived from the classifier output, i.e. MCP in case of LR and Margin score in case of SVM. The results are summarized in Figure 2. It is seen that the MCP uncertainty computed from the estimated $p(y \mid x)$ constitutes much stronger baseline than the Margin score of fully discriminative SVM model. The relative improvement of scores learned on top the LR is only moderate in contrast to the SVM classifier where the improvements are more significant and consistent. It is also seen that on majority of datasest the performance of the learned scores is similar taking into account the statistical error of the AuRC estimate. It is also worth mentioning that results of SELE score have the lowest variance of the AuRC estimates as seen from error bars in Figure 2.



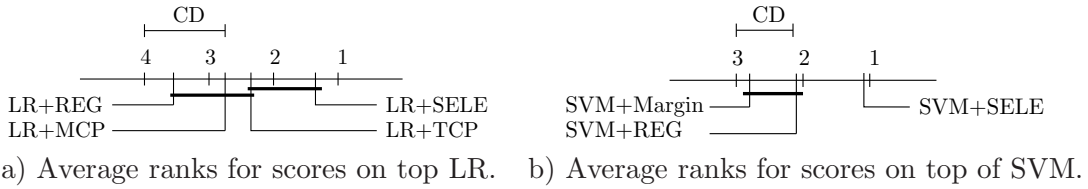a) Average ranks for scores on top LR.   b) Average ranks for scores on top of SVM.

Figure 1: Comparison of all uncertainty scores against each other with the Nemenyi test. The test is computed separately for the LR classifier (4 scores compared) and the SVM classifier (3 scores compared). The figures show the average ranks for each score and the critical distance (CD). Groups of scores that are not significantly different at $p$-value 0.10 are connected.

---

12. The score with smallest AuRC is ranked 1, the second smallest 2 and so on.

| | LR+MCP<br>AuRC | LR+SELE<br>AuRC | LR+REG<br>AuRC | LR+TCP<br>AuRC | LR<br>R@100 |
|---|---|---|---|---|---|
| AVILA | 27.18±0.55 | **25.79±0.44** | 26.62±0.74 | 26.85±0.78 | 43.71±0.42 |
| CODRNA | 0.88±0.05 | **0.65±0.03** | 0.82±0.06 | 0.78±0.04 | 4.81±0.08 |
| COVTYPE | **16.49±0.06** | 17.58±0.07 | 17.62±0.09 | 17.19±0.07 | 27.56±0.17 |
| IJCNN | 1.26±0.04 | **1.00±0.03** | 1.16±0.08 | 1.14±0.06 | 7.54±0.15 |
| LETTER | 7.43±0.40 | **6.42±0.34** | 7.44±0.59 | 6.71±0.42 | 23.32±0.60 |
| MARKETING | 2.60±0.31 | **1.88±0.11** | 1.97±0.12 | 1.90±0.11 | 9.88±0.29 |
| PENDIGIT | **0.69±0.04** | 1.55±0.19 | 1.97±0.55 | 1.47±0.39 | 5.29±0.40 |
| PHISHING | 0.76±0.10 | **0.75±0.10** | 0.91±0.31 | 0.85±0.25 | 6.29±0.44 |
| SATTELITE | 3.83±0.26 | **3.68±0.27** | 4.93±1.07 | 4.52±0.85 | 15.06±0.46 |
| SENSORLESS | 2.03±0.11 | **1.82±0.08** | 2.69±0.09 | 2.37±0.22 | 8.23±0.45 |
| SHUTTLE | 0.59±0.09 | **0.26±0.07** | 1.24±0.51 | 0.58±0.13 | 3.36±0.25 |
| average rank | 2.73 | 1.36 | 3.55 | 2.36 | |

(a) Uncertainty scores on top of LR classifier.

| | SVM+MARGIN<br>AuRC | SVM+SELE<br>AuRC | SVM+REG<br>AuRC | SVM<br>R@100 |
|---|---|---|---|---|
| AVILA | 31.65±0.83 | **25.26±0.67** | 25.95±0.75 | 43.34±0.70 |
| CODRNA | 0.89±0.05 | **0.65±0.03** | 0.82±0.05 | 4.78±0.08 |
| COVTYPE | 25.71±0.81 | 17.79±0.21 | **17.77±0.14** | 27.41±0.11 |
| IJCNN | 1.40±0.04 | **1.01±0.04** | 1.18±0.08 | 7.56±0.16 |
| LETTER | 10.20±0.22 | **6.05±0.65** | 7.15±0.65 | 22.06±0.69 |
| MARKETING | 2.24±0.20 | **1.97±0.10** | 2.04±0.20 | 10.48±0.39 |
| PENDIGIT | 2.79±0.40 | **1.57±0.21** | 2.16±0.43 | 4.88±0.57 |
| PHISHING | 0.84±0.12 | **0.72±0.12** | 0.90±0.30 | 6.37±0.44 |
| SATTELITE | 4.75±0.60 | **3.82±0.27** | 5.44±0.68 | 15.36±0.37 |
| SENSORLESS | 3.68±0.20 | **1.56±0.08** | 2.46±0.29 | 6.92±0.17 |
| SHUTTLE | 1.31±0.47 | **0.24±0.07** | 0.55±0.15 | 2.02±0.15 |
| average rank | 2.82 | 1.09 | 2.09 | |

(b) Uncertainty scores on top of SVM classifier.

Table 2: Performance of the uncertainty scores on 11 classification problems. The scores are constructed on top of the LR classifier and the SVM classifier measured in terms of AuRC. For each score we show the mean and the standard deviation of the test AuRC computed over 5 random splits. We compare the performance of scores learned from examples (SEL, REG, TCP) and the baseline scores derived from the classifiers output (MCP and Margin score). The last column shows the risk of the base (non-selective) classifier. All the values correspond to percentage of misclassification. The best results for each dataset are shown in bold. The last row shows the average rank.
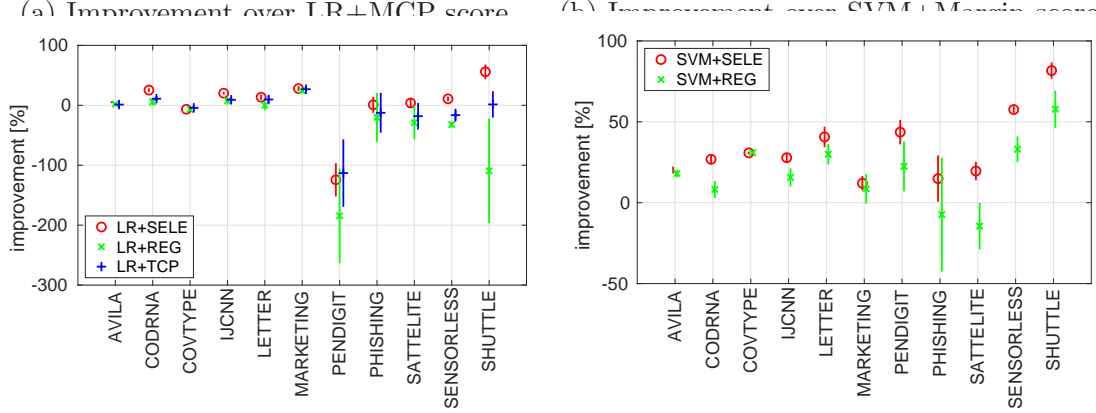
Figure 2: Relative improvement gained by using the uncertainty scores (SELE, REG and TCP) that are learned from examples over the baseline scores (MCP for LR and Margin score for SVM) constructed from the classifier output. The relative improvement is computed as $100 \times (\text{AuRC}_{\text{baseline}} - \text{AuRC}_{\text{method}}) / \text{AuRC}_{\text{baseline}}$. We show the mean and the standard deviation (error bar) of the relative improvement computed over the random 5 splits.

### 5.3.2 ORDINAL REGRESSION

In case of SVOR classifier we compared the baseline Margin score (sec 5.2.2) and the scores learned from examples including SELE (sec 5.1.2) and REG score (sec 5.1.1). We used exactly the same evaluation protocol as for the classification task, however, instead of 0/1-loss the errors were evaluated by the MAE loss. The results are summarized in Table 3.

We again ranked the methods according to the AuRC and summarized their performance by the average rank:

- We applied the Friedman test checking whether the measured average ranks are significantly different from the mean rank. The null hypothesis, stating that the compared scores are equivalent, is rejected for p-value 0.05, i.e., the *performance of compared methods is significantly different.*

- We used post-hoc Nemenyi test to check for each pair whether their average ranks are significantly different. Considering $K = 3$ compared methods using $N = 11$ dataset yields the critical difference for p-value 0.10 is $CD = 0.98$. We conclude that *both SELE and REG scores are significantly better than the baseline Margin score.* The data is not sufficient to reach any conclusion about comparisons of SELE and REG. The result of the Nemenyi test is visualized in Figure 3(a).

The relative improvement gained by using SELE and REG scores learned from examples w.r.t. the baseline Margin score is shown in Figure 3(b). It is seen that the performance of the learned scores is similar and that they consistently outperform the baseline by a significant margin.

| | SVOR+MARGIN | SVOR+SELE | SVOR+REG | SVOR |
|---|---|---|---|---|
| | AuRC | AuRC | AuRC | R@100 |
| CALIFORNIA | 0.98±0.03 | **0.82±0.02** | 0.84±0.02 | 1.18±0.01 |
| ABALONE | 1.48±0.10 | **1.19±0.09** | 1.21±0.05 | 1.54±0.02 |
| BANK | 1.07±0.04 | 0.99±0.04 | **0.98±0.03** | 1.50±0.03 |
| CPU | 0.41±0.01 | **0.36±0.02** | 0.36±0.02 | 0.64±0.03 |
| BIKESHARE | 1.60±0.07 | **1.25±0.01** | 1.27±0.01 | 1.70±0.03 |
| CCPP | 0.46±0.02 | **0.41±0.02** | 0.42±0.02 | 0.58±0.02 |
| FACEBOOK | 0.51±0.01 | 0.37±0.01 | **0.36±0.01** | 1.11±0.01 |
| GPU | 1.43±0.02 | **0.85±0.03** | 0.86±0.03 | 1.49±0.02 |
| METRO | 2.20±0.07 | **1.97±0.01** | 1.98±0.03 | 2.37±0.03 |
| MSD | 6.23±0.07 | 4.26±0.03 | **4.25±0.03** | 6.22±0.03 |
| SUPERCONDUCT | 0.98±0.02 | **0.75±0.01** | 0.77±0.01 | 1.07±0.01 |
| average rank | 3.00 | 1.27 | 1.73 | |

Table 3: Performance of the uncertainty scores on 11 ordinal regression problems. The scores are constructed on top of the SVOR classifier. For each score we show the mean and the standard deviation of the test AuRC computed over 5 random splits. We compare the performance of scores learned from examples (SEL, REG) and the baseline Margin score derived from the SVOR classifier output. The last column shows the risk of the base (non-selective) classifier. All the values correspond to the Mean Absolute Error (MAE). The best results for each dataset are shown in bold. The last row shows the average rank.

### 5.3.3 Structured Output Classification

We trained SELE score (sec 5.1.2) and REG score (sec 5.1.1) on top of the DLIB detector and compared them with the baseline which uses the DLIB face detector score (sec 5.2.3) as an uncertainty measure. The Risk-Coverage curves of the three methods and their corresponding AuRC are shown in Figure 4(a). Both the learned scores, SELE and REG, are significantly better than the baseline face detector score. The SELE is slightly better than REG score. The largest difference between the three scores are seen for low values of coverage where SELE most outperforms the other two methods. High selective risk for low values of coverage means that faces with very bad landmark predictions are assigned the lowest uncertainty scores. SELE score does not suffer from this problem. This can be seen in Figure 5 where we show examples of 10 test faces with the lowest uncertainty and the highest uncertainty predicted by SELE.

Unlike the experiments in the previous section, the number of parameters to be learned ($m = 2,448$) relative to the number of training examples ($n = 3,484$) is much higher. To see whether the number of examples is sufficient we trained SELE and REG scores from increasingly bigger training set. Figure 4(b) shows the test AuRC as a function of the the number of training examples. It is seen that AuRC of the SELE is not yet saturated and

a) Average rank for scores on top of SVOR.

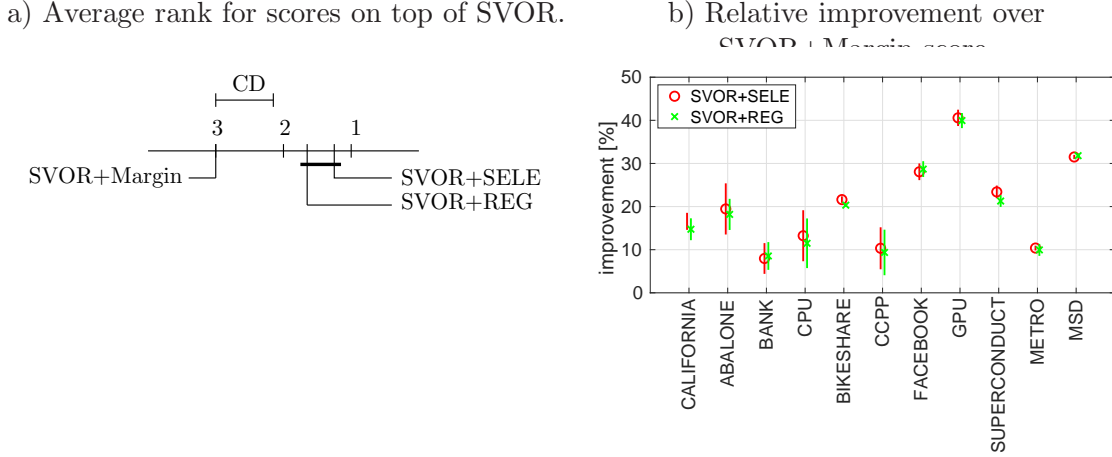b) Relative improvement over SVOR+Margin score.



Figure 3: Statistics derived from results obtained on 11 ordinal regression problems. Figure (a) shows pair-wise comparison of uncertainty scores with the Nemenyi test. The figure shows the average ranks for each score and the critical distance (CD). Groups of scores that are not significantly different at $p$-value 0.10 are connected. Figure (b) shows relative improvement gained by using the SELE and REG uncertainty scores learned from examples over the baseline Margin score.

would most likely converge to a significantly lower value relative to the REG score provided 300-W dataset had more training examples.

## 6. Conclusions

The standard cost-based rejection model introduced by Chow (1970) requires explicit definition of the rejection cost which is difficult in applications when the reject cost and the label loss have different nature or physical units. Pietraszek (2005) proposed the bounded-improvement model which avoids the problem by defining an optimal prediction strategy in terms of the coverage and the selective risk. We have coined a symmetric definition, the bounded-coverage model, which is useful when defining the target coverage is easier than defining the target selective risk. Our main result is a formal proof that despite their different objectives the three rejection models are equivalent in the sense that they lead to the same prediction strategy: the Bayes classifier and the randomized Bayes selection function. Thanks to the common optimal solution it is possible to convert between parameters of different rejection models. For example, for any target risk defining the bounded-improvement model there exists a corresponding reject cost so that both models have the same optimal strategy.

The explicit characterization of the optimal strategies provides a recipe to construct plug-in rules solving the bounded-improvement and the bounded-coverage models. Any method estimating the class posterior probabilities can be thus turned into an algorithm for learning the selective classifier that solves the bounded-improvement and the bounded-coverage model.
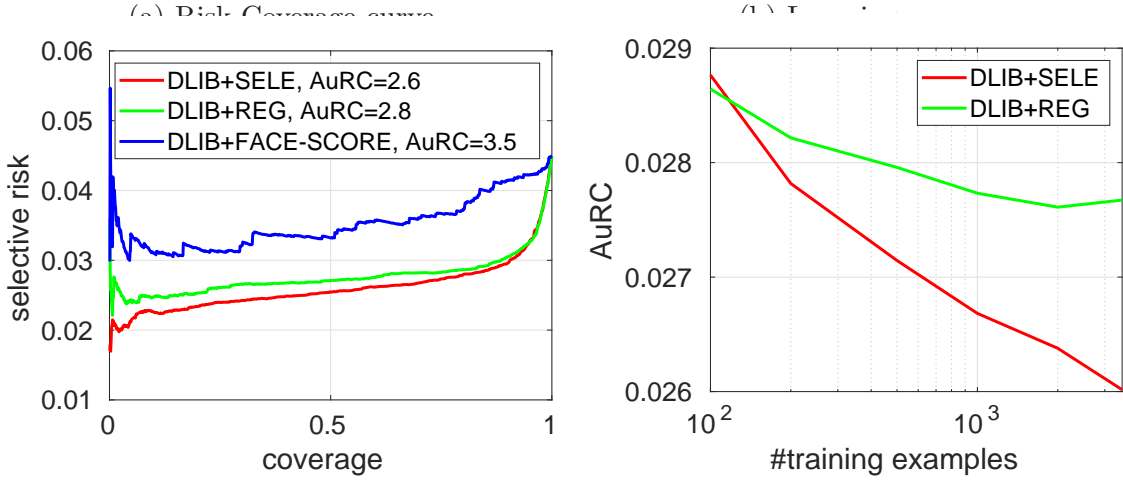
Figure 4: Evaluation of uncertainty scores on top of DLIB landmark detector using the 300-W benchmark. Figure (a) shows the RC curve and AuRC of computed from the predictions of DLIB detector endowed with the three compared uncertainty scores: the proposed SELE and REG scores and the DLIB face detector score used as a baseline. Figure (b) shows the test AuRC for SELE and REG scores as a function of the number of training examples.

We have defined a notion of a proper uncertainty score which is sufficient to construct the randomized Bayes selection function. We proposed two algorithms to learn a proper uncertainty score from examples for a given classifier. We have shown that both algorithms provide a Fisher consistent estimate of the proper uncertainty score. As a proof of concept we evaluated the proposed algorithms on different types of prediction problems. We have shown that the proposed algorithm based on minimization of the SELE loss outperforms existing approaches tailored for a particular prediction model and it works on par with the recently published state-of-the-art TCP score (Corbiere et al., 2019). Unlike the TCP score which requires the classifier to output the class posterior probabilities the proposed algorithms are applicable for an arbitrary black-box classifier.

We have drawn a connection between the proposed bounded-coverage model and the RC curve. Namely, the RC curve represents quality of all admissible solutions of the bounded-coverage model that can be constructed from a pair of classifier and uncertainty score. The AuRC is then the expected quality of the selective classifier constructed from the pair when the target coverage is selected uniformly at random. This connection sheds light on many published methods which do not explicitly define the target objective but use the RC curve and the AuRC as evaluation metrics.

Finally let us mentioned some topics for future work. Firstly, the proposed algorithms consider two-stage scenario when the classifier and the uncertainty score are learned separately from independent training sets. Although the scenario is useful in practice, an algorithm learning the classifier and the uncertainty score simultaneously from a single training set constitute an interesting topic to be solved. Secondly, we have shown how

to learn the proper uncertainty score but have not discussed how to set up the decision threshold and the acceptance probability that are also needed to construct the selective classifier. It is straightforward to tune these parameters on empirical data using the RC curve. Analysis of the generalization error of this empirical approach is an open issue which has been solved only for the decision threshold of the bounded-improvement model by Geifman and El-Yaniv (2017). Thirdly, the empirical evaluation is limited to uncertainty scores linear in the parameters to be learned. Efficient implementations of the algorithms applicable to non-linear models, like e.g. the neural networks, is an another topic left for future.

## Acknowledgments

Figure 5: Figure shows examples of 10 test faces from 300-W database with the lowest and 10 faces with the highest value of the SELE uncertainty score. The ground-truth landmark positions (red) and the landmark positions predicted by DLIB detector (blue) are superimposed into the image. The image title shows the rank induced by the SELE score and the normalized localization error which is used as the classification loss in this application.

## Appendix A. Proofs of theorems from Section 2

### A.1 Proof of Theorem 1

The Bayes classifier reads

$$h_B(x) \in \underset{\hat{y} \in \mathcal{Y}}{\mathrm{argmin}} \sum_{y \in \mathcal{Y}} p(y \mid x) \, \ell(y, \hat{y}) \tag{2}$$

**Problem 2 (Bounded-improvement model)** *Given a* target risk $\lambda > 0$, *the optimal selective classifier* $(h_I, c_I)$ *is a solution to the problem*

$$\max_{h,c} \phi(c) \quad s.t. \quad R_S(h, c) \leq \lambda, \tag{4}$$

*where we assume that both maximizers exist.*

**Theorem 1** *Let* $(h, c)$ *be an optimal solution to (4). Then,* $(h_B, c)$, *where* $h_B$ *is the Bayes classifier (2), is also optimal to (4).*

**Proof** It is sufficient to show that $(h_B, c)$ is feasible to (4), i.e., that $R_S(h_B, c) \leq \lambda$. Then $(h_B, c)$ attains the same maximum objective value $\phi(c)$ as $(h, c)$. Derive

$$
\begin{aligned}
R_S(h_B, c) &= \frac{1}{\phi(c)} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \, \ell(y, h_B(x)) \, c(x) \, dx \\
&= \frac{1}{\phi(c)} \int_{\mathcal{X}} p(x) c(x) \left( \sum_{y \in \mathcal{Y}} p(y \mid x) \, \ell(y, h_B(x)) \right) dx \\
&\overset{(2)}{\leq} \frac{1}{\phi(c)} \int_{\mathcal{X}} p(x) c(x) \left( \sum_{y \in \mathcal{Y}} p(y \mid x) \, \ell(y, h(x)) \right) dx \\
&= \frac{1}{\phi(c)} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \, \ell(y, h(x)) \, c(x) \, dx \\
&= R_S(h, c) \leq \lambda.
\end{aligned}
$$

$\blacksquare$

### A.2 Proof of Theorem 2

The presented proof of the theorem uses Lemmas 13 and 14, both derived based on Lemma 12 bellow.

**Lemma 12** *For a set* $\mathcal{X}$, *let* $f : \mathcal{X} \to \mathbb{R}_+{}^{13}$ *and* $g : \mathcal{X} \to \mathbb{R}$ *be measurable functions such that* $\int_{\mathcal{X}} f(x) dx > 0$ *and* $g(x) > 0$ *for all* $x \in \mathcal{X}$. *Then it holds* $\int_{\mathcal{X}} g(x) f(x) dx > 0$.

---

13. We use $\mathbb{R}$, $\mathbb{R}_+$ and $\mathbb{N}_+$ to denote the set of real numbers, non-negative real numbers and positive integers, respectively.

**Proof** For $n \in \mathbb{N}_+$, define functions

$$f_n(x) = \begin{cases} f(x) & \text{if } g(x) \geq \frac{1}{n}, \\ 0 & \text{otherwise.} \end{cases}$$

The sequence $\{f_n\}_{n=1}^{\infty}$ is monotone and converges to $f$. Using the monotone convergence theorem (Stein and Shakarchi, 2009), derive

$$0 < \int_{\mathcal{X}} f(x)dx = \int_{\mathcal{X}} \lim_{n \to \infty} f_n(x)dx = \lim_{n \to \infty} \int_{\mathcal{X}} f_n(x)dx \,.$$

This means that there is a $k \in \mathbb{N}_+$ such that $\int_{\mathcal{X}} f_k(x)dx > 0$, hence we conclude

$$\int_{\mathcal{X}} g(x)f(x)dx \geq \int_{\mathcal{X}} g(x)f_k(x)dx \geq \int_{\mathcal{X}} \frac{1}{k} f_k(x)dx > 0.$$

■

**Lemma 13** *For a set $\mathcal{X}$, let $f : \mathcal{X} \to \mathbb{R}_+$ and $g : \mathcal{X} \to \mathbb{R}$ be measurable functions such that $\int_{\mathcal{X}} f(x)dx > 0$ and $g(x) > b$ for all $x \in \mathcal{X}$ and some $b \in \mathbb{R}$. Then it holds $\int_{\mathcal{X}} g(x)f(x)dx > b \int_{\mathcal{X}} f(x)dx$.*

**Proof** By Lemma 12, we have

$$\int_{\mathcal{X}} (g(x) - b)f(x)dx > 0,$$

thus

$$\int_{\mathcal{X}} g(x)f(x)dx = \int_{\mathcal{X}} (g(x) - b)f(x)dx + \int_{\mathcal{X}} bf(x)dx > b \int_{\mathcal{X}} f(x)dx.$$

■

**Lemma 14** *For a set $\mathcal{X}$, let $f : \mathcal{X} \to \mathbb{R}_+$ and $g : \mathcal{X} \to \mathbb{R}$ be measurable functions such that $\int_{\mathcal{X}} g(x)f(x)dx > 0$ and $g(x) < 1$ for all $x \in \mathcal{X}$. Then it holds $\int_{\mathcal{X}} f(x)dx > \int_{\mathcal{X}} g(x)f(x)dx$.*

**Proof** $\int_{\mathcal{X}} g(x)f(x)dx > 0$ implies $\int_{\mathcal{X}} f(x)dx > 0$. Since it holds $\forall x \in \mathcal{X} : (1 - g(x)) > 0$, Lemma 12 yields

$$0 < \int_{\mathcal{X}} (1 - g(x))f(x)dx = \int_{\mathcal{X}} f(x)dx - \int_{\mathcal{X}} g(x)f(x)dx,$$

and $\int_{\mathcal{X}} f(x)dx > \int_{\mathcal{X}} g(x)f(x)dx$ is obtained as a direct consequence. ■

**Problem 3 (Bounded-improvement model for known $h(x)$)** *Given a classifier $h(x)$, the optimal selection function $c^*(x)$ is a solution to*

$$\max_{c \in [0,1]^{\mathcal{X}}} \phi(c) \quad s.t. \quad R_S(h,c) \le \lambda. \tag{5}$$

**Theorem 2** *A selection function $c^* : \mathcal{X} \to [0,1]$ is an optimal solution to Problem 3 if and only if it holds*

$$\int_{\mathcal{X}_{\overline{r}(x) < b}} p(x) c^*(x) dx = \int_{\mathcal{X}_{\overline{r}(x) < b}} p(x) dx, \tag{7}$$

$$\int_{\mathcal{X}_{\overline{r}(x) = b}} p(x) c^*(x) dx = \begin{cases} -\dfrac{\rho(\mathcal{X}_{\overline{r}(x) < b})}{b} & if \quad b > 0, \\ \int_{\mathcal{X}_{\overline{r}(x) = 0}} p(x) dx & if \quad b = 0, \end{cases} \tag{8}$$

$$\int_{\mathcal{X}_{\overline{r}(x) > b}} p(x) c^*(x) dx = 0, \tag{9}$$

*where $\overline{r}(x) = r(x) - \lambda$ measures how much the conditional risk $r(x)$ of the classifier $h(x)$ exceeds the target $\lambda$,*

$$\rho(\mathcal{X}') = \int_{\mathcal{X}'} p(x) \overline{r}(x) \, dx \tag{10}$$

*is the expectation of $\overline{r}(x)$ restricted to inputs in $\mathcal{X}'$, and*

$$b = \sup \{ a \mid \rho(\mathcal{X}_{\overline{r}(x) \le a}) \le 0 \} \ge 0. \tag{11}$$

**Proof** Observe that $b \ge 0$, because $\rho(\mathcal{X}_{\overline{r}(x) \le 0}) \le 0$. Next, observe that Problem 3 can be rewritten into the form

$$\max_{c \in [0,1]^{\mathcal{X}}} \int_{\mathcal{X}} p(x) c(x) dx \quad s.t. \quad \int_{\mathcal{X}} p(x) c(x) \overline{r}(x) dx \le 0 \tag{35}$$

since

$$R_S(h,c) - \lambda = \frac{\int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \, \ell(y, h(x)) \, c(x) \, dx - \lambda \phi(c)}{\phi(c)} \tag{36}$$

$$= \frac{\int_{\mathcal{X}} p(x) c(x) r(x) \, dx - \lambda \int_{\mathcal{X}} p(x) c(x)}{\phi(c)} = \frac{\int_{\mathcal{X}} p(x) c(x) \overline{r}(x) \, dx}{\phi(c)}. \tag{37}$$

Let $F(c) = \phi(c) = \int_{\mathcal{X}} p(x) c(x) dx$ denote the objective function of (35).

**Case 1** $b > 0$.
***Claim I*** *Each $c^* : \mathcal{X} \to [0,1]$ which fulfils (7), (8) and (9) is feasible to (35) and*

$$F(c^*) = \int_{\mathcal{X}_{\overline{r}(x) < b}} p(x) dx - \frac{1}{b} \rho(\mathcal{X}_{\overline{r}(x) < b}). \tag{38}$$

*Proof of Claim I.*

31

Equality (38) is simply obtained by summing LHS and RHS of (7), (8) and (9). To verify the constraint of (35), observe that, since $\overline{r}$ is a bounded function and

$$\int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)(c^*(x)-1)dx \stackrel{(7)}{=} 0\,, \tag{39}$$

it holds that

$$\int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)(c^*(x)-1)\overline{r}(x)dx = 0\,, \tag{40}$$

which implies

$$\int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)\overline{r}(x)dx = \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)\overline{r}(x)dx \stackrel{(10)}{=} \rho(\mathcal{X}_{\overline{r}(x)<b}). \tag{41}$$

If $b < \infty$, then

$$\int_{\mathcal{X}} p(x)c^*(x)\overline{r}(x)dx \stackrel{(9)}{=} \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)\overline{r}(x)dx + \int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)\overline{r}(x)dx$$

$$\stackrel{(41)}{=} \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)\overline{r}(x)dx + b\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx \tag{42}$$

$$\stackrel{(8),(10),(41)}{=} \rho(\mathcal{X}_{\overline{r}(x)<b}) - \rho(\mathcal{X}_{\overline{r}(x)<b}) = 0. \tag{43}$$

If $b = \infty$, then

$$\int_{\mathcal{X}} p(x)c^*(x)\overline{r}(x)dx = \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)\overline{r}(x)dx \stackrel{(41)}{=} \rho(\mathcal{X}_{\overline{r}(x)<b}) \leq 0.$$

**Claim II** Let $c : \mathcal{X} \rightarrow [0,1]$ be a feasible solution to (35) that violates at least one of the constraints (7), (8) and (9). Then, $F(c) < F(c^*)$, where $c^* : \mathcal{X} \rightarrow [0,1]$ is a confidence function satisfying (7), (8), (9), and, without loss of generality,

$$\forall x \in \mathcal{X}_{\overline{r}(x)<b} : c^*(x) = 1\,. \tag{44}$$

Proof of Claim II.

Distinguish three cases.

**Case 1.1** Condition (9) is violated (note that this is possible only if $b < \infty$), i.e.

$$\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx > 0. \tag{45}$$

Inequality (45) and Lemma 13 (applied to $f(x) = p(x)c(x)$ and $g(x) = \overline{r}(x)$) yield

$$\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)\overline{r}(x)dx > b\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx.$$

*Therefore, we can write*

$$\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)\overline{r}(x)dx = b' \int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx \tag{46}$$

*for a suitable $b' \in \mathbb{R}_+$ such that*

$$b' > b > 0. \tag{47}$$

*Based on the constraint of (35), derive*

$$\int_{\mathcal{X}} p(x)c(x)\overline{r}(x)dx \overset{(46)}{=} \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)\overline{r}(x)dx + b\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx + b'\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx$$

$$\overset{(35)}{\leq} 0 \overset{(43)}{=} \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)\overline{r}(x)dx + b\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx. \tag{48}$$

*Let $\sigma(x) = \frac{1}{b}\overline{r}(x)$. Inequality (48) can be rearranged and upper bounded as*

$$\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx - \int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx + \frac{b'}{b}\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx \tag{49}$$

$$\overset{(48)}{\leq} \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)(c^*(x)-c(x))\sigma(x)dx \leq \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)dx - \int_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)dx$$

*where the second inequality follows from $\forall x \in \mathcal{X}_{\overline{r}(x)<b} : \sigma(x) \leq 1$. From this we get*

$$\int_{\mathcal{X}_{\overline{r}(x)\leq b}} p(x)c(x)dx \overset{(49)}{\leq} \int_{\mathcal{X}_{\overline{r}(x)\leq b}} p(x)c^*(x)dx - \frac{b'}{b}\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx. \tag{50}$$

*Now, derive*

$$F(c) = \int_{\mathcal{X}_{\overline{r}(x)\leq b}} p(x)c(x)dx + \int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx$$

$$\overset{(50)}{\leq} \int_{\mathcal{X}_{\overline{r}(x)\leq b}} p(x)c^*(x)dx - \left(\frac{b'}{b}-1\right)\int_{\mathcal{X}_{\overline{r}(x)>b}} p(x)c(x)dx$$

$$\overset{(45),(47)}{<} \int_{\mathcal{X}_{\overline{r}(x)\leq b}} p(x)c^*(x)dx = F(c^*).$$

**Case 1.2** *Condition (9) holds, condition (8) is violated.*

*If $\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx < -\frac{\rho(\mathcal{X}_{\overline{r}(x)<b})}{b}$, then obviously $F(c) < F(c^*)$. Hence, assume*

$$\int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx > -\frac{\rho(\mathcal{X}_{\overline{r}(x)<b})}{b}. \tag{51}$$

33

*Analogically to (48), derive*

$$\int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)\overline{r}(x)dx + b \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx \overset{(35)}{\leq} 0 \tag{52}$$

$$\overset{(43)}{=} \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)\overline{r}(x)dx + b \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx, \tag{53}$$

*and*

$$\int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)\sigma(x)dx + \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx \leq \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)\sigma(x)dx + \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx \tag{54}$$

*where $\sigma(x) = \frac{1}{b}\overline{r}(x) < 1$ for all $x \in \mathcal{X}_{\overline{r}(x)<b}$.*

Denote and derive

$$\Delta = \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx - \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx \overset{(8)}{=} \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx + \frac{\rho(\mathcal{X}_{\overline{r}(x)<b})}{b} \overset{(51)}{>} 0. \tag{55}$$

*Then, (54) can be rewritten as*

$$\int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)(c^*(x) - c(x))\sigma(x)dx \geq \Delta \overset{(55)}{>} 0. \tag{56}$$

*Inequality (56) and Lemma 14 (applied to $g(x) = \sigma(x) < 1$ and $f(x) = p(x)(c^*(x)-c(x)) \overset{(44)}{\geq} 0$ over $\mathcal{X}_{\overline{r}(x)<b}$) yield*

$$\int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)(c^*(x) - c(x))dx > \Delta. \tag{57}$$

*Now, combine and rearrange (55) and (57) to obtain*

$$F(c^*) = \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)dx + \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c^*(x)dx \overset{(57)}{>} \Delta \tag{58}$$

$$\overset{(55)}{=} \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)dx + \int\limits_{\mathcal{X}_{\overline{r}(x)=b}} p(x)c(x)dx = F(c). \tag{59}$$

**Case 1.3** *Conditions (8) and (9) hold, condition (7) is violated, i.e.*

$$\int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)dx < \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)dx. \tag{60}$$

*Then,*

$$F(c^*) = \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c^*(x)dx - \frac{\rho(\mathcal{X}_{\overline{r}(x)<b})}{b} \overset{(60)}{>} \int\limits_{\mathcal{X}_{\overline{r}(x)<b}} p(x)c(x)dx - \frac{\rho(\mathcal{X}_{\overline{r}(x)<b})}{b} = F(c).$$

**Case 2** $b = 0$.

This occurs only if $\int_{\mathcal{X}_{\overline{r}(x)<0}} p(x)\overline{r}(x)dx = 0$. The constraint of (35) implies

$$\int_{\mathcal{X}_{\overline{r}(x)>0}} p(x)c(x)\overline{r}(x)dx = 0\,,$$

thus

$$\int_{\mathcal{X}_{\overline{r}(x)>0}} p(x)c(x)dx = 0\,,$$

which confirms condition (9).

Finally, the obvious equations

$$\max_{c:\mathcal{X}\rightarrow[0,1]} \int_{\mathcal{X}_{\overline{r}(x)<0}} p(x)c(x)dx = \int_{\mathcal{X}_{\overline{r}(x)<0}} p(x)dx\,, \text{ and}$$

$$\max_{c:\mathcal{X}\rightarrow[0,1]} \int_{\mathcal{X}_{\overline{r}(x)=0}} p(x)c(x)dx = \int_{\mathcal{X}_{\overline{r}(x)=0}} p(x)dx$$

confirm condition (7) and (8), respectively.

■

## A.3 Proof of Theorem 3

**Theorem 3** *Let $r\colon \mathcal{X} \rightarrow \mathbb{R}$ be the conditional risk (6) of a classifier $h\colon \mathcal{X} \rightarrow \mathcal{Y}$, $\gamma = b + \lambda$ the rejection threshold given by the target risk $\lambda$ and a constant $b$ computed by (11). Then the selection function*

$$c^*(x) = \begin{cases} 1 & \text{if } r(x) < \gamma\,, \\ \tau & \text{if } r(x) = \gamma\,, \\ 0 & \text{if } r(x) < \gamma\,, \end{cases} \tag{12}$$

*where $\tau$ is the acceptance probability given by*

$$\tau = \begin{cases} 1 & \text{if } \rho(\mathcal{X}_{r(x)=\gamma}) = 0\,, \\ -\frac{\rho(\mathcal{X}_{r(x)<\gamma})}{\rho(\mathcal{X}_{r(x)=\gamma})} & \text{if } \rho(\mathcal{X}_{r(x)=\gamma}) > 0\,, \end{cases} \tag{13}$$

*satisfies the optimality condition of Theorem 2, and hence it is a solution to Problem 3.*

**Proof** The optimality conditions (7) and (9) given in Theorem 2 are equivalent to a probabilistic statement $\mathbb{P}_{x\sim p(x)}[c^*(x) = 0 \wedge \overline{r}(x) < b] = 0$ and $\mathbb{P}_{x\sim p(x)}[c^*(x) = 1 \wedge \overline{r}(x) > b] = 0$, respectively. Hence the two conditions are satisfied by a selection function which predicts, $c^*(x) = 1$, whenever $\overline{r}(x) < b$ and rejects, $c^*(x) = 0$, whenever $\overline{r}(x) > b$. Or equivalently, using the identity $\overline{r}(x) = r(x) - \lambda$ and a threshold $\gamma = b + \lambda$, by $c^*(x) = 1$ when $r(x) < \gamma$ and $c^*(x) = 0$ when $r(x) > \gamma$. Finally, if we opt for a selection function that is constant

$c^*(x) = \tau$ inside the boundary region $\mathcal{X}_{\overline{r}(x)=b}$, then the condition (8) implies $\tau = -\frac{\rho(\mathcal{X}_{r(x)<\gamma})}{b \cdot \rho_0}$ if $b > 0$, where $\rho_0 = \int_{\mathcal{X}_{\overline{r}(x)=b}} p(x)\,dx$, and $\tau = 1$ if $b = 0$. Using $\mathcal{X}_{\overline{r}(x)<b} = \mathcal{X}_{r(x)<\gamma}$ and $b \cdot \rho_0 = \rho(\mathcal{X}_{\overline{r}(x)=b}) = \rho(\mathcal{X}_{r(x)=\gamma})$, we derive (13). ∎

## A.4 Proof of Theorem 4

**Theorem 4** *Let $(h, c)$ be an optimal solution to (14). Then, $(h_B, c)$, where $h_B$ is the optimal Bayes classifier (2), is also optimal to (14).*

**Proof** The theorem follows from the fact that $R_S(h_B, c) \leq R_S(h, c)$ for any $(h, c)$ feasible to (14), which is derived as follows:

$$
\begin{aligned}
R_S(h_B, c) &= \frac{1}{\phi(c)} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y)\, \ell(y, h_B(x))\, c(x)\, dx \\
&= \frac{1}{\phi(c)} \int_{\mathcal{X}} p(x) c(x) \left( \sum_{y \in \mathcal{Y}} p(y \mid x)\, \ell(y, h_B(x)) \right) dx \\
&\overset{(2)}{\leq} \frac{1}{\phi(c)} \int_{\mathcal{X}} p(x) c(x) \left( \sum_{y \in \mathcal{Y}} p(y \mid x)\, \ell(y, h(x)) \right) dx \\
&= \frac{1}{\phi(c)} \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y)\, \ell(y, h(x))\, c(x)\, dx \\
&= R_S(h, c).
\end{aligned}
$$

∎

## A.5 Proof of Theorem 5

**Problem 5 (Bounded-coverage model for known $h(x)$)** *Given a classifier $h(x)$ and a target coverage $0 < \omega \leq 1$, the optimal selection function $c^*(x)$ is a solution to the problem*

$$
\min_{c \in [0,1]^{\mathcal{X}}} R_S(h, c) \quad s.t. \quad \phi(c) \geq \omega, \tag{15}
$$

*where we assume that the minimizer exists.*

**Theorem 5** *A selection function $c^* : \mathcal{X} \to [0, 1]$ is an optimal solution to Problem 5 if and only if it holds*

$$
\int_{\mathcal{X}_{r(x)<\beta}} p(x) c^*(x)\, dx = \int_{\mathcal{X}_{r(x)<\beta}} p(x)\, dx, \tag{16}
$$

$$
\int_{\mathcal{X}_{r(x)=\beta}} p(x) c^*(x)\, dx = \omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)\, dx, \tag{17}
$$

$$
\int_{\mathcal{X}_{r(x)>\beta}} p(x) c^*(x)\, dx = 0, \tag{18}
$$

*where*

$$\beta = \inf \left\{ a \mid \int_{\mathcal{X}_{r(x) < a}} p(x) dx \geq \omega \right\}. \tag{19}$$

**Proof** By substituting the definitions of $R_S(h, c)$ and $\phi(c)$ to (15), we rewrite the problem into the form

$$\min_{c \in [0,1]^{\mathcal{X}}} \frac{\int_{\mathcal{X}} p(x) c(x) r(x) dx}{\int_{\mathcal{X}} p(x) c(x) dx} \quad \text{s.t.} \quad \int_{\mathcal{X}} p(x) c(x) dx \geq \omega. \tag{61}$$

Let $F(c) = \frac{\int_{\mathcal{X}} p(x) c(x) r(x) dx}{\int_{\mathcal{X}} p(x) c(x) dx}$ denote the objective function of (61). Whenever $c^* : \mathcal{X} \to [0, 1]$ fulfils (16), (17) and (18), it is feasible to (61) and

$$F(c^*) = \beta + \frac{1}{\omega} \int_{\mathcal{X}_{r(x) < \beta}} p(x) r(x) dx - \frac{\beta}{\omega} \int_{\mathcal{X}_{r(x) = \beta}} p(x) dx, \tag{62}$$

which is a value independent of $c^*$.

We will prove the theorem by showing that any $c : \mathcal{X} \to [0, 1]$ feasible to (61) that violates at least one of conditions (16), (17), (18) is not an optimal solution. Three cases will be examined.

**Case 1** Condition (16) is violated, i.e.,

$$\int_{\mathcal{X}_{r(x) < \beta}} p(x) c(x) dx < \int_{\mathcal{X}_{r(x) < \beta}} p(x) dx. \tag{63}$$

This means that there is a subset $X \subseteq \mathcal{X}$ such that

$$\forall x \in X \ : \ r(x) \geq \beta \tag{64}$$

and

$$\int_{X} p(x) c(x) dx = \int_{\mathcal{X}_{r(x) < \beta}} p(x) dx - \int_{\mathcal{X}_{r(x) < \beta}} p(x) c(x) dx \overset{(63)}{>} 0. \tag{65}$$

Define $c' : \mathcal{X} \to [0, 1]$ as follows.

$$c'(x) = \begin{cases} 1 & \text{if} & r(x) < \beta, \\ 0 & \text{if} & x \in X, \\ c(x) & \text{otherwise}. \end{cases} \tag{66}$$

$c'$ is feasible to (61) as $\phi(c') = \phi(c)$. Derive

$$\phi(c)\left(F(c) - F(c')\right) = \int_X p(x)c(x)r(x)dx - \int_{\mathcal{X}_{r(x)<\beta}} p(x)r(x)dx + \int_{\mathcal{X}_{r(x)<\beta}} p(x)c(x)r(x)dx$$

$$(67)$$

$$\overset{(64),(65)}{\geq} \int_{\mathcal{X}_{r(x)<\beta}} \beta \cdot p(x)(1 - c(x))dx - \int_{\mathcal{X}_{r(x)<\beta}} p(x)(1 - c(x))r(x)dx$$

$$(68)$$

$$= \int_{\mathcal{X}_{r(x)<\beta}} p(x)(1 - c(x))(\beta - r(x))dx > 0 \tag{69}$$

where the inequality in (69) is obtained from Lemma 12 applied to $f(x) = p(x)(1 - c(x))$, $g(x) = \beta - r(x)$, and the set $\mathcal{X}_{r(x)<\beta}$. This shows that $c$ is not an optimal solution.

**Case 2** Condition (16) is satisfied, condition (17) is violated and

$$\int_{\mathcal{X}_{r(x)=\beta}} p(x)c(x)dx < \omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx > 0. \tag{70}$$

In this case, there is a $c' : \mathcal{X} \to [0, 1]$ such that

$$c'(x) = c(x) \quad \text{if} \ \ r(x) < \beta, \tag{71}$$

$$c'(x) = 0 \qquad \text{if} \ \ r(x) > \beta, \tag{72}$$

and

$$\int_{\mathcal{X}_{r(x)=\beta}} p(x)c'(x)dx = \int_{\mathcal{X}_{r(x)=\beta}} p(x)c(x)dx + \int_{\mathcal{X}_{r(x)>\beta}} p(x)c(x)dx. \tag{73}$$

If Lemma 13 is applied to $f(x) = p(x)c(x)$, $g(x) = r(x)$, and the set $\mathcal{X}_{r(x)>\beta}$, we get

$$\int_{\mathcal{X}_{r(x)>\beta}} p(x)c(x)r(x)dx > \beta \int_{\mathcal{X}_{r(x)>\beta}} p(x)c(x)dx. \tag{74}$$

It also holds that $\phi(c') = \phi(c)$. Now, deriving

$$\phi(c)\left(F(c) - F(c')\right) \overset{(73)}{=} \int_{\mathcal{X}_{r(x)>\beta}} p(x)c(x)r(x)dx - \beta \int_{\mathcal{X}_{r(x)>\beta}} p(x)c(x)dx \tag{75}$$

$$\overset{(74)}{>} \beta \int_{\mathcal{X}_{r(x)>\beta}} p(x)c(x) - \beta \int_{\mathcal{X}_{r(x)>\beta}} p(x)x(x)dx = 0 \tag{76}$$

shows that $c$ is not an optimal solution.

**Case 3** $\phi(c) > \omega$, which occurs if

$$\int_{\mathcal{X}_{r(x)=\beta}} p(x)c(x)dx > \omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx > 0 \tag{77}$$

38

(implying that condition (17) is violated), or if condition (18) is violated.

Observe that $F(c) = F(\alpha \cdot c)$ for any $a \in \mathbb{R}_+$. Let $c' = \frac{\omega}{\phi(c)} \cdot c$. Since $\phi(c') = \omega$, the selection function $c'$ is feasible to (61). Because

$$\int_{\mathcal{X}_{r(x)<\beta}} p(x)c'(x)dx = \frac{\omega}{\phi(c)} \int_{\mathcal{X}_{r(x)<\beta}} p(x)c(x)dx < \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx , \qquad (78)$$

$c'$ violates condition (16) and is therefore not an optimal solution (see Case 1). This implies that $c$ is not an optimal solution too. ∎

## A.6 Proof of Theorem 6

**Theorem 6** *Let $r \colon \mathcal{X} \to \mathbb{R}$ be the conditional risk (6) of a classifier $h \colon \mathcal{X} \to \mathcal{Y}$, $1 \geq \omega > 0$ be a target coverage and $\beta$ be the constant computed by (19). Then the selection function*

$$c^*(x) = \begin{cases} 1 & if \quad r(x) < \beta \,, \\ \kappa & if \quad r(x) = \beta \,, \\ 0 & if \quad r(x) > \beta \,, \end{cases} \qquad (20)$$

*where $\kappa$ is the acceptance probability given by*

$$\kappa = \begin{cases} 0 & if \ \int_{\mathcal{X}_{r(x)=\beta}} p(x)dx = 0 \,, \\ \frac{\omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx}{\int_{\mathcal{X}_{r(x)=\beta}} p(x)dx} & otherwise \,, \end{cases} \qquad (21)$$

*satisfies the optimality condition of Theorem 6, and hence it is a solution of Problem 5.*

**Proof** It is easy to see that $c^*$ satisfies conditions (16) and (18). The validity of condition (17) is proved as follows. If $\int_{\mathcal{X}_{r(x)=\beta}} p(x)dx = 0$, then $\int_{\mathcal{X}_{r(x)=\beta}} p(x)dx = \omega$, and condition (17) is met. If $\int_{\mathcal{X}_{r(x)=\beta}} p(x)dx > 0$, we derive

$$\int_{\mathcal{X}_{r(x)=\beta}} p(x)c^*(x)dx = \frac{\omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx}{\int_{\mathcal{X}_{r(x)=\beta}} p(x)dx} \int_{\mathcal{X}_{r(x)=\beta}} p(x)dx = \omega - \int_{\mathcal{X}_{r(x)<\beta}} p(x)dx . \quad (79)$$

∎

## Appendix B. Proofs of theorems from Section 3

### B.1 Proof of Theorem 7

The expectation of the squared loss deviation reads

$$E_{\text{reg}}(s) = \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \Big( \ell(y, h(x)) - s(x) \Big)^2 dx . \qquad (80)$$

**Theorem 7** *The conditional risk $r(x)$ defined by (6) is an optimal solution to $\min\limits_{s:\mathcal{X}\to\mathbb{R}} E_{\mathrm{reg}}(s)$.*

**Proof** We can rewrite $E_{\mathrm{reg}}(s)$ as

$$E_{\mathrm{reg}}(s) = \int_{\mathcal{X}} p(x) \sum_{y\in\mathcal{Y}} p(y\mid x)\Big(\ell(y,h(x))^2 - 2\,\ell(y,h(x))\,s(x) + s(x)^2\Big)dx = \int_{\mathcal{X}} p(x)f(s(x))dx.$$

Due to additivity we can solve $\min_{s:\,\mathcal{X}\to\mathbb{R}} E_{\mathrm{reg}}(s)$ for each $x\in\mathcal{X}$ separately by setting derivative of $f(s)$ to zero and solving for $s$ which yields

$$f'(s) = -2\sum_{y\in\mathcal{Y}} p(y\mid x)\ell(y,h(x)) + 2\sum_{y\in\mathcal{Y}} p(y\mid x)s(x) = 0 \Rightarrow s^*(x) = \sum_{y\in\mathcal{Y}} p(y\mid x)\ell(y,h(x))\,.$$

∎

## B.2 Proof of Theorem 8

**Lemma 15** *For an integer $n \geq 1$, let $\{a_i\}_{i=1}^n$ and $\{b_i\}_{i=1}^n$ be sequences of positive real numbers such that $\{\frac{a_i}{b_i}\}_{i=1}^n$ is a non-increasing sequence. For each non-decreasing sequence $\{\ell_i\}_{i=1}^n$ of non-negative real numbers with a positive sum it holds that*

$$\frac{\sum_{i=1}^n a_i\ell_i}{\sum_{i=1}^n b_i\ell_i} \leq \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i}\,.$$

**Proof** By induction on $n$. Base case: If $n = 1$, then $\frac{a_1\ell_1}{b_1\ell_1} = \frac{a_1}{b_1}$.

Induction step: Let $n > 1$. The fact that $\{\frac{a_i}{b_i}\}_{i=1}^n$ is non-increasing implies $a_i \leq \frac{a_1}{b_1}\cdot b_i$ for all $i = 1,\ldots,n$, hence

$$\frac{\sum_{i=2}^n a_i}{\sum_{i=2}^n b_i} \leq \frac{\sum_{i=2}^n \frac{a_1}{b_1}\cdot b_i}{\sum_{i=2}^n b_i} = \frac{a_1}{b_1}\,. \tag{81}$$

The lemma is obviously satisfied if $0 < \ell_1 = \ell_2 = \ldots = \ell_n$. Assume that $\ell_n > \ell_1$. Then, the sequence $\{\ell_i - \ell_n\}_{i=2}^n$ is non-decreasing with positive sum, hence the induction hypothesis yields that

$$\frac{\sum_{i=2}^n a_i(\ell_i - \ell_1)}{\sum_{i=2}^n b_i(\ell_i - \ell_1)} \leq \frac{\sum_{i=2}^n a_i}{\sum_{i=2}^n b_i} \leq \frac{a_1}{b_1}\,. \tag{82}$$

The induction hypotheses also ensures that

$$\frac{\sum_{i=2}^n a_i\ell_i}{\sum_{i=2}^n b_i\ell_i} \leq \frac{\sum_{i=2}^n a_i}{\sum_{i=2}^n b_i}\,. \tag{83}$$

40

We can thus derive the following sequence of equivalent inequalities:

$$b_1 \sum_{i=2}^{n} a_i (\ell_i - \ell_1) \overset{(82)}{\le} a_1 \sum_{i=2}^{n} b_i (\ell_i - \ell_1)$$

$$b_1 \sum_{i=2}^{n} a_i \ell_i + a_1 \ell_1 \sum_{i=2}^{n} b_i \le a_1 \sum_{i=2}^{n} b_i \ell_i + b_1 \ell_1 \sum_{i=2}^{n} a_i$$

$$a_1 \ell_1 b_1 + b_1 \sum_{i=2}^{n} a_i \ell_i + a_1 \ell_1 \sum_{i=2}^{n} b_i + \sum_{i=2}^{n} b_i \sum_{i=2}^{n} a_i \ell_i \le a_1 \ell_1 b_1 + a_1 \sum_{i=2}^{n} b_i \ell_i + b_1 \ell_1 \sum_{i=2}^{n} a_i + \sum_{i=2}^{n} b_i \sum_{i=2}^{n} a_i \ell_i$$

$$a_1 \ell_1 b_1 + b_1 \sum_{i=2}^{n} a_i \ell_i + a_1 \ell_1 \sum_{i=2}^{n} b_i + \sum_{i=2}^{n} b_i \sum_{i=2}^{n} a_i \ell_i \overset{(83)}{\le} a_1 \ell_1 b_1 + a_1 \sum_{i=2}^{n} b_i \ell_i + b_1 \ell_1 \sum_{i=2}^{n} a_i + \sum_{i=2}^{n} a_i \sum_{i=2}^{n} b_i \ell_i$$

$$a_1 \ell_1 \sum_{i=1}^{n} b_i + \sum_{i=1}^{n} b_i \sum_{i=2}^{n} a_i \ell_i \le b_1 \ell_1 \sum_{i=1}^{n} a_i + \sum_{i=1}^{n} a_i \sum_{i=2}^{n} b_i \ell_i$$

$$\sum_{i=1}^{n} a_i \ell_i \sum_{i=1}^{n} b_i \le \sum_{i=1}^{n} b_i \ell_i \sum_{i=1}^{n} a_i$$

$$\frac{\sum_{i=1}^{n} a_i \ell_i}{\sum_{i=1}^{n} b_i \ell_i} \le \frac{\sum_{i=1}^{n} a_i}{\sum_{i=1}^{n} b_i}.$$

∎

For $i = 1, \ldots, n$ and a permutation $\pi$ on $\{1, \ldots, n\}$, let $a_{\pi_i} = \sum_{j=i}^{n} \frac{n}{j}$ and $b_{\pi_i} = n - i + 1$. We can write

$$\Delta_{\text{sele}}(s, \mathcal{T}_n) = \frac{1}{n^2} \sum_{i=1}^{n} b_{\pi_i} \ell_{\pi_i}$$

and

$$\text{AuRC}(s, \mathcal{T}_n) = \frac{1}{n^2} \sum_{i=1}^{n} a_{\pi_i} \ell_{\pi_i},$$

where $\ell_{\pi_i} = \ell(y_{\pi_i}, h(x_{\pi_i}))$.

Without loss of generality, assume that $\ell_1 \le \ell_2 \le \ldots \le \ell_n$ and $\pi_i = i$ for all $i = 1, \ldots, n$. Let $H_k = \sum_{i=1}^{k} \frac{1}{i}$ denote the $k$-th harmonic number. It fulfils

$$\ln(k) + \gamma + \frac{1}{2k + 1} \le H_k \le \ln(k) + \gamma + \frac{1}{2k - 1} \tag{84}$$

where $\gamma \approx 0.5772156649$ is the Euler–Mascheroni constant. Moreover, let us define $H_0 = 0$.

**Lemma 16** *Let $n \ge 2$ be an integer. For each $i = 1, \ldots, n$, let $a_i = \sum_{j=i}^{n} \frac{n}{j}$ and $b_i = n - i + 1$. Then, $\frac{a_i}{b_i} \ge \frac{a_{i+1}}{b_{i+1}}$ holds for all $i = 1, \ldots, n - 1$.*

**Proof** Since

$$\frac{a_i}{b_i} - \frac{a_{i+1}}{b_{i+1}} = \frac{n \cdot (H_n - H_{i-1})}{n - i + 1} - \frac{n \cdot (H_n - H_i)}{n - i}$$

it suffices to show that

$$(n - i)(H_n - H_{i-1}) \geq (n - i + 1)(H_n - H_i)$$

and this is equivalent to showing that

$$(n - i + 1)H_i - (n - i)H_{i-1} \geq H_n \,.$$

If we substitute $H_{i-1} = H_i - \frac{1}{i}$, the inequality further reduces to

$$H_n - H_i \leq \frac{n - i}{i} = \frac{n}{i} - 1 \,. \tag{85}$$

Now we derive

$$H_n - H_i \overset{(84)}{\leq} \ln\left(\frac{n}{i}\right) + \frac{1}{2n - 1} - \frac{1}{2i + 1} < \frac{n}{i} - 1$$

where the second inequality follows from the fact that $\ln(x) < x - 1$ for all $x > 1$. This confirms inequality (85).

■

**Lemma 17** $\sum_{i=1}^{n}(H_n - H_{i-1}) = n$.

**Proof** By induction on $n$. The lemma trivially holds for $n = 1$. Let $n > 1$. Then, using the induction hypothesis, we derive

$$\sum_{i=1}^{n}(H_n - H_{i-1}) = H_n - H_{n-1} + \sum_{i=1}^{n-1}\left(\frac{1}{n} + H_{n-1} - H_{i-1}\right)$$

$$= \frac{1}{n} + \frac{n - 1}{n} + \sum_{i=1}^{n-1}(H_{n-1} - H_{i-1}) = n \,.$$

■

**Theorem 8** *The inequality* $\mathrm{AuRC}(s, \mathcal{T}_n) \leq 2 \cdot \Delta_{\mathrm{sele}}(s, \mathcal{T}_n)$ *holds true for any* $s \colon \mathcal{X} \to \mathbb{R}$ *and* $\mathcal{T}_n = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \mid i = 1, \ldots, n\}$.

**Proof** The theorem trivially holds if $\sum_{i=1}^{n}\ell_i = 0$. If $\sum_{i=1}^{n}\ell_i > 0$, we apply Lemmas 16, 15 and 17 to derive

$$\frac{\mathrm{AuRC}(s, \mathcal{T}_n)}{\Delta_{\mathrm{sele}}(s, \mathcal{T}_n)} \leq \frac{\sum_{i=1}^{n}a_i}{\sum_{i=1}^{n}b_i} = \frac{\sum_{i=1}^{n}n \cdot (H_n - H_{i-1})}{\sum_{i=1}^{n}(n - i + 1)} = \frac{n^2}{\frac{n}{2}(n + 1)} = \frac{2n}{n + 1} < 2 \,.$$

■

## B.3 Proof of Theorem 9

**Remark 18** *For the sake of simplicity, for predicates $\varphi_1(x, z), \ldots, \varphi_k(x, z)$ and a function $f : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, we write*

$$\int\limits_{\mathcal{X}} \int\limits_{\substack{\varphi_1(x,z) \\ \vdots \\ \varphi_k(x,z)}} f(x, z) dz\, dx$$

*to represent*

$$\int\limits_{\mathcal{X}} \int\limits_{\mathcal{X}} f(x, z) [\![ \varphi_1(x, z) \wedge \ldots \wedge \varphi_k(x, z) ]\!] dz\, dx\,.$$

**Theorem 9** *A function $s^* : \mathcal{X} \to \mathbb{R}$ is an optimal solution to $\min_{s:\mathcal{X}\to\mathbb{R}} E_{\mathrm{sele}}(s)$ iff*

$$\int_{\mathcal{X}} \int\limits_{\substack{z\neq x \\ s^*(z)=s^*(x)}} \max\{r(x), r(z)\} p(x) p(z) dz\, dx = 0\,, \text{ and} \tag{29}$$

$$\int_{\mathcal{X}} \int\limits_{\substack{r(z)<r(x) \\ s^*(z)>s^*(x)}} (r(x) - r(z))\, p(x) p(z) dz\, dx = 0\,. \tag{30}$$

**Proof** We first present four equalities to be used later. We assume that $s : \mathcal{X} \to \mathbb{R}$ is any measurable function. The validity of the equalities can be easily verified.

$$\int\limits_{\mathcal{X}} r(x)p(x) \int\limits_{\substack{r(z)>r(x) \\ s(z)<s(x)}} p(z) dz\, dx = \int\limits_{\mathcal{X}} p(z) \int\limits_{\substack{r(x)<r(z) \\ s(x)>s(z)}} r(x)p(x) dx\, dz = \int\limits_{\mathcal{X}} p(x) \int\limits_{\substack{r(z)<r(x) \\ s(z)>s(x)}} r(z)p(z) dz\, dx\,,$$
$$\tag{86}$$

$$\int\limits_{\mathcal{X}} \int\limits_{\substack{r(z)<r(x) \\ s(z)=s(x)}} r(x)p(x)p(z) dz\, dx = \frac{1}{2} \int\limits_{\mathcal{X}} \int\limits_{\substack{z\neq x \\ s(z)=s(x)}} \max\{r(x), r(z)\} p(x) p(z) dz\, dx$$
$$- \frac{1}{2} \int\limits_{\mathcal{X}} \int\limits_{\substack{z\neq x \\ r(z)=r(x) \\ s(z)=s(x)}} \max\{r(x), r(z)\} p(x) p(z) dz\, dx\,, \tag{87}$$

$$\int\limits_{\mathcal{X}} \int\limits_{\substack{r(z)=r(x) \\ s(z)<s(x)}} r(x)p(x)p(z) dz\, dx = \frac{1}{2} \int\limits_{\mathcal{X}} \int\limits_{\substack{r(z)=r(x)}} r(x)p(x)p(z) dz\, dx - \frac{1}{2} \int\limits_{\mathcal{X}} \int\limits_{\substack{r(z)=r(x) \\ s(z)=s(x)}} r(x)p(x)p(z) dz\, dx\,,$$
$$\tag{88}$$

$$\int\limits_{\mathcal{X}} \int\limits_{\substack{r(z)=r(x) \\ s(z)=s(x)}} r(x)p(x)p(z) dz\, dx = 2 \int\limits_{\mathcal{X}} \int\limits_{\substack{z>x \\ r(z)=r(x) \\ s(z)=s(x)}} r(x)p(x)p(z) dz\, dx + \int\limits_{\mathcal{X}} \int\limits_{\substack{z=x}} r(x)p(x)p(z) dz\, dx\,.$$
$$\tag{89}$$

Since $\operatorname{argmin}_{s:\mathcal{X}\to\mathbb{R}} E(s) = \operatorname{argmin}_{s:\mathcal{X}\to\mathbb{R}} (E(s) - E(r))$, it suffices to analyze minimizers of $E(s) - E(r)$ instead of $E(s)$. Derive

$$E(s) - E(r) = \int_{\mathcal{X}} \int_{s(z)\geq s(x)} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{r(z)\geq r(x)} r(x)p(x)p(z)dz\,dx$$

$$= \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)\geq s(x)}} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{\substack{r(z)\geq r(x)\\s(z)<s(x)}} r(x)p(x)p(z)dz\,dx$$

$$= \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)>s(x)}} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{\substack{r(z)>r(x)\\s(z)<s(x)}} r(x)p(x)p(z)dz\,dx$$

$$+ \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)=s(x)}} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{\substack{r(z)=r(x)\\s(z)<s(x)}} r(x)p(x)p(z)dz\,dx$$

$$= F_1(s) + F_2(s)$$

where

$$F_1(s) = \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)>s(x)}} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{\substack{r(z)>r(x)\\s(z)<s(x)}} r(x)p(x)p(z)dz\,dx$$

$$\overset{(86)}{=} \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)>s(x)}} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)>s(x)}} r(z)p(x)p(z)dz\,dx$$

$$= \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)>s(x)}} (r(x) - r(z))\, p(x)p(z)dz\,dx$$

and

$$F_2(s) = \int_{\mathcal{X}} \int_{\substack{r(z)<r(x)\\s(z)=s(x)}} r(x)p(x)p(z)dz\,dx - \int_{\mathcal{X}} \int_{\substack{r(z)=r(x)\\s(z)<s(x)}} r(x)p(x)p(z)dz\,dx$$

$$\overset{(87,88,89)}{=} \frac{1}{2}\int_{\mathcal{X}} \int_{\substack{z\neq x\\s(z)=s(x)}} \max\{r(x),r(z)\}p(x)p(z)dz\,dx + \frac{1}{2}\int_{\mathcal{X}} \int_{z=x} r(x)p(x)p(z)dz\,dx$$

$$- \frac{1}{2}\int_{\mathcal{X}} \int_{r(z)=r(x)} r(x)p(x)p(z)dz\,dx\,.$$

Observe that

$$\min_{s:\mathcal{X}\to\mathbb{R}} F_1(s) = 0,$$

44

$$\min_{s:\mathcal{X}\to\mathbb{R}} F_2(s) = \frac{1}{2}\int_\mathcal{X}\int_{z=x} r(x)p(x)p(z)dz\,dx - \frac{1}{2}\int_\mathcal{X}\int_{r(z)=r(x)} r(x)p(x)p(z)dz\,dx\,,$$

and both minima are attained by a scoring function $s^*: \mathcal{X} \to \mathbb{R}$ if and only if conditions (29) and (30) hold for $s^*$. Also note that the conditions can be fulfilled, e.g. by any $s^*$ such that

$$(\forall x, z \in \mathcal{X})\,(x \neq z \Rightarrow s^*(x) \neq s^*(z) \ \wedge \ r(x) < r(z) \Rightarrow s^*(x) < s^*(z))\,.$$

∎

## B.4 Proof of Theorem 11

The expectation of $\psi_\text{sele}$ reads

$$E_\text{proxy}(s) = \frac{n^2 - n}{n^2}\int_\mathcal{X} p(x)r(x)\left(\int_\mathcal{X} p(z)\log\left(1 + \exp(s(z) - s(x))\right)dz\right)dx + \frac{\log(2)}{n}\int_\mathcal{X} p(x)r(x)dx.$$

**Theorem 11** *Let $s^*: \mathcal{X} \to \mathbb{R}$ be an optimal solution to $\min_{s\colon\mathcal{X}\to\mathbb{R}} E_\text{proxy}(s)$. Then, the condition*

$$\forall(x, x') \in \mathcal{X} \times \mathcal{X} : r(x) < r(x') \Rightarrow s^*(x) < s^*(x')$$

*is satisfied almost surely.*

**Proof** For every $a \in \mathcal{X}$, $E_\text{proxy}(s)$ can be seen as a function of one variable $s(a)$, where the others $s(b)$, $b \in \mathcal{X} \setminus \{a\}$ are fixed. Hence if $s$ is a minimizer of $E_\text{proxy}(s)$, then for every $s(a), a \in \mathcal{X}$, the partial derivative w.r.t. $s(a)$ must be zero, i.e.,

$$\begin{aligned}
0 &= \frac{\partial}{\partial s(a)}\int_\mathcal{X} p(x)\,r(x)\left(\int_\mathcal{X} p(z)\log\left(1 + \exp(s(z) - s(x))\right)dz\right)dx \\
&= p(a)r(a)\int_\mathcal{X} p(z)\frac{-\exp(s(z) - s(a))}{1 + \exp(s(z) - s(a))}dz + p(a)\int_\mathcal{X} p(x)r(x)\frac{\exp(s(a) - s(x))}{1 + \exp(s(a) - s(x))}dx \\
&= p(a)\int_\mathcal{X}\frac{r(a)\,p(x) + r(x)\,p(x)}{1 + \exp(s(a) - s(x))}dx - p(a)\int_\mathcal{X} p(x)\,r(x)\,dx \\
&= f(r(a), s(a)) - C\,.
\end{aligned}$$

It shows that $f(r(a), s(a)) = C$ for any $a \in \mathcal{X}$ in order to guarantee that $s$ is a minimizer of $E_\text{proxy}(s)$. We prove by contradiction that the condition $r(a) < r(b) \Rightarrow s(a) < s(b)$ is satisfied up to a set $\{(a, b)|(a, b) \in \mathcal{X}^2\}$ of zero measure. Assume $s$ is optimal and the condition is violated, i.e. $r(a) < r(b) \wedge s(a) \geq s(b)$ holds for a pair $(a, b) \in \mathcal{X}^2$. Since $s$ is optimal then $f(r(a), s(a)) = C$. Since $r(b) > r(a)$ and $s(b) \leq s(a)$ then $f(r(b), s(b)) > f(r(a), s(b))$ because the function $f(u, v)$ is strictly increasing in $u$ and strictly decreasing in $v$. Combined it implies that $f(r(b), s(b)) > 0$ which leads to a contradiction because an optimal $s$ requires $f(r(b), s(b)) = C$.

∎

## References

P. L. Bartlett and M. H. Wegkamp. Classification with a reject option using a hinge loss. *Journal of Machine Learning Research*, 9(59):1823–1840, 2008.

C.C. Chang and C.J.Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. URL http://www.csie.ntu.edu.tw/~cjlin/libsvm.

C. Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*, 16(1):41–46, 1970.

W. Chu and S. S. Keerthi. New approaches to support vector ordinal regression. In *Proceedings of the International Conference on Machine Learning*, pages 145–152, 2005.

C. Corbiere, N. Thome, A. Bar-Hen, M. Cord, and P. Perez. Addressing failure prediction by learning model confidence. In *Advances in Neural Information Processing Systems*, volume 32, pages 2902–2913, 2019.

C. Cortes, G. DeSalvo, and M. Mohri. Boosting with abstention. In *Advances in Neural Information Processing Systems*, volume 29, pages 1660–1668, 2016.

N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proceedings of Conference on Computer Vision and Patter Recognition*, volume 1, pages 886–893, 2005.

J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7(1):1–30, 2006.

D. Dua and E. Karra Taniskidou. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.

R. El-Yaniv and Y. Wiener. On the foundations of noise-free selective classification. *Journal of Machine Learning Research*, 11(53):1605–1641, 2010.

L. Fischer, B. Hammer, and H. Wersing. Optimal local rejection for classifiers. *Neurocomputing*, 214:445–457, 2016.

L. Fisher, B. Hammer, and H. Wersing. Efficient rejection strategies for prototype-based classification. *Neurocomputing*, 169:334 – 342, 2015.

V. Franc and D. Prusa. On discriminative learning of prediction uncertainty. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pages 1963–1971, 2019.

G. Fumera and F. Roli. Support vector machines with embedded reject option. In *Pattern Recognition with Support Vector Machines, Lecture Notes in Computer Science*, volume 2388. Springer, 2002.

G. Fumera, F. Roli, and G. Giacinto. Multiple reject thresholds for improving classification reliability. In *Advances in Pattern Recognition*, pages 863–871, 2000.

Y. Geifman and R. El-Yaniv. Selective classification for deep neural networks. In *Advances in Neural Information Processing Systems 30*, pages 4878–4887, 2017.

Y. Grandvalet, A. Rakotomamonjy, J. Keshet, and S. Canu. Support vector machines with a reject option. In *Advances in Neural Information Processing Systems*, volume 21, pages 537–544, 2008.

B. Hanczar and E. R. Dougherty. Classification with reject option in gene expression data. *Bioinformatics*, 24:1889–1895, 2008.

T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: data mining, inference and prediction.* Springer, 2009.

R. Herbei and M.H. Wegkamp. Classification with reject option. *The Canadian Journal of Statistics / La Revue Canadienne de Statistique*, 34(4):709–721, 2006.

H. Jiang, B. Kim, M. Y. Guan, and M. Gupta. To trust or not to trust a classifier. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, page 5546–5557, 2018.

V. Kazemi and J. Sullivan. One millisecond face alignment with an ensemble of regression trees. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1867–1874, 2014.

D. E. King. Dlib-ml: A machine learning toolkit. *Journal of Machine Learning Research*, 10:1755–1758, 2009.

J. Kummert, B. Paassen, J. Jensen, C. Göpfert, and B. Hammer. Local reject option for deterministic multi-class SVM. In *Artificial Neural Networks and Machine Learning – ICANN, Lecture Notes in Computer Science*, volume 9887. Springer, 2016.

B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, volume 30, pages 6402–6413, 2017.

Y. LeCun, B. Boser, J.S. Denker, D. Henderson, R.E. Howard, W. Hubbard, and L.D. Jakel. Handwritten digit recognition with a back-propagation networks. In *Advances in Neural Information Processing Systems*, volume 2, pages 396–404, 1990.

J. Lei. Classification with confidence. *Bimetrika*, 101:755–769, 2014.

T. Pietraszek. Optimizing abstaining classifiers using ROC analysis. In *Proceedings of the 22nd International Conference on Machine Learning*, page 665–672, 2005.

C. Sagonas, E. Antonakos, G. Tzimiropoulos, S. Zafeiriou, and M. Pantic. 300 faces in-the-wild challenge: database and results. *Image and Vision Computing*, 47:3 – 18, 2016.

M.I. Schlesinger and V. Hlaváč. *Ten lectures on statistical and structural pattern recognition.* Kluwer Academic Publishers, 2002.

E.M. Stein and R. Shakarchi. *Real Analysis: Measure Theory, Integration, and Hilbert Spaces*. Princeton University Press, 2009.

C. H. Teo, S.V.N. Vishwanthan, A. J. Smola, and Q. V. Le. Bundle methods for regularized risk minimization. *Journal of Machine Learning Research*, 11(10):311–365, 2010.

F. Tortorella. An optimal reject rule for binary classifiers. In *Advances in Pattern Recognition, Lecture Notes in Computer Science*, volume 1876. Springer, 2000.

V.N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, Inc., 1998.

T. Villman, M. Kaden, A. Bohnsack, J. M. Villman, T. Drogies, S. Saralajew, and B. Hammer. Self-adjusting reject options in prototype based classification. In *Advances in Intelligent Systems and Computing*, volume 428. Springer, 2016.

M. Yuan and M. Wegkamp. Classification methods with reject option based on convex risk minimization. *Journal of Machine Learning Research*, 11(5):111–130, 2010.

H. Zaragoza and F. d'Alche Buc. Confidence measures for neural network classifiers. In *7th Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems*, 1998.