

---

# Reject option models comprising out-of-distribution detection

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The optimal prediction strategy for out-of-distribution (OOD) setups is a funda-  
2 mental question in machine learning. In this paper, we address this question and  
3 present several contributions. We propose three reject option models for OOD  
4 setups: the Cost-based model, the Bounded TPR-FPR model, and the Bounded  
5 Precision-Recall model. These models extend the standard reject option models  
6 used in non-OOD setups and define the notion of an optimal OOD selective classi-  
7 fier. We establish that all the proposed models, despite their different formulations,  
8 share a common class of optimal strategies. Motivated by the optimal strategy, we  
9 introduce double-score OOD methods that leverage uncertainty scores from two  
10 chosen OOD detectors: one focused on OOD/ID discrimination and the other on  
11 misclassification detection. The experimental results consistently demonstrate the  
12 superior performance of this simple strategy compared to state-of-the-art methods.  
13 Additionally, we propose novel evaluation metrics derived from the definition of  
14 the optimal strategy under the proposed OOD rejection models. These new metrics  
15 provide a comprehensive and reliable assessment of OOD methods without the  
16 deficiencies observed in existing evaluation approaches.

## 17 1 Introduction

18 Most methods for learning predictors from data are based on the closed-world assumption, i.e., the  
19 training and the test samples are generated i.i.d. from the same distribution, so-called in-distribution  
20 (ID). However, in real-world applications, ID test samples can be contaminated by samples from  
21 another distribution, the so-called Out-of-Distribution (OOD), which is not represented in training  
22 examples. A trustworthy prediction model should detect OOD samples and reject to predict them,  
23 while simultaneously minimizing the prediction error on accepted ID samples.

24 In recent years, the development of deep learning models for handling OOD data has emerged as a  
25 critical challenge in the field of machine learning, leading to an explosion of research papers dedicated  
26 to developing effective OOD detection methods (OODD) [10, 11, 4, 3, 12, 8, 1, 17, 16, 19, 20].  
27 Existing methods use various principles to learn a classifier of ID samples and a selective function  
28 that accepts the input for prediction or rejects it to predict. We further denote the pair of ID classifier  
29 and the selective function as OOD selective classifier, borrowing terminology from the non-OOD  
30 setup [7]. There is an agreement that a good OOD selective classifier should reject OOD samples  
31 and simultaneously achieve high classification accuracy on ID samples that are accepted [22]. To  
32 our knowledge, there is surprisingly no formal definition of an optimal OOD selective classifier.  
33 Consequently, there is also no consensus on how to evaluate the OODD methods. The commonly used  
34 metrics [21] evaluate only one aspect of the OOD selective classifier, either the accuracy of the ID  
35 classifier or the performance of the selective function as an OOD/ID discriminator. Such evaluation  
36 is inconclusive and usually inconsistent; e.g., the two most commonly used metrics, AUROC and  
37 OSCR, often lead to a completely reversed ranking of evaluated methods (see Sec. 3.4).

38 In this paper, we ask the following question: What would be the optimal prediction strategy for  
39 the OOD setup in the ideal case when ID and OOD distributions were known? To this end, we  
40 offer the contributions: (i) We propose three reject option models for the OOD setup: Cost-based  
41 model, bounded TPR-FPR model, and Bounded Precision-Recall model. These models extend the  
42 standard rejection models used in the non-OOD setup [2, 15] and define the notion of an optimal OOD  
43 classifier. (ii) We establish that all the proposed models, despite their different formulations, share  
44 a common class of optimal strategies. The optimal OOD selective classifier combines a Bayes ID  
45 classifier with a selective function based on a linear combination of the conditional risk and likelihood  
46 ratio of the OOD and ID samples. This selective function enables a trade-off between distinguishing  
47 ID from OOD samples and detecting misclassifications. (iii) Motivated by the optimal strategy,  
48 we introduce double-score OOD methods that leverage uncertainty scores from two chosen OOD  
49 detectors: one focused on OOD/ID discrimination and the other on misclassification detection. We  
50 show experimentally that this simple strategy consistently outperforms the state-of-the-art. (iv) We  
51 review existing metrics for evaluation of OOD methods and show that they provide incomplete  
52 view, if used separately, or inconsistent view of the evaluated methods, if used together. We propose  
53 novel evaluation metrics derived from the definition of optimal strategy under the proposed OOD  
54 rejection models. These new metrics provide a comprehensive and reliable assessment of OOD  
55 methods without the deficiencies observed in existing approaches.

## 56 2 Reject option models for OOD setup

57 The terminology of ID and OOD samples comes from the setups when the training set contains only  
58 ID samples, while the test set contains a mixture of ID and OOD samples. In this paper, we analyze  
59 which prediction strategies are optimal on the test samples, but we do not address the problem of  
60 learning such strategy. We follow the OOD setup from [5]. Let  $\mathcal{X}$  be a set of observable inputs (or  
61 features), and  $\mathcal{Y}$  a finite set of labels that can be assigned to in-distribution (ID) inputs. ID samples  
62  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  are generated from a joint distribution  $p_I: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ . Out-of-distribution (OOD)  
63 samples  $x$  are generated from a distribution  $p_O: \mathcal{X} \rightarrow \mathbb{R}_+$ . ID and OOD samples share the same  
64 input space  $\mathcal{X}$ . Let  $\emptyset$  be a special label to mark the OOD sample. Let  $\bar{\mathcal{Y}} = \mathcal{Y} \cup \{\emptyset\}$  be an extended  
65 set of labels. In the testing stage the samples  $(x, \bar{y}) \in \mathcal{X} \times \bar{\mathcal{Y}}$  are generated from the joint distribution  
66  $p: \mathcal{X} \times \bar{\mathcal{Y}} \rightarrow \mathbb{R}_+$  defined as a mixture of ID and OOD:

$$p(x, \bar{y}) = \begin{cases} p_O(x) \pi & \text{if } \bar{y} = \emptyset \\ p_I(x, \bar{y}) (1 - \pi) & \text{if } \bar{y} \in \mathcal{Y} \end{cases}, \quad (1)$$

67 where  $\pi \in [0, 1)$  is the probability of observing the OOD sample. Our OOD setup subsumes the  
68 standard non-OOD setup as a special case when  $\pi = 0$ , and the reject option models that will be  
69 introduced below will become for  $\pi = 0$  the known reject option models for the non-OOD setup.

70 Our goal is to design OOD selective classifier  $q: \mathcal{X} \rightarrow \mathcal{D}$ , where  $\mathcal{D} = \mathcal{Y} \cup \{\text{reject}\}$ , which either  
71 predicts a label,  $q(x) \in \mathcal{Y}$ , or it rejects the prediction,  $q(x) = \text{reject}$ , when (i) input  $x \in \mathcal{X}$  prevents  
72 accurate prediction of  $y \in \mathcal{Y}$  because it is noisy, or (ii) comes from OOD. We represent the selective  
73 classifier by the ID classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$ , and a stochastic selective function  $c: \mathcal{X} \rightarrow [0, 1]$  that  
74 outputs a probability that the input is accepted [7], i.e.,

$$q(x) = (h, c)(x) = \begin{cases} h(x) & \text{with probability } c(x) \\ \text{reject} & \text{with probability } 1 - c(x) \end{cases}. \quad (2)$$

75 In the following sections, we propose three reject option models that define the notion of the optimal  
76 OOD selective classifier of the form (2) applied to samples generated by (1).

### 77 2.1 Cost-based rejection model for OOD setup

78 A classical approach to define an optimal classifier is to formulate it as a loss minimization problem.  
79 This requires defining a loss  $\bar{\ell}: \bar{\mathcal{Y}} \times \mathcal{D} \rightarrow \mathbb{R}_+$  for each combination of the label  $\bar{y} \in \bar{\mathcal{Y}} = \mathcal{Y} \cup \{\emptyset\}$   
80 and the output of the classifier  $q(x) \in \mathcal{D} = \mathcal{Y} \cup \{\text{reject}\}$ . Let  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  be some application-  
81 specific loss on ID samples, e.g., 0/1-loss or MAE. Furthermore, we need to define the loss for the  
82 case where the input is OOD sample  $\bar{y} = \emptyset$  or the classifier rejects  $q(x) = \text{reject}$ . Let  $\varepsilon_1 \in \mathbb{R}_+$  be  
83 the loss for rejecting the ID sample,  $\varepsilon_2 \in \mathbb{R}_+$  loss for prediction on the OOD sample, and  $\varepsilon_3 \in \mathbb{R}_+$   
84 loss for correctly rejecting the OOD sample.  $\ell, \varepsilon_1, \varepsilon_2$  and  $\varepsilon_3$  can be arbitrary, but we assume that

85  $\varepsilon_2 > \varepsilon_3$ . The loss  $\bar{\ell}$  is then:

$$\bar{\ell}(\bar{y}, q) = \begin{cases} \ell(\bar{y}, q) & \text{if } \bar{y} \in \mathcal{Y} \wedge q \in \mathcal{Y} \\ \varepsilon_1 & \text{if } \bar{y} \in \mathcal{Y} \wedge q = \text{reject} \\ \varepsilon_2 & \text{if } \bar{y} = \emptyset \wedge q \in \mathcal{Y} \\ \varepsilon_3 & \text{if } \bar{y} = \emptyset \wedge q = \text{reject} \end{cases} \quad (3)$$

86 Having the loss  $\bar{\ell}$ , we can define the optimal OOD selective classifier as a minimizer of the expected  
87 risk  $R(h, c) = \mathbb{E}_{x, y \sim p(x, \bar{y})} \bar{\ell}(\bar{y}, (h, c)(x))$ .

88 **Definition 1** (*Cost-based OOD model*) An optimal OOD selective classifier  $(h_C, c_C)$  is a solution to  
89 the minimization problem  $\min_{h, c} R(h, c)$  where we assume that both minimizers exist.

90 An optimal solution of the cost-based OOD model requires three components: The Bayes ID classifier

$$h_B(x) \in \underset{y' \in \mathcal{Y}}{\text{Argmin}} \sum_{y \in \mathcal{Y}} p_I(y | x) \ell(y, y'), \quad (4)$$

91 its conditional risk  $r_B(x) = \sum_{y \in \mathcal{Y}} p_I(y | x) \ell(y, h_B(x))$ , and the likelihood ratio of the OOD and  
92 ID inputs,  $g(x) = \frac{p_O(x)}{p_I(x)}$ , which we defined to be  $g(x) = \infty$  for  $p_I(x) = 0$ .

93 **Theorem 1** An optimal selective classifier  $(h_C, c_C)$  under the cost-based OOD model is composed  
94 of the Bayes classifier (4),  $h_C = h_B$ , and the selective function

$$c_C(x) = \begin{cases} 1 & \text{if } s_C(x) < \varepsilon_1 \\ \tau & \text{if } s_C(x) = \varepsilon_1 \\ 0 & \text{if } s_C(x) > \varepsilon_1 \end{cases} \quad \text{using the score } s_C(x) = r_B(x) + (\varepsilon_2 - \varepsilon_3) \frac{\pi}{1 - \pi} g(x) \quad (5)$$

95 where  $\tau$  is an arbitrary number in  $[0, 1]$ , and  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  are losses defining the extended loss (3).

96 Note that  $\tau$  can be arbitrary and therefore a deterministic selective function  $c_C(x) = \llbracket s_C(x) \leq \varepsilon_1 \rrbracket$  is  
97 also optimal. An optimal selective function accepts inputs based on the score  $s_C(x)$ , which is a linear  
98 combination of two functions, conditional risk  $r_B(x)$  and the likelihood ratio  $g(x) = p_O(x)/p_I(x)$ .

99 **Relation to cost-based model for Non-OOD setup** For  $\pi = 0$ , the cost-based OOD model reduces  
100 to the standard cost-based model of the reject option classifier in a non-OOD setup [2]. In the  
101 non-OOD setup, we do not need to specify the losses  $\varepsilon_2$  and  $\varepsilon_3$  and the risk  $R(h, c)$  simplifies  
102 to  $R'(h, c) = \mathbb{E}_{x, y \sim p_I(x, y)} [\ell(y, h(x)) c(x) + \varepsilon_1 (1 - c(x))]$ . The well-known optimal solution  
103 is composed of the Bayes classifier  $h_B(x)$  as in the OOD case; however, the selection function  
104  $c'_C(x) = \llbracket r(x) \leq A \rrbracket$  accepts the input solely based on the conditional risk  $r(x)$ .

## 105 2.2 Bounded TPR-FPR rejection model

106 The cost-based OOD model requires the classification loss  $\ell$  for ID samples and defining the costs  $\varepsilon_1,$   
107  $\varepsilon_2, \varepsilon_3$  which is difficult in practice because the physical units of  $\ell$  and  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  are often different.  
108 In this section, we propose an alternative approach which requires only the classification loss  $\ell$  while  
109 costs  $\varepsilon_1, \varepsilon_2, \varepsilon_3$  are replaced by constraints on the performance of the selective function.

110 The selective function  $c: \mathcal{X} \rightarrow [0, 1]$  can be seen as a discriminator of OOD/ID samples. Let  
111 us consider ID and OOD samples as positive and negative classes, respectively. We introduce  
112 three metrics to measure the performance of the OOD selective classifier  $(h, c)$ . We measure the  
113 performance of selective function by the True Positive Rate (TPR) and the False Positive Rate (FPR).  
114 The TPR is defined as the probability that ID sample is accepted by the selective function  $c$ , i.e.,

$$\phi(c) = \int_{\mathcal{X}} p(x | \bar{y} \neq \emptyset) c(x) dx = \int_{\mathcal{X}} p_I(x) c(x) dx. \quad (6)$$

115 The FPR is defined as the probability that OOD sample is accepted by the selective function  $c$ , i.e.,

$$\rho(c) = \int_{\mathcal{X}} p(x | \bar{y} = \emptyset) c(x) dx = \int_{\mathcal{X}} p_O(x) c(x) dx. \quad (7)$$

116 The second identity in (6) and (7) is obtained after substituting the definition of  $p(x, \bar{y})$  from (1).  
 117 Lastly, we characterize the performance of the ID classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$  by the selective risk

$$\mathbb{R}^S(h, c) = \frac{\int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} p_I(x, y) \ell(h(x), y) c(x) dx}{\phi(c)}$$

118 defined for non-zero  $\phi(c)$ , i.e., the expected loss of the classifier  $h$  calculated on the ID samples  
 119 accepted by the selective function  $c$ .

120 **Definition 2 (Bounded TPR-FPR model)** Let  $\phi_{\min} \in [0, 1]$  be the minimal acceptable TPR and  
 121  $\rho_{\max} \in [0, 1]$  maximal acceptable FPR. An optimal OOD selective classifier  $(h_T, c_T)$  under the  
 122 bounded TPR-FPR model is a solution of the problem

$$\min_{h \in \mathcal{Y}^{\mathcal{X}}, c \in [0, 1]^{\mathcal{X}}} \mathbb{R}^S(h, c) \quad \text{s.t.} \quad \phi(c) \geq \phi_{\min} \quad \text{and} \quad \rho(c) \leq \rho_{\max}, \quad (8)$$

123 where we assume that both minimizers exist.

124 **Theorem 2** Let  $(h, c)$  be an optimal solution to (8). Then  $(h_B, c)$ , where  $h_B$  is the Bayes ID  
 125 classifier (4), is also optimal to (8).

126 According to Theorem 2, the Bayes ID classifier  $h_B$  is an optimal solution to (8) that defines the  
 127 bounded TPR-FPR model. This is not surprising, but it is a practically useful result, because it allows  
 128 one to solve (8) in two consecutive steps: First, set  $h_T$  to the Bayes ID classifier  $h_B$ . Second, when  
 129  $h_T$  is fixed, the optimal selection function  $c_T$  is obtained by solving (8) only w.r.t.  $c$  which boils  
 130 down to:

131 **Problem 1 (Bounded TPR-FPR model for known  $h(x)$ )** Given ID classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$ , the opti-  
 132 mal selective function  $c^*: \mathcal{X} \rightarrow [0, 1]$  is a solution to

$$\min_{c \in [0, 1]^{\mathcal{X}}} \mathbb{R}^S(h, c) \quad \text{s.t.} \quad \phi(c) \geq \phi_{\min}, \quad \text{and} \quad \rho(c) \leq \rho_{\max}.$$

133 Problem 1 is meaningful even if  $h$  is not the Bayes ID classifier  $h_B$ . We can search for an optimal  
 134 selective function  $c^*(x)$  for any fixed  $h$ , which in practice is usually our best approximation of  $h_B$   
 135 learned from the data.

136 **Theorem 3** Let  $h: \mathcal{X} \rightarrow \mathcal{Y}$  be ID classifier and  $r: \mathcal{X} \rightarrow \mathbb{R}$  its conditional risk  $r(x) = \sum_{y \in \mathcal{Y}} p_I(y |$   
 137  $x) \ell(y, h(x))$ . Let  $g(x) = p_I(x) / p_I(x)$  be the likelihood ratio of ID and OOD samples. Then, the set  
 138 of optimal solutions of Problem 1 contains the selective classifier

$$c^*(x) = \begin{cases} 0 & \text{if } s(x) > \lambda \\ \tau(x) & \text{if } s(x) = \lambda \\ 1 & \text{if } s(x) < \lambda \end{cases} \quad \text{using score } s(x) = r(x) + \mu g(x) \quad (9)$$

139 where decision threshold  $\lambda \in \mathbb{R}$ , and multiplier  $\mu \in \mathbb{R}$  are constants and  $\tau: \mathcal{X} \rightarrow [0, 1]$  is a function  
 140 implicitly defined by the problem parameters.

141 The optimal  $c^*(x)$  is based on the score composed of a linear combination of  $r(x)$  and  $g(x)$  as in the  
 142 case of the cost-based model (5). Unlike the cost-based model, the acceptance probability  $\tau(x)$  for  
 143 boundary inputs  $\mathcal{X}_{s(x)=\lambda} = \{x \in \mathcal{X} \mid s(x) = \lambda\}$  cannot be arbitrary, in general. However, if  $\mathcal{X}$  is  
 144 continuous, the set  $\mathcal{X}_{s(x)=\lambda}$  has probability measure zero, up to some pathological cases, and  $\tau(x)$   
 145 can be arbitrary, i.e., the deterministic  $c^*(x) = \mathbb{I}[s(x) \leq \lambda]$  is optimal. If  $\mathcal{X}$  is finite, the value of  
 146  $\tau(x)$  can be found by linear programming. The linear program and more details on the form of  $\tau(x)$   
 147 are in the Appendix.

148 **Relation to Bounded-Abstention model for the non-OOD setup** For  $\pi = 0$ , the bounded TPR-  
 149 FPR model reduces to the bounded-abstention option model for non-OOD setup [15]. Namely,  
 150  $\rho(c) \leq \rho_{\max}$  can be removed because there are no OOD samples, and (8) becomes the bounded-  
 151 abstention model:  $\min_{h, c} \mathbb{R}^S(h, c)$ , s.t.  $\phi(c) \geq \phi_{\min}$ , which seeks the selective classifier with  
 152 guaranteed TPR and minimal selective risk. In the non-OOD setup, TPR is called *coverage*. An  
 153 optimal solution of the bounded abstention model [6], is composed of the Bayes ID classifier  $h_B$ , and  
 154 the same optimal selective function as the TPR-FPR model (9), however, with  $\mu = 0$  and  $\tau(x) = \tau$ ,  
 155  $\forall x \in \mathcal{X}$ , i.e., the score depends only on  $r(x)$  and an identical randomization is applied in all edge  
 156 cases [6]. Therefore,  $r(x)$  is the optimal score to detect misclassified ID samples in non-OOD setup  
 157 as it allows to achieve the minimal selective risk  $\mathbb{R}^S$  for any fixed coverage (TPR,  $\phi$ ).

158 **2.3 Bounded Precision-Recall rejection model**

159 The optimal selective classifier under the bounded TPR-FPR model does not depend on the prior of  
 160 the OOD samples  $\pi$ , which is useful, e.g., when  $\pi$  is unknown in the testing stage. In the case  $\pi$  is  
 161 known, it might be more suitable to constrain the precision rather than the FPR, while the constraint  
 162 on TPR remains the same. In the context of precision, we denote  $\phi(c)$  as recall instead of TPR. The  
 163 precision  $\kappa(c)$  is defined as the portion of samples accepted by  $c(x)$  that are actual ID samples, i.e.,

$$\kappa(c) = \frac{(1 - \pi) \int_{\mathcal{X}} p(x | \bar{y} \neq \emptyset) c(x) dx}{\int_{\mathcal{X}} p(x) c(x) dx} = \frac{(1 - \pi) \phi(c)}{\rho(c) \pi + \phi(c) (1 - \pi)}.$$

164 **Definition 3 (Bounded Precision-Recall model)** Let  $\kappa_{\min} \in [0, 1]$  be a minimal acceptable pre-  
 165 cision and  $\phi_{\min} \in [0, 1]$  minimal acceptable recall (a.k.a. TPR). An optimal selective classifier  
 166  $(h_P, c_P)$  under the bounded Precision-Recall model is a solution of the problem

$$\min_{h \in \mathcal{Y}^{\mathcal{X}}, c \in [0, 1]^{\mathcal{X}}} R^S(h, c) \quad \text{s.t.} \quad \phi(c) \geq \phi_{\min} \quad \text{and} \quad \kappa(c) \geq \kappa_{\min} \quad (10)$$

167 where we assume that both minimizers exist.

168 **Theorem 4** Let  $(h, c)$  be an optimal solution to (10). Then  $(h_B, c)$ , where  $h_B$  is the Bayes ID  
 169 classifier (4), is also optimal to (10).

170 Theorem 4 ensures that the Bayes ID classifier is an optimal solution to (10). After fixing  $h_P = h_B$ ,  
 171 the search for an optimal selective function  $c$  leads to:

172 **Problem 2 (Bounded Prec-Recall model for known  $h(x)$ )** Given ID classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$ , the  
 173 optimal selective function  $c^*: \mathcal{X} \rightarrow [0, 1]$  is a solution to

$$\min_{c \in [0, 1]^{\mathcal{X}}} R^S(h, c) \quad \text{s.t.} \quad \phi(c) \geq \phi_{\min} \quad \text{and} \quad \kappa(c) \geq \kappa_{\min}.$$

174 **Theorem 5** Let  $h: \mathcal{X} \rightarrow \mathcal{Y}$  be ID classifier and  $r: \mathcal{X} \rightarrow \mathbb{R}$  its conditional risk  $r(x) = \sum_{y \in \mathcal{Y}} p_I(y |$   
 175  $x) \ell(y, h(x))$ . Let  $g(x) = p_O(x) / p_I(x)$  be the likelihood ratio of OOD and ID samples. Then, the  
 176 set of optimal solutions of Problem 2 contains the selective function

$$c^*(x) = \begin{cases} 0 & \text{if } s(x) > \lambda \\ \tau(x) & \text{if } s(x) = \lambda \\ 1 & \text{if } s(x) < \lambda \end{cases} \quad \text{using the score } s(x) = r(x) + \mu g(x) \quad (11)$$

177 where detection threshold  $\lambda \in \mathbb{R}$ , and multiplier  $\mu \in \mathbb{R}$  are constants and  $\tau: \mathcal{X} \rightarrow [0, 1]$  is a function  
 178 implicitly defined by the problem parameters.

179 **2.4 Summary**

180 We proposed three rejection models for OOD setup which define the notion of optimal OOD selective  
 181 classifier: Cost-based model, Bounded TRP-FPR model, and Bounded Precision-Recall model. We  
 182 established that all three models, despite different formulation, share the class of optimal prediction  
 183 strategies. Namely, the optimal OOD selective classifier  $(h^*, c^*)$  is composed of the Bayes ID  
 184 classifier (4),  $h^* = h_B$ , and the selective function

$$c^*(x) = \begin{cases} 0 & \text{if } s(x) > \lambda \\ \tau(x) & \text{if } s(x) = \lambda \\ 1 & \text{if } s(x) < \lambda \end{cases} \quad \text{where } s(x) = r(x) + \mu g(x) \quad (12)$$

185 where  $\lambda$ ,  $\mu$ , and  $\tau(x)$  are specific for the used rejection model. However, in all cases, the optimal  
 186 uncertainty score  $s(x)$  for accepting the inputs is based on a linear combination of the conditional  
 187 risk  $r(x)$  of the ID classifier  $h^*$  and the OOD/ID likelihood ratio  $g(x) = p_O(x) / p_I(x)$ . On the other  
 188 hand, from the optimal solution of the well-known Neyman-Person problem [14], it follows that the  
 189 likelihood ratio  $g(x)$  is the optimal score of OOD/ID discrimination. Our results thus show that the  
 190 optimal OOD selective function needs to trade-off the ability to detect the misclassification of ID  
 191 samples and the ability to distinguish ID from OOD samples.

192 **Single-score vs. double-score OOD methods** The existing OOD methods, which we further  
193 call *single-score methods*, produce a classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$  and an uncertainty score  $s: \mathcal{X} \rightarrow \mathbb{R}$ . The  
194 score  $s(x)$  is used to construct a selective function  $c(x) = \llbracket s(x) \leq \lambda \rrbracket$  where  $\lambda \in \mathbb{R}$  is a decision  
195 threshold chosen in post-hoc evaluation. Hence, the existing methods effectively produce a set of  
196 selective classifiers  $\mathcal{Q} = \{(h, c) \mid c(x) = \llbracket s(x) \leq \lambda \rrbracket, \lambda \in \mathbb{R}\}$ . In contrast to existing methods, we  
197 established that the optimal selective function is always based on a linear combination of two scores:  
198 conditional risk  $r(x)$  and likelihood ratio  $g(x)$ . Therefore, we propose the *double-score method*,  
199 which in addition to a classifier  $h(x)$ , produces two scores,  $s_r: \mathcal{X} \rightarrow \mathbb{R}$  and  $s_g: \mathcal{X} \rightarrow \mathbb{R}$ , and uses  
200 their combination  $s(x) = s_r(x) + \mu s_g(x)$  to accept inputs. Formally, the double-score method  
201 produces a set of selective classifiers  $\mathcal{Q} = \{(h, c) \mid c(x) = \llbracket s_r(x) + \mu s_g(x) \leq \lambda \rrbracket, \mu \in \mathbb{R}, \lambda \in \mathbb{R}\}$ .  
202 The double-score strategy can be used to leverage uncertainty scores from two chosen OOD  
203 methods: one focused on OOD/ID discrimination and the other on misclassification detection.

### 204 3 Post-hoc tuning and evaluation metrics

205 Let  $\mathcal{T} = ((x_i, \bar{y}_i) \in \mathcal{X} \times \bar{\mathcal{Y}} \mid i = 1, \dots, n)$  be a set of validation examples i.i.d. drawn from a  
206 distribution  $p(x, \bar{y})$ . Given a set of selective classifiers  $\mathcal{Q}$ , trained by the single-score or double-score  
207 OOD method, the goal of the post-hoc tuning is to use  $\mathcal{T}$  to select the best selective classifier  
208  $(h_n, c_n) \in \mathcal{Q}$  and estimate its performance on unseen samples generated from the same  $p(x, \bar{y})$ . This  
209 task requires a notion of an optimal selective classifier which we defined by the proposed rejection  
210 models. In Sec 3.2 and Sec 3.3, we propose the post-hoc tuning and evaluation metrics for the  
211 Bounded TPR-FPR and Bounded Precision-Recall models, respectively. In Sec 3.4 we review the  
212 existing evaluation metrics for OOD methods and point out their deficiencies. We will exemplify  
213 the proposed metrics on synthetic data and OOD methods described in Sec 3.1.

#### 214 3.1 Synthetic data and exemplar single-score and double-score OOD methods

215 Let us consider a simple 1-D setup. The input space is  $\mathcal{X} = \mathbb{R}$  and there are three ID labels  
216  $\mathcal{Y} = \{1, 2, 3\}$ . ID samples are generated from  $p_I(x, 1) = 0.3\mathcal{N}(x; -1, 1)$ ,  $p_I(x, 2) = 0.3\mathcal{N}(x; 1, 1)$ ,  
217  $p_I(x, 3) = 0.4\mathcal{N}(x; 3, 1)$ , where  $\mathcal{N}(x; \mu, \sigma)$  is normal distribution with mean  $\mu$  and variance  $\sigma$ .  
218 OOD is the normal distribution  $p_O(x) = \mathcal{N}(x; 3, 0.2)$ , and the OOD prior  $\pi = 0.25$ . We use  
219 0/1-loss  $\ell(y, y') = \llbracket y \neq y' \rrbracket$ , i.e.,  $\mathbb{R}^S$  is the classification error on accepted inputs. The known ID  
220 and OOD allows us to evaluate the Bayes ID classifier  $h_B(x)$  by (4), its conditional risk  $r_B(x) =$   
221  $\min_{y' \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} p_I(y \mid x) \ell(y, y')$  and the OOD/ID likelihood ratio  $g(x) = p_O(x)/p_I(x)$ .

222 We consider 3 exemplar single-score OOD methods A, B, C. The methods produce the same optimal  
223 classifier  $h^*(x)$  and the selective functions  $c(x) = \llbracket r_B(x) + \mu g(x) \leq \lambda \rrbracket$  with a different setting of  
224  $\mu$ . I.e., the method  $k \in \{A, B, C\}$  produces the set of selective classifiers  $\mathcal{Q}_k = \{(h^*(x), c(x)) \mid$   
225  $c(x) = \llbracket r_B(x) + \mu_k g(x) \leq \lambda \rrbracket, \lambda \in \mathbb{R}\}$ , where the constant  $\mu_k$  is defined as follows:

- 226 • Method A( $\infty$ ):  $\mu = \infty$ ,  $s(x) = g(x)$ . This corresponds to the optimal OOD/ID discriminator.
- 227 • Method B(0.2):  $\mu = 0.2$ ,  $s(x) = r_B(x) + 0.2g(x)$ . Combination of method A and C.
- 228 • Method C(0):  $\mu = 0$ ,  $s(x) = r_B(x)$ . This corresponds to the optimal misclassification detector.

229 We also consider a double-score method, Method D( $\mathbb{R}$ ), which outputs the same optimal classifier  
230  $h_*(x)$ , and scores  $s_r(x) = r(x)$  and  $s_g(x) = g(x)$ . I.e., Method D( $\mathbb{R}$ ) produces the set of selective  
231 classifiers  $\mathcal{Q}_D = \{(h^*(x), c(x)) \mid c(x) = \llbracket r(x) + \mu g(x) \leq \lambda \rrbracket, \mu \in \mathbb{R}, \lambda \in \mathbb{R}\}$ . Note that we have  
232 shown that  $\mathcal{Q}_D$  contains an optimal selective classifier regardless of the reject option model used.

#### 233 3.2 Bounded TPR-FPR rejection model

234 The bounded TPR-FPR model is defined using the selective risk  $\mathbb{R}^S(h, c)$ , TPR  $\phi(c)$  and FPR  $\rho(c)$   
235 the value of which can be estimated from the validation set  $\mathcal{T}$  as follows:

$$\mathbb{R}_n^S(h, c) = \frac{\sum_{i \in \mathcal{I}_I} \ell(y_i, h(x_i)) c(x_i)}{\sum_{i \in \mathcal{I}_I} c(x_i)}, \quad \phi_n(h, c) = \frac{1}{|\mathcal{I}_I|} \sum_{i \in \mathcal{I}_I} c(x_i), \quad \rho_n(h, c) = \frac{1}{|\mathcal{I}_O|} \sum_{i \in \mathcal{I}_O} c(x_i)$$

236 where  $\mathcal{I}_I = \{i \in \{1, \dots, n\} \mid \bar{y}_i \neq \emptyset\}$  and  $\mathcal{I}_O = \{i \in \{1, \dots, n\} \mid \bar{y}_i = \emptyset\}$  are indices of ID and  
237 OOD samples in  $\mathcal{T}$ , respectively.

Method	Proposed metrics		Existing metrics		
	TPR-FPR model	Prec-Recall model	↑ Existing metrics		
	↓ Selective risk at TPR(0.7),FPR(0.2)	↓ Selective risk at Prec(0.9),Recall(0.7)	AUROC	AUPR	OSCR
A( $\infty$ )	0.157	0.157	<b>0.88</b>	<b>0.96</b>	0.82
B(0.2)	0.143	0.143	0.86	0.95	0.83
C(0)	unable	unable	0.76	0.92	<b>0.86</b>
D( $\mathbb{R}$ ) proposed	<b>0.133</b>	<b>0.129</b>	<b>0.88</b>	<b>0.96</b>	<b>0.86</b>

Table 1: Evaluation of the exemplar single-score methods A, B, C and the proposed double-score method D on synthetic data using the proposed metrics and the existing ones. The selective risk corresponds to the classification error on accepted ID samples.

238 Given the target TPR  $\phi_{min} \in (0, 1]$  and FPR  $\rho_{max} \in (0, 1]$ , the best selective classifier  $(h_n, c_n)$  out  
239 of  $\mathcal{Q}$  is found by solving:

$$(h_n, c_n) \in \underset{(h,c) \in \mathcal{Q}}{\text{Argmin}} R_n^S(h, c) \quad \text{s.t.} \quad \phi_n(h, c) \geq \phi_{min}, \quad \text{and} \quad \rho_n(h, c) \leq \rho_{max}. \quad (13)$$

240 **Proposed evaluation metric** If problem (13) is feasible,  $R_n^S(h_n, c_n)$  is reported as the performance  
241 estimator of OODD method producing  $\mathcal{Q}$ . Otherwise, the method is marked as unable to achieve  
242 the target TPR and FPR. Tab. 1 shows the selective risk for the methods A-D at the target TPR  
243  $\phi_{min} = 0.7$  and FPR  $\rho_{max} = 0.2$ . The minimal  $R_n^S$  is achieved by method D( $\mathbb{R}$ ), followed by B(0.2)  
244 and A( $\infty$ ), while C(0) is unable to achieve the target TPR and FPR. One can visualize  $R_n^S$  in a range  
245 of operating points while bounding only  $\rho_{max}$  or  $\phi_{min}$ . E.g., by fixing  $\rho_{max}$  we can plot  $R_n^S$  as  
246 a function of attainable values of  $\phi_n$  by which we obtain the Risk-Coverage curve, known from  
247 non-OOD setup, at  $\rho_{max}$ . Recall that TPR is coverage. See Appendix for Risk-Coverage curve at  
248  $\rho_{max}$  for methods A-D.

249 **ROC curve** The problem (13) can be infeasible. To choose a feasible target on  $\phi_{min}$  and  $\rho_{max}$ , it  
250 is advantageous to plot the ROC curve, i.e., values of TPR and FPR attainable by the classifiers in  $\mathcal{Q}$ .  
251 For single-score methods, the ROC curve is a set of points obtained by varying the decision threshold:  
252  $\text{ROC}(\mathcal{Q}) = \{(\phi_n(h, c), \rho_n(h, c)) \mid c(x) = \llbracket s(x) \leq \lambda \rrbracket, \lambda \in \mathbb{R}\}$ . In case of double-score methods,  
253 we vary  $\rho_{max} \in [0, 1]$  and for each  $\rho_{max}$  we choose the maximal feasible  $\phi_n$ . I.e., ROC curve  
254 is  $\text{ROC}(\mathcal{Q}) = \{(\phi, \rho_{max}) \mid \phi = \max_{(h,c) \in \mathcal{Q}} \phi_n(h, c) \text{ s.t. } \rho_n(h, c) \leq \rho_{max}, \rho_{max} \in [0, 1]\}$ .  
255 See Appendix for ROC curve of the methods A-D. In Tab. 1 we report the Area Under ROC curve  
256 (AUROC) which is a commonly used summary of the entire ROC curve. The highest AUROC  
257 achieved Methods A( $\infty$ ) and E( $\mathbb{R}$ ). Recall that Method A( $\infty$ ) uses the optimal ID/OOD discriminator  
258 and the proposed Method E( $\mathbb{R}$ ) subsumes A( $\infty$ ).

### 259 3.3 Bounded Precision-Recall rejection model

260 Let  $\kappa_n(c) = (1 - \pi) \phi_n(c) / ((1 - \pi) \phi_n(c) + \pi \rho_n(c))$  be the sample precision of the selective  
261 function  $c$ . Given the target recall  $\phi_{min} \in (0, 1]$  and precision  $\kappa_{min} \in (0, 1]$ , the best selective  
262 classifier  $(h_n, c_n)$  out of  $\mathcal{Q}$  is found by solving

$$(h_n, c_n) \in \underset{(h,c) \in \mathcal{Q}}{\text{Argmin}} R_n^S(h, c) \quad \text{s.t.} \quad \phi_n(h, c) \geq \phi_{min}, \quad \kappa_n(h, c) \geq \kappa_{min}. \quad (14)$$

263 **Proposed evaluation metric** If problem (14) is feasible,  $R_n^S(h_n, c_n)$  is reported as the performance  
264 estimator of OODD method which produced  $\mathcal{Q}$ . Otherwise, the method is marked as unable to achieve  
265 the target Precision/Recall. Tab. 1 shows the selective risk for the methods A-D at the Precision  
266  $\kappa_{min} = 0.9$  and recall  $\phi_{max} = 0.7$ . The minimal  $R_n^S$  is achieved by the proposed method D( $\mathbb{R}$ ),  
267 followed by B(0.2) and A( $\infty$ ), while method C(0) is unable to achieve the target Precision/Recall.  
268 Note that single-score methods A-C achieve the same  $R_n^S$  under both TPR-FPR and Prec-Recall  
269 models while the results for double-score method D( $\mathbb{R}$ ) differ. The reason is that both models share  
270 the same constraint  $\phi_n \geq 0.7$  (TPR is Recall) which is active, while the other two constraints are not  
271 active because  $R_n^S$  is a monotonic function w.r.t. the value of the decision threshold.

272 **Precision-Recall (PR) curve** To choose feasible bounds on  $\kappa_{min}$  and  $\phi_{min}$  before solving (14),  
 273 one can plot the PR curve, i.e., the values of precision and recall attainable by the classifiers in  
 274  $\mathcal{Q}$ . For single-score methods, the PR curve is a set of points obtained by varying the decision  
 275 threshold:  $\text{PR}(\mathcal{Q}) = \{(\kappa_n(h, c), \phi_n(h, c)) \mid c(x) = \llbracket s(x) \leq \lambda \rrbracket, \lambda \in \mathbb{R}\}$ . In case of double-  
 276 score methods, we vary  $\phi_{min} \in [0, 1]$  and for each  $\phi_{min}$  we choose the maximal feasible  $\kappa_n$ , i.e.,  
 277  $\text{PR}(\mathcal{Q}) = \{(\kappa, \phi_{min}) \mid \kappa = \max_{(h, c) \in \mathcal{Q}} \kappa_n(h, c) \text{ s.t. } \phi_n(h, c) \geq \phi_{min}, \phi_{min} \in [0, 1]\}$ . See  
 278 Appendix for PR curve of the methods A-D. We compute the Area Under the PR curve and report it  
 279 for Methods A-D in Tab. 1. Rankings of the methods w.r.t AUPR and AUROC are the same.

### 280 3.4 Shortcomings of existing evaluation metrics

281 The most commonly used metrics to evaluate OOD methods are the AUROC and AUPR [10, 13,  
 282 3, 12, 1, 16]. Both metrics measure the ability of the selective function  $c(x)$  to distinguish ID from  
 283 OOD samples. AUROC and AUPR are often the only metrics reported although they completely  
 284 ignore the performance of the ID classifier. Our synthetic example shows that high AUROC/AUPR  
 285 is not a precursor of a good OOD selective classifier. E.g., Method A( $\infty$ ), using optimal OOD/ID  
 286 discriminator, attains the highest (best) AUROC and AUPR (see Tab. 1), however, at the same time  
 287 Method A( $\infty$ ) achieves the highest (worst)  $R_n^S$  under both rejection models, and it is also the worst  
 288 misclassification detector according to the OSCR score defined below.

289 The performance of the ID classifier  $h(x)$  is usually evaluated by the ID classification accuracy  
 290 (a.k.a. closed set accuracy) [13, 3] and by the OSCR score [4, 8, 1]. The ID accuracy measures  
 291 the performance of  $h(x)$  assuming all inputs are accepted, i.e.,  $c(x) = 1, \forall x \in \mathcal{X}$ , hence it says  
 292 nothing about the performance on the actually accepted samples like  $R_n^S$ . E.g., Methods A-D in our  
 293 synthetic example use the same classifier  $h(x)$  and hence have the same ID accuracy, however, they  
 294 perform quite differently in terms of the other more relevant metrics, like  $R_n^S$  or OSCR. The OSCR  
 295 score is defined as the area under CCR versus FPR curve [21], where the CCR stands for the correct  
 296 classification rate on the accepted ID samples; in case of 0/1-loss  $\text{CCR} = 1 - R_n^S$ . The CCR-FPR  
 297 curve evaluates the performance of the ID classifier on the accepted samples, but it ignores the ability  
 298 of  $c(x)$  to discriminate OOD and ID samples as it does not depend on TPR. E.g., Method D(0), using  
 299 the optimal misclassification detector, achieves the highest (best) OSCR score; however, at the same  
 300 time, it has the lowest (worst) AUROC and AUPR.

301 Other, less frequently used metrics involve: F1-score,  $\text{FPR@TPRx}$ ,  $\text{TNR@TPRx}$ ,  $\text{CCR@FPRx}$  [10, 8,  
 302 1, 21, 16]. All these metrics are derived from either ROC, PR or CCR-FPR curve, and hence they  
 303 suffer with the same conceptual problems as AUROC, AUPR and OSCR, respectively.

304 We argue that the existing metrics evaluate only one aspect of the OOD selective classifier, namely,  
 305 either the ability to discriminate ID from OOD samples, or the performance of ID classifier on the  
 306 accepted (or on possibly all) ID samples. We show that in principle there can be methods that are best  
 307 OOD/ID discriminators but the worst misclassification detectors and vice versa. Therefore, using  
 308 individual metrics can (and often does) provide inconsistent ranking of the evaluated methods.

### 309 3.5 Summary

310 We propose novel evaluation metrics derived from the definition of the optimal strategy under the  
 311 proposed OOD rejection models. The proposed metrics simultaneously evaluate the classification  
 312 performance on the accepted ID samples and they guarantee the performance of the OOD/ID discrimi-  
 313 nator, either via constraints in TPR-FPR or Precision-Recall pair. Advantages of the proposed metrics  
 314 come at a price. Namely, we need to specify feasible target TPR and FPR, or Precision and Recall,  
 315 depending on the model used. However, feasible values of TPR-FPR and Prec-Recall pairs can be  
 316 easily read out of the ROC and PR curve, respectively. We argue that setting these extra parameters is  
 317 better than using the existing metrics that provide incomplete, if used separately, or inconsistent, if  
 318 used in combination, view of the evaluated methods.

319 Another issue is solving the problems (13) and (14) to compute the proposed evaluation metrics and  
 320 figures. Fortunately, both problems lead to optimization w.r.t one or two variables in case of the  
 321 single-score and double-score methods, respectively. A simple and efficient algorithm to solve the  
 322 problems in  $\mathcal{O}(n \log n)$  time is provided in Appendix.



Method	OOD: notmnist			OOD: fashionmnist			OOD: cifar10			
	↓ S. risk at TPR(0.80) FPR(0.08)	↑ AUROC	↑ OSCR	↓ S. risk at TPR(0.80) FPR(0.10)	↑ AUROC	↑ OSCR	↓ S. risk at TPR(0.80) FPR(0.29)	↑ AUROC	↑ OSCR	
ID: mnist	MSP [10]	<b>0.00014</b>	0.936	<b>0.996</b>	<b>0.00013</b>	0.956	<b>0.994</b>	<b>0.00013</b>	0.989	0.991
	MLS [9]	0.00139	0.941	0.993	0.00139	0.972	0.991	0.00139	<b>0.993</b>	0.990
	ODIN [11]	0.00069	0.942	0.993	0.00069	0.970	0.991	0.00069	0.993	0.990
	REACT [17]	0.00637	0.962	0.991	0.00637	<b>0.985</b>	0.990	0.00637	0.992	0.989
	KNN [19]	0.00041	0.976	0.991	0.00041	0.947	0.993	0.00041	0.976	0.991
	VIM [20]	0.00193	<b>0.983</b>	0.990	0.00194	0.926	0.993	0.00194	0.860	<b>0.995</b>
	KNN+MSP	<b>0.00000</b>	0.976	<b>0.996</b>	<b>0.00000</b>	0.962	<b>0.994</b>	<b>0.00000</b>	0.991	0.991
VIM+MSP	0.00014	<b>0.987</b>	<b>0.996</b>	0.00013	0.976	<b>0.994</b>	0.00013	0.992	<b>0.995</b>	
Method	OOD: cifar100			OOD: tiny imagenet			OOD: mnist			
	↓ S. risk at TPR(0.80) FPR(0.21)	↑ AUROC	↑ OSCR	↓ S. risk at TPR(0.80) FPR(0.19)	↑ AUROC	↑ OSCR	↓ S. risk at TPR(0.80) FPR(0.19)	↑ AUROC	↑ OSCR	
ID: cifar10	MSP [10]	0.00676	0.871	<b>0.977</b>	0.00676	0.887	<b>0.976</b>	0.00676	0.899	<b>0.976</b>
	MLS [9]	0.00984	0.861	0.973	0.00984	0.885	0.971	0.00984	0.905	0.971
	ODIN [11]	0.01000	0.851	0.975	0.01000	0.864	0.974	0.00995	0.915	0.969
	REACT [17]	0.00856	0.864	0.973	0.00856	0.888	0.971	0.00856	0.883	0.972
	KNN [19]	<b>0.00665</b>	<b>0.896</b>	0.974	<b>0.00665</b>	<b>0.914</b>	0.972	<b>0.00665</b>	<b>0.916</b>	0.973
	VIM [20]	0.01232	0.872	0.972	0.01232	0.888	0.971	0.01236	0.873	0.974
	KNN+MSP	<b>0.00652</b>	<b>0.896</b>	<b>0.977</b>	<b>0.00652</b>	<b>0.914</b>	<b>0.976</b>	<b>0.00652</b>	<b>0.916</b>	<b>0.976</b>
VIM+MSP	0.00676	0.879	<b>0.977</b>	0.00676	0.894	<b>0.976</b>	0.00676	0.900	<b>0.976</b>	

Table 2: Evaluation of existing single-score methods MSP, MLS, ODIN, REACT, KNN and two instances of the proposed double-score strategy: KNN+MSP and VIM+MSP. We use MNIST (top table) and CIFAR10 (bottom table) as ID, and three different datasets as OOD. We report the standard AUROC and OSCR, and the proposed selective risk at target TPR and FPR, where the selective risk corresponds to the classification error on accepted ID samples. Best results are in bold.

## 4 Experiments

In this section, we evaluate single-score OOD methods and the proposed double-score strategy, using the existing and the proposed evaluation metrics. We use MSP [10], MLS [9], ODIN [11] as baselines and REACT [17], KNN [19], VIM [20] as representatives of recent single-score approaches. We evaluate two instances of the double-score strategy. First, we combine the scores of MSP [10] and KNN [18] and, second, scores of MSP and VIM [20]. MSP score is asymptotically the best misclassification detector, while KNN and VIM are two best OOD/ID discriminators according to their AUROC. We always use the ID classifier of the MSP method. The evaluation data and implementations of OOD methods are taken from OpenOOD benchmark [21]. Because the datasets have unrealistically high portion of OOD samples, e.g.,  $\pi > 0.5$ , we use metrics that do not depend on  $\pi$ . Namely, AUROC and OSCR as the most frequently used metrics, and the proposed selective risk at TPR and FPR. We use 0/1-loss, hence the reported selective risk is the classification error on accepted ID samples with guaranteed TPR and FPR. In all experiments we fix the target TPR to 0.8 while FPR is set for each database to the highest FPR attained by all compared methods.

Results are presented in Tab. 2. It is seen that the single-score methods with the highest AUROC and OSCR are always different, which prevents us to create a single conclusive ranking of the evaluated approaches. MSP is almost consistently the best misclassification detector according to OSCR. The best OOD/ID discriminator is, according to AUROC, one of the recent methods: REACT, KNN, or VIM. The proposed double-score strategy, KNN+MSP and VIM+MSP, consistently outperforms the other approaches in all metrics.

## 5 Conclusions

This paper introduces novel reject option models which define the notion of the optimal prediction strategy for OOD setups. We prove that all models, despite their different formulations, share the same class of optimal prediction strategies. The main insight is that the optimal prediction strategy must trade-off the ability to detect misclassified examples and to distinguish ID from OOD samples. This is in contrast to existing OOD methods that output a single uncertainty score. We propose a simple and effective double-score strategy that allows us to boost performance of two existing OOD methods by combining their uncertainty scores. Finally, we suggest improved evaluation metrics for assessing OOD methods that simultaneously evaluate all aspects of the OOD methods and are directly related to the optimal OOD strategy under the proposed reject option models.

353 **References**

- 354 [1] Guanyao Chen, Peixi Peng, Xiangqian Wang, and Yonghong Tian. Adversarial reciprocal points learning  
355 for open set recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8065–  
356 8081, 2022.
- 357 [2] C. Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*,  
358 16(1):41–46, 1970.
- 359 [3] Terrance DeVries and Graham W Taylor. Learning confidence for out-of-distribution detection in neural  
360 networks. *arXiv preprint arXiv:1802.04865*, 2018.
- 361 [4] Akshay Raj Dhamija, Manuel Günther, and Terrance Boult. Reducing network agnostophobia. In S. Bengio,  
362 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural*  
363 *Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 364 [5] Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. Is out-of-distribution detection  
365 learnable? In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in*  
366 *Neural Information Processing Systems*, volume 35, pages 37199–37213. Curran Associates, Inc., 2022.
- 367 [6] Vojtech Franc, Daniel Prusa, and Vaclav Voracek. Optimal strategies for reject option classifiers. *Journal*  
368 *of Machine Learning Research*, 24(11):1–49, 2023.
- 369 [7] Y. Geifman and R. El-Yaniv. Selective classification for deep neural networks. In *Advances in Neural*  
370 *Information Processing Systems 30*, pages 4878–4887, 2017.
- 371 [8] Federica Granese, Marco Romanelli, Daniele Gorla, Catuscia Palamidessi, and Pablo Piantanida. Doctor:  
372 A simple method for detecting misclassification errors. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S.  
373 Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34,  
374 pages 5669–5681. Curran Associates, Inc., 2021.
- 375 [9] Dan Hendrycks, Steven Basart, Mantas Mazeika, Andy Zou, Joseph Kwon, Mohammadreza Mostajabi,  
376 Jacob Steinhardt, and Dawn Song. Scaling out-of-distribution detection for real-world settings. In Kamalika  
377 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings*  
378 *of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine*  
379 *Learning Research*, pages 8759–8773. PMLR, 17–23 Jul 2022.
- 380 [10] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples  
381 in neural networks. In *Proceedings of International Conference on Learning Representations*, 2017.
- 382 [11] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in  
383 neural networks. In *International Conference on Learning Representations*, 2018.
- 384 [12] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In S. Bengio,  
385 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural*  
386 *Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 387 [13] Lawrence Neal, Matthew Olson, Xiaoli Fern, Weng-Keen Wong, and Fuxin Li. Open set learning with  
388 counterfactual images. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors,  
389 *Computer Vision – ECCV 2018*, pages 620–635, Cham, 2018. Springer International Publishing.
- 390 [14] Jerzy Neyman and Egon Person. On the use and interpretation of certain test criteria for purpose of  
391 statistical inference. *Biometrika*, pages 175–240, 1928.
- 392 [15] T. Pietraszek. Optimizing abstaining classifiers using ROC analysis. In *Proceedings of the 22nd Interna-*  
393 *tional Conference on Machine Learning*, page 665–672, 2005.
- 394 [16] Yue Song, Nicu Sebe, and Wei Wang. Rankfeat: Rank-1 feature removal for out-of-distribution detection.  
395 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural*  
396 *Information Processing Systems*, volume 35, pages 17885–17898. Curran Associates, Inc., 2022.
- 397 [17] Yiyu Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. In  
398 A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information*  
399 *Processing Systems*, 2021.
- 400 [18] Yiyu Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations.  
401 In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in*  
402 *Neural Information Processing Systems*, volume 34, pages 144–157. Curran Associates, Inc., 2021.

- 403 [19] Yiyu Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest  
404 neighbors. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan  
405 Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
406 *Proceedings of Machine Learning Research*, pages 20827–20840. PMLR, 17–23 Jul 2022.
- 407 [20] Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-logit  
408 matching. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages  
409 4911–4920, 2022.
- 410 [21] Jingkan Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang,  
411 Guangyao Chen, Bo Li, Yiyu Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Dan Hendrycks, Yixuan  
412 Li, and Ziwei Liu. Openood: Benchmarking generalized out-of-distribution detection. In *Conference on  
413 Neural Information Processing Systems (NeurIPS 2022) Track on Datasets and Benchmar*, 2022.
- 414 [22] Jingkan Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A  
415 survey, 2022.