



# 1 Learning with mitigating random consistency 2 from the accuracy measure

3 Jieting Wang<sup>1</sup> · Yuhua Qian<sup>1</sup> · Feijiang Li<sup>1</sup>

4 Received: 13 April 2020 / Revised: 27 July 2020 / Accepted: 19 September 2020

5 © The Author(s) 2020

## 6 Abstract

7 Human beings may make random guesses in decision-making. Occasionally, their guesses  
8 may generate consistency with the real situation. This kind of consistency is termed ran-  
9 dom consistency. In the area of machine learning, the randomness is unavoidable and ubiq-  
10 uitous in learning algorithms. However, the accuracy (A), which is a fundamental perfor-  
11 mance measure for machine learning, does not recognize the random consistency. This  
12 causes that the classifiers learnt by A contain the random consistency. The random con-  
13 sistency may cause an unreliable evaluation and harm the generalization performance. To  
14 solve this problem, the pure accuracy (PA) is defined to eliminate the random consistency  
15 from the A. In this paper, we mainly study the necessity, learning consistency and leaning  
16 method of the PA. We show that the PA is insensitive to the class distribution of classifier  
17 and is more fair to the majority and the minority than A. Subsequently, some novel gener-  
18 alization bounds on the PA and A are given. Furthermore, we show that the PA is Bayes-  
19 risk consistent in finite and infinite hypothesis space. We design a plug-in rule that maxi-  
20 mizes the PA, and the experiments on twenty benchmark data sets demonstrate that the  
21 proposed method performs statistically better than the kernel logistic regression in terms of  
22 PA and comparable performance in terms of A. Compared with the other plug-in rules, the  
23 proposed method obtains much better performance.

24 **Keywords** Random consistency · Accuracy · Pure accuracy · Bayes-risk consistent

---

A1 Editors: Kee-Eung Kim, Vineeth N. Balasubramanian.

A2 ✉ Yuhua Qian  
A3 jinchengqyh@126.com

A4 Jieting Wang  
A5 jjietingwang@email.sxu.edu.cn

A6 Feijiang Li  
A7 feijiangli@email.sxu.edu.cn

A8 <sup>1</sup> Institute of Big Data Science and Industry, Shanxi University, Taiyuan 030006, Shanxi Province,  
A9 China

## 25 1 Introduction

26 In the process of decision-making, human beings may make random guesses with-  
27 out logical reasoning when they lack sufficient evidence or detailed knowledge. For  
28 instance, intern doctors are likely to diagnose patients with colds during flu season, and  
29 students are likely to choose a lucky option when faced with a difficult multiple-choices  
30 question. Sometimes, these random guesses may generate consistency with the real situ-  
31 ation. We term this consistency the random consistency.

32 In the area of machine learning, randomness is unavoidable and ubiquitous in con-  
33 structing classifiers, such as collecting and labeling data, selecting the structures or  
34 parameters of models and even in setting random operations (Ghahramani 2015). The  
35 prediction results of the learning models may also contain the random consistency. The  
36 random consistency produces dishonest feedback, misleads the decision direction and  
37 harms the improvement of the generalization ability, especially when the tendency of  
38 random guesses coincides with the class distribution of the real situation.

39 Eliminating the random consistency from evaluation measures has been well-studied  
40 in the field of educational psychology, where researchers advocate that the expected  
41 score for the accurate answer with no insight would be zero rather than one. This elimi-  
42 nation has proven helpful in achieving a higher reliability and validity assessment and  
43 increasing the performance of examinees (Sabers and Feldt 1968; Diamond and Evans  
44 1973; Wu et al. 2017; Budescu and Bar-Hillel 1993; Espinosa and Gardeazabal 2010).  
45 In the field of clustering evaluation, eliminating the random consistency has been an  
46 increasingly employed method to improve the quality of clustering evaluation (Hubert  
47 and Arabie 1985; Albatineh et al. 2006; Vinh et al. 2009, 2010; Albatineh and Niewia-  
48 domska-Bugaj 2011; Qian et al. 2016; Li et al. 2018, 2019).

49 In the area of classification, the accuracy ( $A$ ) is a vital performance measure in  
50 model evaluation and learning theory. The original learning theories focus on searching  
51 the generalization bounds for the error probability (one minus accuracy) (Valiant 1984;  
52 Bartlett and Mendelson 2003). The traditional algorithms, including logistic regression,  
53 support vector machine and Adaboost are designed to optimize convex surrogate loss  
54 functions of the error probability (Zhang 2003; Bartlett et al. 2006). In ensemble learn-  
55 ing, accuracy has been used as the preferential measure to evaluate the performance of  
56 integration (Zhou et al. 2002; Martinezmunoz and Suarez 2006). Although it is a funda-  
57 mental performance measure, the accuracy does not recognize the random consistency,  
58 which may limit the performance of the algorithms based on it. In this paper, we aim to  
59 define a performance measure that eliminates the random consistency from the accuracy  
60 and to study the learning performance of the measure theoretically and experimentally.

### 61 1.1 Related work

62 The measure that eliminates the random consistency from the accuracy is referred to  
63 as the pure accuracy (PA). The PA measure is a kind of non-decomposable measures.  
64 The non-decomposable measures cannot be decomposed into each individual instance  
65 (Waegeman et al. 2014; Kotlowski and Dembczynski 2017; Sanyal et al. 2018). Simi-  
66 lar measures include the F-measure, AUC, and balanced error rate (Zhao et al. 2013).  
67 For the non-decomposable measures, many learning theories and algorithms have been  
68 developed.

69 From the aspect of learning theory, Waegeman et al. (2014) investigated the generali-  
70 zation bound in terms of the F-measure when optimizing the Hamming loss and subset  
71 zero-one loss in a multi-label learning setting, and concluded that optimizing such losses  
72 as a surrogate of the F-measure leads to a high worst-case regret. Bayes-risk consistency  
73 guarantees that by increasing the amount of data, a rule can eventually learn the optimal  
74 decision with high probability. Agarwal et al. (2005b) show the Bayes-risk consistency of  
75 the AUC based on a new proposed combinatorial parameter. The key step of their proof is  
76 the symmetrization by a ghost sample that is the same as that for the classification error  
77 rate (Devroye et al. 1996). In this paper, to clarify the surrogate relation of PA and A, we  
78 show the upper bound of PA value for A-optimal rule and the upper bound of A value for  
79 PA-optimal rule. In addition, we give a Bayes-risk consistency analysis for the pure accu-  
80 racy based on the Rademacher complexity in a finite hypothesis space and based on the VC  
81 dimension in an infinite space.

82 In optimizing the non-decomposable measures, Musicant et al. (2003) extended the sup-  
83 port vector machine to optimize the F-measure by setting appropriate parameters in the  
84 standard SVM. Joachims (2005) proposed a large margin machine for maximizing a convex  
85 lower bound of non-decomposable measures. Hazan et al. (2010) and Song et al. (2016)  
86 trained deep neural networks by inferring the gradients of the non-decomposable meas-  
87 ures. Narasimhan and Agarwal (2013) proposed a SVM model for optimizing the AUC  
88 via a tight convex upper bound. Waegeman et al. (2014) proposed an exact algorithm for  
89 optimizing the F-measure in the context of multi-label learning. Gao et al. (2016) proposed  
90 a one-pass AUC optimization algorithm that needed to read the training data only once.  
91 These methods directly optimize the non-decomposable measures. In addition to these  
92 direct methods, the plug-in rule is an effective method that learns a posterior probability  
93 function by the logistic regression method or some other mature methods, and searches a  
94 threshold that optimizes the objective measure. For optimizing non-decomposable meas-  
95 ures, Narasimhan et al. (2015) simply used the bisection method to determine a threshold.  
96 The bisection method require the monotonicity of the function being solved. Then, there  
97 is still much room for improving the effectiveness of the plug-in method. Here, we give an  
98 interval search method to determine the threshold of the plug-in rule for optimizing the PA.

## 99 1.2 Contributions

100 We aim to verify the learning ability and Bayes-risk consistency of the PA in this paper.  
101 First, with regard to the cost-sensitive loss function, we give a non-closed formulation of  
102 the optimal rule w.r.t the PA. Based on this formulation, we illustrate that the PA is insen-  
103 sitive to the class distribution of classifiers and gets a low bias in minority accuracy and  
104 majority accuracy compared with A. Second, we give a novel lower and an upper bound  
105 for the optimal rules w.r.t the A and PA, respectively. These bounds help us clarify the  
106 surrogate relation between the PA and A. Furthermore, the generalization upper bounds of  
107 the PA in the worst case are given to analyze the consistency. The proof of these bounds  
108 employ the same symmetrization technique that was applied to prove the generalization  
109 upper bound of the accuracy (Devroye et al. 1996) and AUC (Agarwal et al. 2005a). How-  
110 ever, the difference is that the PA has fractional formulation. Thus, the consistent analysis  
111 of the PA needs to handle the fractional formulation. Last, we design a plug-in rule in  
112 terms of maximizing the PA and experimentally validate its performance.

113 Briefly, the major contributions of this paper are summarized as follows:

- 114 • Some bounds for the optimal rules w.r.t the PA and A are given. These bounds theoret-  
115 ically show that the PA-optimal rule is capable of approaching a satisfactory A value for all  
116 distributions.
- 117 • Second, we develop an inequality to handle the probability of large deviations of variables  
118 in fractional form. The generalization bounds for the PA are shown in finite and infinite  
119 hypothesis space. These bounds verify the Bayes-risk consistency of learning by PA.
- 120 • We propose a plug-in rule based on the interval search method for optimizing the PA.  
121 Through it, we experimentally verify the fairness and performance of PA in learning.

122 The organization of this paper is presented as follows: We give the definition of the PA in  
123 Sect. 2. In Sect. 3, two examples are given to show the necessity of evaluating classifiers by  
124 the PA. In Sect. 4, a surrogate analysis between the PA and the A is conducted. In Sect. 5, the  
125 generalization upper bounds of the PA are developed. We propose a plug-in rule for optimiz-  
126 ing the PA and experimentally validate its performance in Sect. 6. We form a conclusion and  
127 propose future work in Sect. 7.

128 In this paper, definitions and theorems which are tagged with a literature reference are  
129 taken from the literature, while the original ones come without such a tag. All the proofs are  
130 presented in the “Appendix”.

## 131 2 Preliminaries

132 We consider the task of binary classification. Let  $\mathcal{X} \subset \mathcal{R}^d$  and  $\mathcal{Y} = \{+1, -1\}$  be the fea-  
133 ture space and the label space, respectively. The underlying distribution of  $\mathcal{X} \times \mathcal{Y}$  is usually  
134 unknown, and we only have a collection of empirical data  $\mathcal{S}_N = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  that  
135 are drawn independently from this distribution. The goal of classification is to learn a classi-  
136 fier  $h(\mathbf{x})$  mapping from  $\mathcal{X}$  to  $\mathcal{Y}$  via  $\mathcal{S}_N$ . Let  $\mathcal{H}$  be the hypothesis space, from which the clas-  
137 sifier  $h(\mathbf{x})$  is learnt. To evaluate the performance of classifiers, the confusion matrix is usu-  
138 ally employed. Let  $TP, FP, FN, TN$  denote the true positive  $\mathbb{P}(h(X) = +1, Y = +1)$ , false  
139 positive  $\mathbb{P}(h(X) = +1, Y = -1)$ , false negative  $\mathbb{P}(h(X) = -1, Y = +1)$  and true negative  
140  $\mathbb{P}(h(X) = -1, Y = -1)$ , respectively. Let  $p$  and  $q(h)$  denote the probability of  $\mathbb{P}(Y = +1)$  and  
141  $\mathbb{P}(h(X) = +1)$ , respectively. The confusion matrix is shown in Table 1.

142 Based on the confusion matrix, the accuracy (A) and the error probability (L) are defined  
143 as:

$$144 \quad A(h) = \mathbb{P}(h(X) = Y) = TP + TN, \tag{1}$$

$$146 \quad L(h) = \mathbb{P}(h(X) \neq Y) = FP + FN. \tag{2}$$

**Table 1** Confusion matrix

h(X)	Y		Total (h)
	Y = +1	Y = -1	
$h(X) = +1$	$TP$	$FP$	$q(h)$
$h(X) = -1$	$FN$	$TN$	$1 - q(h)$
Total (Y)	$p$	$1 - p$	1

148 **2.1 The definition of PA**

149 To define the pure accuracy (PA), we begin with giving the definition of random accuracy  
 150 (RA), which aims to measure the random consistency in accuracy. For the classifier  $h(\mathbf{x})$  to  
 151 be evaluated, let  $\mathcal{H}^{q(h)}$  be the set of all possible binary partitions with the same class distri-  
 152 bution as it:

$$153 \quad \mathcal{H}^{q(h)} = \{h' : \mathbb{P}(h'(X) = +1) = q(h), h'(X) \in \{+1, -1\}\}. \quad (3)$$

154  
 155 Considering that the output preference of the classifier (tendency of predicting which  
 156 instances as positive) is unknown in advance, we suppose the partitions in  $\mathcal{H}^{q(h)}$  are uni-  
 157 formly distributed. Because the partitions in  $\mathcal{H}^{q(h)}$  have the same output randomness as the  
 158 classifier to be evaluated, we define RA as the expectation accuracy over the partitions in  
 159  $\mathcal{H}^{q(h)}$ .

160 **Lemma 1** *When the partitions in  $\mathcal{H}^{q(h)}$  are distributed uniformly, the expectation accu-  
 161 racy of partitions in  $\mathcal{H}^{q(h)}$  is:*

$$162 \quad \mathbb{E}_{h' \in \mathcal{H}^{q(h)}} A(h') = pq(h) + (1 - p)(1 - q(h)). \quad (4)$$

163  
 164 **Definition 1** The RA is defined as:

$$165 \quad RA(h) = pq(h) + (1 - p)(1 - q(h)). \quad (5)$$

166  
 167 **Definition 2** The PA is defined as:

$$168 \quad PA(h) = \frac{A(h) - RA(h)}{1 - RA(h)}. \quad (6)$$

169  
 170 **Definition 3** The pure loss (PL) is defined as:

$$171 \quad PL(h) = 1 - PA(h) = \frac{1 - A(h)}{1 - RA(h)}. \quad (7)$$

172  
 173 The denominator of PA guarantees the maximum value to be 1. Note that the formula-  
 174 tion of the PA coincides with the definition of Cohen's  $\kappa$  statistic (Cohen 1960; Scott 1955;  
 175 Goodman and Kruskal 1963). The difference between them is how to define the random  
 176 consistency. In the definition of Cohen's  $\kappa$  statistic, random consistency is called as chance  
 177 agreement. The chance agreement is the agreement degree that the two raters give their rat-  
 178 ings independently. The chance agreement between the classifier  $h(X)$  and the label  $Y$   
 179 is:

$$180 \quad \mathbb{P}(h(X) = Y) = \sum_{l \in \{-1, +1\}} \mathbb{P}(h(X) = Y = l) = pq(h) + (1 - p)(1 - q(h)). \quad (8)$$

181  
 182 The way we define the RA gives a general framework to measure the random consist-  
 183 ency in measures and is helpful to propose new performance measures.

184 Cohen's  $\kappa$  statistic has been successfully used in the area of psychology (Cameron  
 185 et al. 2003) and medicine (Blair and Stanley 2008). The advantage of correction for  
 186 the expected agreement by chance has made Cohen's  $\kappa$  statistic commonly be used as  
 187 a reliable performance measure in the area of machine learning (Ferri et al. 2009;). In  
 188 ensemble learning, Kappa-error diagrams have been used to gain insights about the

189 effectiveness of classifier ensembles (Kuncheva 2013) and to prune classifiers (Margin-  
 190 eantu and Dietterich 1997). In addition, Cohen’s  $\kappa$  statistic has been used for feature  
 191 selection (Vieira et al. 2010).

### 192 3 On the advantages of pure accuracy measure

193 A learning algorithm sensitive to the class distribution may get a decision boundary that  
 194 deviates from the optimal one. Thus, the learning objective should be insensitive to the  
 195 output class distribution. The extensively applied accuracy does not satisfy this prop-  
 196 erty. We employ Example 1 to show that the PA is satisfactory in this respect.

197 **Example 1** (Class distribution insensitivity) In this example, we aim to compare the eval-  
 198 uation result of the A and PA on the prediction results with different class distribution.  
 199 Under the settings of  $N = 100$  and  $p = 0.3$ , we randomly generate a binary vector as the  
 200 true label vector. A partition with a fixed class distribution  $q$  can be generated by Algo-  
 201 rithm 1. The class distribution  $q$  is varied from 0 to 1 with a step of 0.05. Under each  $q$ , we  
 202 run Algorithm 1 1000 times to generate 1000 partitions and use A and PA to evaluate the  
 203 partitions, respectively. The distributions of the A value and PA value are shown in Fig. 1.  
 204 From Fig. 1, it is easy to observe that the value of A decreases with the increase of  $q$ , while  
 205 the value of PA is always near zero. This finding reflects that the A is sensitive to the class  
 206 distribution of classifiers, while the PA is not.

---

#### 207 Algorithm 1 Generator of Partition with a Fixed Class Distribution

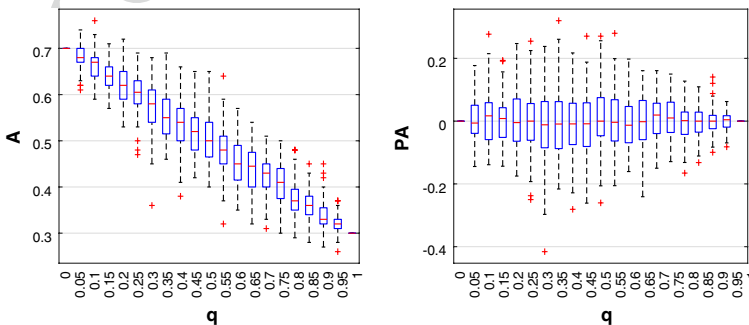
---

**Require:** Data set  $S_N = \{(x_i, y_i), i = 1, 2, \dots, N\}$ , class ratio  $q \in [0, 1]$ .

1: **for** each  $i \in N$  **do**  
 2:     Generating  $q_0 \sim \text{uniform}(0, 1)$ .  
 3:     **if**  $q_0 < q$  **then**  $h(x_i) = +1$ ;  
 4:     **else**  $h(x_i) = -1$   
 5:     **end if**  
 6: **end for**

**Ensure:** The predicted label  $h(x_i), i = 1, 2, \dots, N$ .

---



**Fig. 1** Distribution of A and PA under different  $q$ . Under each  $q$ , the box plot depicts the A values (left panel) and PA values (right panel) of Algorithm 1

209 **Lemma 2** (Devroye et al. 1996) Let  $\eta(\mathbf{x}) = \mathbb{P}(Y = +1|X = \mathbf{x})$  be the conditional class  
210 probability given  $X = \mathbf{x}$ . The classifier that maximizes the A or minimizes the L is:

$$211 \quad h_A^*(\mathbf{x}) = \arg \max_h A(h) = \begin{cases} +1, & \eta(\mathbf{x}) > \frac{1}{2}, \\ -1, & \text{otherwise.} \end{cases} \quad (9)$$

212

213 Correspondingly, the minimal error probability is

$$214 \quad L^* = L(h_A^*) = \mathbb{E}_X \min\{\eta(X), 1 - \eta(X)\}. \quad (10)$$

215

216 **Theorem 1** The classifier that maximizes the PA is

$$217 \quad h_{PA}^*(\mathbf{x}) = \arg \max_h PA(h) \quad (11)$$

218

$$219 \quad = \begin{cases} +1, & \eta(\mathbf{x}) > (\frac{1}{2} - p)PA^* + p, \\ -1, & \text{otherwise.} \end{cases} \quad (12)$$

220

221 where  $PA^* = PA(h_{PA}^*)$  and  $p = \mathbb{P}(Y = +1)$ .

222 For the cost-sensitive loss  $L_\rho = \rho FP + (1 - \rho)FN$ , it is known that when  $\rho$  is smaller,  
223 more attention will be paid to the minority class to get a smaller  $L_\rho$ . According to the proof  
224 of Theorem 1, PA is equivalent to  $L_\rho$  with  $\rho = (1/2 - p)PA^* + p$ . Due to  $PA^* \leq 1$ , a smaller  
225  $p$  value will generate a smaller  $(1/2 - p)PA^* + p$  value. In this case,  $h_{PA}^*$  will pay more  
226 attention to the minority class. Thus,  $h_{PA}^*$  may be insensitive to class distribution.

227 In learning classifiers, the minority class is often overwhelmed by the majority class to  
228 guarantee a higher overall accuracy (He and Garcia 2009). Then the classifiers learnt by  
229 optimizing the accuracy or error probability are usually biased to the majority class. This  
230 phenomenon is particularly desirable to avoid because the minority class is precious and  
231 inadequately represented. We employ Example 2 to show that the pure accuracy can miti-  
232 gate the classification bias.

233 **Example 2** (Fairness) To measure the bias of the classifier  $h(X)$ , we use the absolute differ-  
234 ence of the two class accuracy:

$$235 \quad Bias(h) = |\mathbb{P}(h(X) = -1|Y = -1) - \mathbb{P}(h(X) = +1|Y = +1)| \quad (13)$$

236

237 Assume that two class data are generated from Gaussian distribution:  $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma)$  and  
238  $\mathcal{N}(\boldsymbol{\mu}_2, \Sigma)$ . The label of the minority class is corrupted by the instance-independent noise at  
239 the level  $s_1$ :  $\mathbb{P}(\tilde{Y} = -1|Y = +1) = s_1$ .

240 For this learning task, the bias of  $h_A^*$  is:

$$241 \quad Bias(h_A^*) = \left| \Phi\left(\frac{d_0 + \Delta/2}{\sqrt{\Delta}}\right) - 1 + \Phi\left(\frac{d_0 - \Delta/2}{\sqrt{\Delta}}\right) \right|, \quad (14)$$

242

243 where  $\Phi(\bullet)$  is the cumulative distribution function of the standard normal distribution,  
244  $\Delta = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$  and  $d_0 = \ln \frac{1-p}{p} \frac{1}{1-2s_1}$ . Due to the formulation of  $h_{PA}^*$  is non-  
245 closed, the bias of it is simulate through a large number of instances. First, a sample that  
246 obey the distribution of this task are generated with a size of  $10^4$ . Then, the threshold that  
247 optimizes the PA is searched from the range  $[0, 1]$  with a step  $10^{-4}$ , and the bias of  $h_{PA}^*$  is  
248 calculated through the sample.

249 Let  $\mu_1 = -1$ ,  $\Sigma = 1$  and  $\mu_2$  vary from 0 to 2,  $p$  vary from 0.05 to 0.35 and the one-side  
 250 noise level  $s_1$  vary from 0 to 0.5. The bias curve of  $h_A^*$  (the dashed line) and that of  $h_{PA}^*$  (the  
 251 solid line) are shown in Fig. 2. As Fig. 2 shown, the dashed line is consistently lower than  
 252 the solid line in each case, which demonstrates that learning by PA is more fair than learn-  
 253 ing by A under different imbalance degree, overlap degree and noise level.

#### 254 4 Surrogate analysis of the optimal rules

255 The task of classification is to predict the labels of future observations. The optimal clas-  
 256 sifier is usually obtained by minimizing a loss function. From the same hypothesis space,  
 257 different loss functions usually obtain different optimal classifiers. In this section, we focus  
 258 on giving some novel bounds for  $h_{PA}^*(x)$  and  $h_A^*(x)$  to clarify the substitution relationship  
 259 between them in learning classifiers. Theorem 2 (derived by Lemma 3) and Theorem 3  
 260 (derived by Lemma 4) are major results of this section.

261 **Lemma 3** For all distributions, the plug-in rule with  $\rho$  as the decision threshold

$$262 \quad h_\rho(x) = \begin{cases} +1, & \eta(x) > \rho, \\ -1, & \text{otherwise,} \end{cases} \quad \text{where } \rho \in (0, \frac{1}{2}), \quad (15)$$

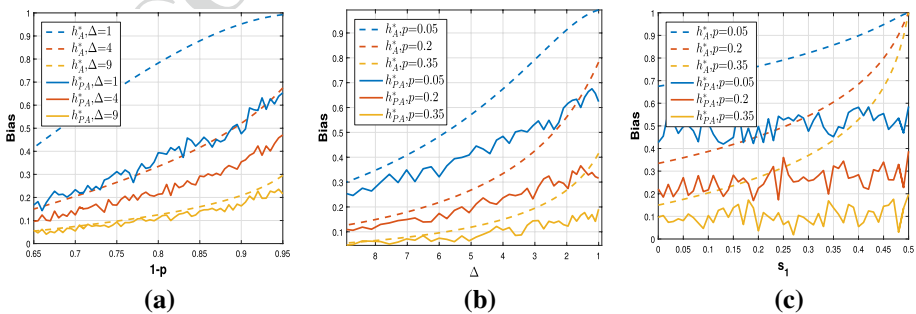
263 satisfies:

$$265 \quad L(h_\rho) \leq \frac{1-\rho}{\rho} L^*, \quad (16)$$

266 when  $\rho = 1/2$ , the equality holds.

268 Lemma 3 gives an upper bound on the error probability of the plug-in rule. Accord-  
 269 ing to Lemma 3, we have:

270 **Theorem 2** For all distributions, suppose that  $p = \mathbb{P}(Y = +1) \leq \frac{1}{2}$ , the error probability  
 271 of  $h_{PA}^*$  satisfies:



**Fig. 2** Bias of  $h_A^*$  and  $h_{PA}^*$  as a function of the ratio of majority class (left panel), the Mahalanobis distance of the two distributions (middle panel) and the one-side noise level (right panel). The dashed line is the bias curve of  $h_A^*$ , and the solid line is that of  $h_{PA}^*$



272 
$$L^* \leq L(h_{PA}^*) \leq \left( \frac{1}{(\frac{1}{2} - p)PA^* + p} - 1 \right) L^*. \tag{17}$$

273  
 274 From Theorem 2, we can conclude that the error probability of the optimal classifier  
 275 learnt by PA satisfies  $L(h_{PA}^*) \rightarrow L(h_A^*)$  as  $PA^* \rightarrow 1$  for all distributions.

276 **Lemma 4** For all distributions, suppose that  $p = \mathbb{P}(Y = +1) \leq \frac{1}{2}$ , the pure loss of  $h_A^*$   
 277 satisfies:

278 
$$PL(h_A^*) \leq \frac{L^*}{p\left(\frac{3}{2} - p\right) - L^*\left(\frac{1}{2} - p\right)}. \tag{18}$$

279  
 280 Lemma 4 gives the upper bound of the pure loss of  $h_A^*$  with respect to  $L^*$ . To obtain  
 281 the convergence relation between  $PL(h_A^*)$  with  $PL(h_{PA}^*)$ , we further amplifying  $L^*$  in  
 282 Theorem 3.

283 **Theorem 3** For all distributions, suppose  $p \leq \frac{1}{2}$ , the pure loss of  $h_A^*$  satisfies:

284 
$$PL(h_{PA}^*) \leq PL(h_A^*) \tag{19}$$

285  
 286 
$$\leq \frac{2(1 - p)}{p(3 - 2p) - L^*(1 - 2p)} PL(h_{PA}^*). \tag{20}$$

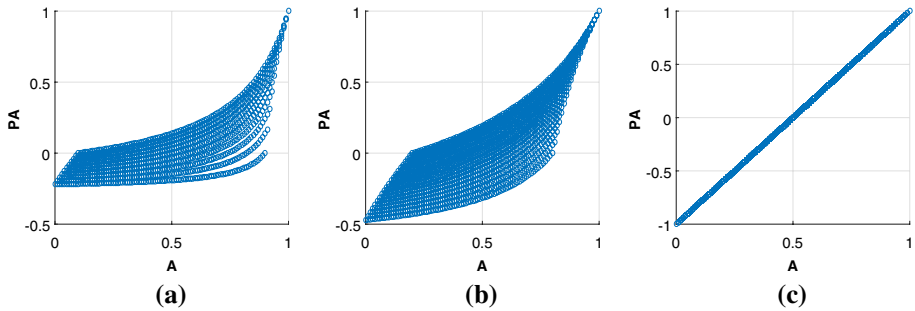
287  
 288 From Theorem 3, we can conclude that the pure loss of the optimal classifier learnt by  
 289 A satisfies  $PL(h_A^*) \rightarrow PL(h_{PA}^*)$  as  $L^* \rightarrow 0$  only when  $p = \frac{1}{2}$ . Based on Theorems 2 and 3, we  
 290 can infer that learning by PA can obtain a satisfactory A for all distributions, while learn-  
 291 ing by A can obtain a satisfactory PA only when the class distribution is balanced. We also  
 292 employ Example 3 to reflect this phenomenon.

293 **Example 3** (Surrogate analysis) In this example, we aim to analyse the surrogate relation  
 294 of A and PA. Under the settings of  $N = 100$  and  $p = \{0.1, 0.2, 0.5\}$ , we enumerate all pos-  
 295 sible values of FP and FN and calculate the A values and PA values. The A value and PA  
 296 value of each pair of (FP, FN) under different  $p$  are shown in Fig. 3. From Fig. 3, we can  
 297 observe that under the settings  $p = \{0.1, 0.2\}$ , when the PA value tends to 1, most of the A  
 298 values tends to 1, while when the A value tends to 1, most of the PA values are low. When  
 299  $p = 0.5$ , the relation between A and PA is linear.

300 **5 Bayes-risk consistency analysis of learning by the pure accuracy**  
 301 **measure**

302 The underlying distribution of  $\mathcal{X} \times \mathcal{Y}$  is usually unknown, and we only have a collection of  
 303 the empirical data  $\mathcal{S}_N = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  that is drawn independently from the distri-  
 304 bution. In machine learning, the classifier is generally obtained by the principle of empiri-  
 305 cal risk minimization (ERM). The feasibility of the ERM is guaranteed by the property of  
 306 Bayes-risk consistency. The corresponding loss function of PA is PL. Therefore, in this  
 307 section, we validate the learnability of PA by analyzing the Bayes-risk consistency of PL.

Author Proof



**Fig. 3** Surrogate analysis of A and PA when  $p = 0.1$  (left panel),  $p = 0.2$  (middle panel) and  $p = 0.5$  (right panel)

308 For the risk function  $R$ , let  $\hat{R}_N(h)$  be the empirical risk calculated on  $\mathcal{S}_N$ :  
 309  $\hat{R}_N(h) = \mathbb{E}_{X \times Y \in \mathcal{S}_N} R(h(X), Y)$ . ERM obtains the optimal rule  $h_{\hat{R}_N}^*$  from a hypothesis space  $\mathcal{H}$   
 310 by minimizing  $\hat{R}_N(h)$ :

$$311 \quad h_{\hat{R}_N}^* = \arg \min_{h \in \mathcal{H}} \hat{R}_N(h). \quad (21)$$

312  
 313 To guarantee the feasibility of the ERM, the property of Bayes-risk consistency is  
 314 defined as:

315 **Definition 4** (Devroye et al. 1996) The rule  $h_{\hat{R}_N}^*$  is Bayes-risk consistent, if for any small  
 316 enough  $\varepsilon$ , it satisfies

$$317 \quad \lim_{N \rightarrow \infty} \mathbb{P}(|R(h_{\hat{R}_N}^*) - \inf_h R(h)| > \varepsilon) = 0. \quad (22)$$

318  
 319 The Bayes-risk consistency requires that the empirical optimal hypothesis  $h_{\hat{R}_N}^*$  has a  
 320 large probability of converging to the universal optimal hypothesis as the number of empiri-  
 321 cal data tends to infinite.

322 To analysis the Bayes-risk consistency, the gap between  $R(h_{\hat{R}_N}^*)$  and  $\inf_h R(h)$  is usually  
 323 upper bounded by (Devroye et al. 1996):

$$324 \quad R(h_{\hat{R}_N}^*) - \inf_h R(h) \leq 2 \sup_{h \in \mathcal{H}} |\hat{R}_N(h) - R(h)|, \quad (23)$$

325  
 326 which is known as the estimation error. The estimation error measures the performance gap  
 327 between the empirical data and the underlying distribution. The convergence of the estima-  
 328 tion error ensures that the rule learnt finite samples can be generalized to infinite samples.  
 329 The bound of the estimation error, the so-called generalization bound, is the key factor in  
 330 studying the property of the Bayes-risk consistency.

331 The Rademacher complexity (Bartlett and Mendelson 2003) and the VC Dimension (Vap-  
 332 nik and Chervonenkis 1971) are two complexity measures of the hypothesis space; they have  
 333 a crucial role in bounding the estimation error in the sense of accuracy. Here, we use the gen-  
 334 eralization bounds based on them to analyse the Bayes-risk consistency of learning by PA. To  
 335 save space, we omit the definitions of the Rademacher complexity, the VC dimension and the  
 336 corresponding generalization bounds.

### 337 5.1 The Bayes-risk consistency of the pure loss measure in a finite hypothesis space

338 The fractional form of the pure loss leads to that the empirical value of it is not an unbiased  
 339 estimation of the expected value. Therefore, the techniques in deriving the generalization  
 340 bounds of the error probability (Theorem 8 in Bartlett and Mendelson (2003) and Theorem 2  
 341 in Vapnik and Chervonenkis (1971)) cannot be directly applied. Here, we establish a bridge  
 342 between the estimation error of the pure loss and that of the error probability; and then obtain  
 343 the Bayes-risk consistency of the pure loss in finite hypothesis space and infinite hypothesis  
 344 space based on Theorem 8 in Bartlett and Mendelson (2003) and Theorem 2 in Vapnik and  
 345 Chervonenkis (1971), respectively.

346 First, we give the formulation of the empirical error probability  $\widehat{L}_N(h)$  and the empirical  
 347 random accuracy  $\widehat{RA}_N(h)$  to analysis the Bayes-risk consistency:

$$348 \quad \widehat{L}_N(h) = \sum_{i=1}^N \mathbf{I}\{h(x_i) \neq y_i\}, \quad (24)$$

349

$$350 \quad \widehat{RA}_N(h) = \frac{1}{N|\mathcal{H}^{q(h)}|} \sum_{j=1}^{|\mathcal{H}^{q(h)}|} \sum_{i=1}^N \mathbf{I}\{h_j(x_i) = y_i\}, \quad (25)$$

351

352 where  $h_j \in \mathcal{H}^{q(h)}$ ,  $|\cdot|$  is the cardinality of a set and  $\mathbf{I}\{\cdot\}$  is the indicator function. Then, the  
 353 empirical pure loss  $\widehat{PL}_N(h)$  is

$$354 \quad \widehat{PL}_N(h) = \frac{\widehat{L}_N(h)}{1 - \widehat{RA}_N(h)}. \quad (26)$$

355

356 In practice, according to Lemma 1, the empirical random accuracy is computed by:

$$357 \quad \widehat{RA}_N(h) = 1 - \widehat{p}_N - (1 - 2\widehat{p}_N)\widehat{q}(h)_N, \quad (27)$$

358

359 where

$$360 \quad \widehat{p}_N = \sum_{i=1}^N \mathbf{I}\{y_i = +1\}, \quad (28)$$

361

$$362 \quad \widehat{q}(h)_N = \sum_{i=1}^N \mathbf{I}\{h(x_i) = +1\}. \quad (29)$$

363

364 **Lemma 5** For two random variables  $Z_1, Z_2 \in [0, 1]$ , any  $\varepsilon \in (0, 1]$ , let  
 365  $\alpha = \mathbb{E}Z_1\mathbb{E}Z_2 / (2\mathbb{E}Z_1 + \mathbb{E}Z_2)$ , we have

$$366 \quad \mathbb{P}\left(\left|\frac{Z_1}{Z_2} - \frac{\mathbb{E}Z_1}{\mathbb{E}Z_2}\right| > \varepsilon\right) \leq \mathbb{P}(|Z_1 - \mathbb{E}Z_1| > \alpha\varepsilon) + 3\mathbb{P}(|Z_2 - \mathbb{E}Z_2| > \alpha\varepsilon). \quad (30)$$

367

368 Lemma 5 links the probability of the large estimation error of the fractional variable to  
 369 that of the numerator and denominator. Based on Lemma 5, we obtain Theorems 4 and 5.

370 **Theorem 4** Suppose the cardinality of  $\mathcal{H}$  is finite:  $|\mathcal{H}| < \infty$ , then for every  $h \in \mathcal{H}$ , any  
371  $\varepsilon \in (0, 1]$ , we have

$$372 \quad \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \varepsilon \right\} \\ 373 \quad \leq 8|\mathcal{H}| \exp \left\{ -2N \left( \alpha\varepsilon - \frac{Rc(\mathcal{H})}{2} \right)^2 \right\}, \quad (31)$$

374 where  $\alpha = \min_{h \in \mathcal{H}} \frac{L(h)}{2+PL(h)}$  and  $Rc(\mathcal{H})$  is the Rademacher complexity of  $\mathcal{H}$ .

375 Theorem 4 provides the probability of the large estimation error in terms of the number  
376 of the empirical data in finite hypothesis space. From Theorem 4, we can conclude that  
377 learning by the PA is Bayes-risk consistency in a finite hypothesis space.

## 378 5.2 The Bayes-risk consistency of the pure loss measure in an infinite hypothesis 379 space

380 In this section, we consider the Bayes-risk consistency in an infinite hypothesis space. For  
381 an infinite hypothesis space, the union bound cannot be utilized. We utilize the symmetri-  
382 zation technical to bound the estimation error of the pure loss. Next, we divide the hypoth-  
383 esis space into  $N + 1$  subspaces according to the class probability of hypothesis functions,  
384 to ensure that each hypothesis subspace has the same degree of random accuracy. Then,  
385 we employ the VC bound of the error probability to bound the estimation error of the pure  
386 loss.

387 **Lemma 6** Let  $\mathcal{S}'_N = \{(\mathbf{x}'_{1,N}, y'_{1,N}), \dots, (\mathbf{x}'_{N,N}, y'_{N,N})\}$  be an independent and identically distributed  
388 collection as  $\mathcal{S}_N$  and  $PL'_N(h)$  is the corresponding empirical pure loss. Suppose  
389  $N \geq 5(6 + 4\alpha\varepsilon)\alpha^{-2}\varepsilon^{-2}$ , where  $\alpha = \min_{h \in \mathcal{H}} \frac{L(h)}{2PL(h)+1}$ ,  $\varepsilon \in (0, 1]$ , then we have

$$390 \quad \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \varepsilon \right\} \\ 391 \quad \leq 2\mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - \widehat{PL}'_N(h) \right| > \frac{\varepsilon}{2} \right\}. \quad (32)$$

392 **Theorem 5** As the same condition as Lemma 6 and suppose the VC dimension of  $\mathcal{H}$  is  
393 finite:  $d_{vc}(\mathcal{H}) < \infty$ , we have:

$$394 \quad \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \varepsilon \right\} \\ 395 \quad \leq 4(N + 1) \exp \left\{ - \left( \frac{\varepsilon^2(1 - |\widehat{2p}_N - 1|)^2}{16} - \frac{d_{vc}(\mathcal{H}) \ln(2eN/d_{vc}(\mathcal{H}))}{N} \right) N \right\}. \quad (33)$$

396 Theorem 5 provides the probability of the large estimation error in terms of the number  
397 of the empirical data in infinite hypothesis space. From Theorem 5, we can conclude that  
398 learning by the PA is Bayes-risk consistent in an infinite hypothesis space.

## 399 6 Performance validation of learning by the pure accuracy measure

400 By the Bayes-risk consistency, we have shown that the PA can be utilized to learn clas-  
 401 sifiers through minimizing PL. However, due to the fractional form, optimizing PL is a  
 402 challenging task. To handle this challenge, we introduce the plug-in rule and propose an  
 403 interval search method.

404 The plug-in rule refers to a rule with a formulation of  $h_{\delta^*}(\mathbf{x}) = \text{sign}(\hat{\eta}(\mathbf{x}) - \delta^*)$ , where  
 405  $\hat{\eta}(\mathbf{x})$  is an estimator of the posterior probability  $\eta(\mathbf{x}) = \mathbb{P}(Y = +1|X = \mathbf{x})$  and  $\delta^*$  is a thresh-  
 406 old (Koyejo et al. 2014). The plug-in method mainly contains the following steps: first,  
 407 randomly split the training data  $\mathcal{S}_N$  into  $\mathcal{S}_1$  and  $\mathcal{S}_2$ ; second, learn  $\hat{\eta}(\mathbf{x})$  by minimizing a loss  
 408 function on  $\mathcal{S}_1$ ; third, determine  $\delta^*$  by maximizing the learning objective on  $\mathcal{S}_2$ .

409 In Narasimhan et al. (2014), it has been proved that assigning an empirical threshold to  
 410 a suitable posterior probability estimate can optimize the performance measures expressed  
 411 as a function of the *TP* and *TN* and  $p$ . That is, the plug-in method can optimize a com-  
 412 plex performance measure through searching a decision threshold that optimizes the meas-  
 413 ure for the posterior probability estimate. The major focus of this section is developing an  
 414 method to search the threshold that optimizes PL rather than to learn the posterior prob-  
 415 ability  $\hat{\eta}(\mathbf{x})$ .

416 In this section, first, we introduce the method to learn the posterior probability. Then,  
 417 we discuss some methods of determining the threshold that optimizes PL and propose a  
 418 interval search method. Finally, we experimentally validate the performance of the interval  
 419 search method and the classifier learnt by the PA.

### 420 6.1 Learning $\hat{\eta}(\mathbf{x})$

421 Many methods can be employed to learn  $\hat{\eta}(\mathbf{x})$ . Here, we introduce the kernel logistic  
 422 regression model, which is proven to be a suitable posterior probability estimate (Ingo  
 423 2005; Narasimhan et al. 2014; Menon et al. 2013). The kernel logistic regression model is:

$$424 \max_{\alpha_j} \sum_{i=1}^{|\mathcal{S}_1|} \sum_{j=1}^{|\mathcal{S}_1|} \alpha_j y_j K(\mathbf{x}_j, \mathbf{x}_i) y_i - \sum_{i=1}^{|\mathcal{S}_1|} \log \left( 1 + \exp \left( \sum_{j=1}^{|\mathcal{S}_1|} \alpha_j y_j K(\mathbf{x}_j, \mathbf{x}_i) \right) \right) \quad (34)$$

425 where  $\alpha_i$  are the variables to be solved and  $K(\bullet, \bullet)$  is kernel function. With the optimal  $\alpha_i^*$  is  
 426 obtained by the gradient descent method, we have

$$428 \hat{\eta}(\mathbf{x}) = \frac{1}{1 + \exp(-\sum_{j=1}^{|\mathcal{S}_1|} \alpha_j^* y_j K(\mathbf{x}_j, \mathbf{x}))}. \quad (35)$$

### 430 6.2 The interval search method

431 As for determining  $\delta^*$ , different threshold settings correspond to optimizing different learning  
 432 objective functions.

433 To optimize the accuracy, the threshold  $\delta^*$  of the plug-in rule is 0.5, and this is the so-  
 434 called kernel logistic regression (KLR) method. To optimize the balanced accuracy (BA),  
 435 the threshold  $\delta^*$  of the plug-in rule is  $p$  (Menon et al. 2013).

436 For the measures in a fractional form, search strategies are effective and simple. An  
 437 intuitive approach to determine the optimal threshold is the point-wise search method,

438 namely, evaluating the fractional measure at each possible threshold and outputting the best  
 439 performing threshold. There is no doubt that exhausting all possible thresholds is impos-  
 440 sible. The grid search is a method to handle this, which divide the range of the threshold  
 441 into multiple equal intervals and set the end points as the candidate thresholds. Besides,  
 442 the posterior probabilities on  $\mathcal{S}_2$  can also be set to the candidate thresholds. We term this  
 443 search strategy the  $\mathcal{S}_2$ -search. The grid search method and the  $\mathcal{S}_2$ -search method search the  
 444 threshold in a limited range. In addition to the point-wise search methods, the bisection  
 445 method transforms the fractional measure to a one-dimensional function and obtains the  
 446 optimal threshold by solving the zero root of the one-dimensional function in binary (Nar-  
 447 asimhan et al. 2015). The bisection requires that the objective function be monotone on  
 448 the interval, while the fractional performance measures are usually non-monotonic with  
 449 respect to the threshold.

450 In this subsection, we develop a method for searching the optimal threshold via the  
 451 interval search method, and use this method to minimize PL. The interval search method  
 452 is an effective way to search the local minimum of a unimodal function (Chong and Zak  
 453 2011). For a unimodal one-dimensional function  $f(r)$  defined in  $[\alpha, \beta]$ , to obtain the mini-  
 454 mum  $r^*$ , the interval search method is based on the idea that it produces a series of inter-  
 455 vals  $[\alpha_k, \beta_k]$ , where  $[\alpha_{k+1}, \beta_{k+1}] \subset [\alpha_k, \beta_k]$  and  $\lim_{k \rightarrow \infty} \beta_k = \lim_{k \rightarrow \infty} \alpha_k = r^*$ . Specifically, the  
 456 interval search method inserts two points in each iteration and produces  $[\alpha_k, \lambda_k, \mu_k, \beta_k]$ . If  
 457  $f(\lambda_k) < f(\mu_k)$ , then  $\alpha_{k+1} = \alpha_k$  and  $\beta_{k+1} = \mu_k$ ; otherwise,  $\alpha_{k+1} = \lambda_k$  and  $\beta_{k+1} = \beta_k$ . When the  
 458 interval length is reduced by the ratio of  $1 - (\sqrt{5} - 1)/2$ , the interval search method is so-  
 459 called gold section method.

460 For any plug-in rule  $h_\delta(\mathbf{x}) = \text{sign}(\hat{\eta}(\mathbf{x}) - \delta)$ , we briefly discuss about whether the PL is a  
 461 unimodal function of the threshold  $\delta$ . According to the proof of Theorem 1, the PL is con-  
 462 sistent to the cost-sensitive loss with the optimal threshold as the cost weight:

$$463 \quad L_{\delta^*}(\delta) = \delta^* FP(\delta) + (1 - \delta^*) FN(\delta), \quad (36)$$

464 where  $\delta^*$  is the minimum of  $L_{\delta^*}(\delta)$ :

$$465 \quad \delta^* = \underset{\delta}{\text{argmin}} L_{\delta^*}(\delta), \quad (37)$$

466 and  $FP(\delta) = \mathbb{P}(\eta(X) > \delta, Y = -1)$ ,  $FN(\delta) = \mathbb{P}(\eta(X) \leq \delta, Y = +1)$ . Because

$$469 \quad FP(\delta) = \mathbb{P}(\eta(X) > \delta, Y = -1) = \mathbb{P}(\eta(X) > \delta) - \mathbb{P}(\eta(X) > \delta, Y = +1), \quad (38)$$

470 we have:

$$472 \quad L_{\delta^*}(\delta) = \delta^* FP(\delta) + (1 - \delta^*) FN(\delta) \\ 473 \quad = \delta^* \mathbb{P}(\eta(X) > \delta) - \delta^* \mathbb{P}(Y = +1) + \mathbb{P}(\eta(X) \leq \delta, Y = +1). \quad (39)$$

474 For  $\delta_1 < \delta_2$ , we have:

$$475 \quad L_{\delta^*}(\delta_1) - L_{\delta^*}(\delta_2) \\ 476 \quad = \delta^* \mathbb{P}(\eta(X) \in (\delta_1, \delta_2]) - \mathbb{P}(\eta(X) \in (\delta_1, \delta_2], Y = +1). \quad (40)$$

477 Thus, if

$$478 \quad \frac{\mathbb{P}(\eta(X) \in (\delta_1, \delta_2], Y = +1)}{\mathbb{P}(\eta(X) \in (\delta_1, \delta_2])} < \delta^* \quad (41)$$

479 we have  $L_{\delta^*}(\delta_1) > L_{\delta^*}(\delta_2)$ ; otherwise,  $L_{\delta^*}(\delta_1) < L_{\delta^*}(\delta_2)$ .

481 The unimodality of PL requires that for  $\delta_1 < \delta_2 < \delta^*$ ,  $L_{\delta^*}(\delta_1) > L_{\delta^*}(\delta_2)$  and for  
 482  $\delta^* < \delta_1 < \delta_2$ ,  $L_{\delta^*}(\delta_1) < L_{\delta^*}(\delta_2)$ . Thus, when  $\delta^* < \delta_1 < \delta_2$ , the unimodality of  $L_{\delta^*}(\delta)$   
 483 requires that the posteriori probability should satisfy condition (41), which signifies that  
 484 there exist a small number of positive objects in the objects with small posterior prob-  
 485 abilities. When  $\delta_1 < \delta_2 < \delta^*$ , the unimodal of  $L_{\delta^*}(\delta)$  requires that the posteriori probability  
 486 should satisfy the contrary case of condition (41), which signifies that there exist a large  
 487 number of positive objects in the objects with large posterior probabilities.

488 According to the above discussion, if the posteriori probability is sufficiently good, PL  
 489 is a unimodal function of  $\delta$ . The interval search method is applied to obtain  $\delta^*$ . From Theo-  
 490 rem 1, we have

$$491 \quad \delta^* = \left(\frac{1}{2} - p\right)PA^* + p = \frac{1}{2} - \left(\frac{1}{2} - p\right)PL^*. \quad (42)$$

492 Then, we express the plug-in rule as:

$$493 \quad h_r(\mathbf{x}) = \text{sign}\left[\hat{\eta}(\mathbf{x}) - \left(\frac{1}{2} - \left(\frac{1}{2} - p\right)r\right)\right], \quad (43)$$

494 and apply the interval search method to finding the optimal  $r$  that minimizes  $\widehat{PL}_{|S_2|}(h_r(\mathbf{x}))$ .

495 A fixed reduction of the interval is employed. In each iteration, the interval length is  
 496 reduced by the  $\tau \in (0, 0.5)$  ratio. The interval search method is thus called as  $\tau$ -interval  
 497 search method and the ratio  $\tau$  is a parameter to be tuned. The interval search method for  
 498 minimizing the PL is shown as Algorithm 2. The time complexity of the  $\tau$ -interval search  
 499 method contains two parts, which are learning  $\hat{\eta}(\mathbf{x})$  and searching  $\delta^*$ . The time complex-  
 500 ity of learning  $\hat{\eta}(\mathbf{x})$  is the same as the gradient descent method, and the time complexity  
 501 of search  $\delta^*$  is  $\mathcal{O}(N \log_r \epsilon)$ , where  $N$  is the number of training data,  $\tau$  is the reduction ratio  
 502 of the interval and  $\epsilon$  is the threshold of the stop condition. Learning  $\hat{\eta}(\mathbf{x})$  is the main time  
 503 consuming part. When handling large number of samples, it is suggested to utilize effective  
 504 gradient descent method.

505

---

#### Algorithm 2 The $\tau$ -Interval Search Method for Minimizing the PL

---

**Require:** The training data  $S_N$

Randomly split the training data  $S_N$  into  $S_1$  and  $S_2$  with a ratio of 8 : 2 and use  $S_1$  to estimate  $\hat{\eta}(\mathbf{x})$

Set  $\alpha = 0, \beta = 1, t = 0, \epsilon = 0.0001$

Let  $\lambda = \alpha + \tau(\beta - \alpha)$  and  $\mu = \beta - \tau(\beta - \alpha)$ ,

Obtain  $h_\lambda(\mathbf{x}) = \text{sign}[\hat{\eta}(\mathbf{x}) - (\frac{1}{2} - (\frac{1}{2} - p)\lambda)]$ ,  $h_\mu(\mathbf{x}) = \text{sign}[\hat{\eta}(\mathbf{x}) - (\frac{1}{2} - (\frac{1}{2} - p)\mu)]$  and

calculate  $\widehat{PL}_{|S_2|}(h_\lambda, Y)$  and  $\widehat{PL}_{|S_2|}(h_\mu, Y)$  on  $S_2$

**while**  $\beta - \alpha > \epsilon$ , **do**

**IF**  $\widehat{PL}_{|S_2|}(h_\lambda(\mathbf{x})) \leq \widehat{PL}_{|S_2|}(h_\mu(\mathbf{x}))$ , **THEN** update  $\beta = \mu$  **ELSE** update  $\alpha = \lambda$ ;

$\lambda = \alpha + \tau(\beta - \alpha)$ ,  $\mu = \beta - \tau(\beta - \alpha)$  and calculate  $\widehat{PL}_{|S_2|}(h_\lambda(\mathbf{x}))$  and  $\widehat{PL}_{|S_2|}(h_\mu(\mathbf{x}))$  on

$S_2$ ;

$t = t + 1$ ;  $\delta^t = \frac{1}{2} - (\frac{1}{2} - p)\lambda$ .

**end while**

**Ensure:** The optimal threshold  $\delta^*$ .

---

### 508 6.3 Experiments

509 We validate the performance of the  $\tau$ -interval search on a variety of benchmark data sets.  
 510 By the benchmark data sets, we show that learning by PA is more fair in majority accuracy

511 and minority accuracy than A and compare the  $\tau$ -interval search method with some other  
512 plug-in rules to show its effectiveness.

513 The benchmark data sets are downloaded from the KEEL Data Set Repository  
514 (Alcalafdez et al. 2008) and the UCI Machine Learning Repository (Dua and Graff  
515 2017). These data sets are briefly described in Table 2, including data ID, name, size,  
516 number of attributes and the imbalance ratio(IR). The posterior probability is gen-  
517 erated by the kernel logistic regression and the kernel function is the RBF kernel  
518  $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$ .

519 Each data set is randomly divided into a training set, a validation set and a test set at  
520 a ratio of 3:1:1. The methods are compared in the same division. We randomly divide  
521 the data set 30 times to obtain an average performance. The parameter  $\gamma$  is chosen from  
522  $\{2^{-4}, 2^{-2}, 2^0, 2^2, 2^4, 2^6\}$  and the  $\tau$  is chosen from  $\{0.1, 0.2, 0.3, 0.4\}$  via the validation set.  
523 Each attribute is linearly scaled to the range  $[0, 1]$  using the maximum and minimum  
524 values in the training data. For each data, we also add 3% and 5% random uniform label  
525 noise to increase the complexity of the data.

526 First, to show that learning by PA is more fair than A, we compare the bias [refer  
527 to Eq. (13)] of KLR and the  $\tau$ -interval method. Figure 4 shows the comparison result.  
528 Each bar of Fig. 4 is the difference of the mean bias over 30 times between KLR and  
529 the  $\tau$ -interval method on each benchmark data set. As shown in Fig. 4, we observe that  
530 16/20, 17/20, 16/20 bars are greater than zero under 0% noise, 3% noise, 5% noise,  
531 respectively. That is, the bias of KLR is large than that of the  $\tau$ -interval method, which  
532 reflects the classifiers learnt by PA is more fair than the classifiers learnt by A.

**Table 2** Description of data sets

Data ID	Data name	Attribute	Instance	IR	Download
1	First-order theorem proving1	51	6118	1.02	UCI
2	First-order theorem proving2	51	6118	1.28	UCI
3	First-order theorem proving3	51	6118	1.16	UCI
4	First-order theorem proving4	51	6118	1.14	UCI
5	First-order theorem proving5	51	6118	1.27	UCI
6	Crx	15	653	1.21	KEEL
7	Heart	13	270	1.25	KEEL
8	Australian	14	690	1.25	KEEL
9	Wdbc	30	569	1.68	KEEL
10	Bands	19	365	1.70	KEEL
11	Ionosphere	33	351	1.79	KEEL
12	Wisconsin	9	683	1.86	KEEL
13	Pima	8	768	1.87	KEEL
14	Titanic	3	2201	2.10	KEEL
15	German	20	1000	2.33	KEEL
16	Segment	19	2308	6.02	KEEL
17	Dermatology	34	358	16.90	KEEL
18	Wilt	5	4839	17.54	UCI
19	Flare	11	1066	23.79	KEEL
20	Winequality-red	11	1599	29.17	KEEL



533 Second, to validate the performance the proposed method, the A and PA are  
 534 employed as evaluation measures. The benchmark methods are KLR, p-cut (with the  
 535 proportion of the minority class in  $S_2$  as the threshold), grid-search,  $S_2$ -search and  
 536 bisection method. The KLR aims to optimize the A, and p-cut aims to optimize the  
 537 balanced accuracy. The grid-search and  $S_2$ -search aim to optimize the PA. The bisection  
 538 is used to optimize the  $F_1$ -measure and PA, which are noted as Bisection- $F_1$  and  
 539 Bisection-PA, respectively. Tables 3, 5 and 7 show the mean and the standard deviation  
 540 of A over 30 time comparisons with 0%, 3% and 5% label noise, respectively. Tables 4,  
 541 6 and 8 show the mean and the standard deviation of PA over 30 time comparisons with  
 542 0%, 3% and 5% label noise, respectively. In each row of the tables, the method with the  
 543 maximal evaluation value is underlined and printed in bold type, and the method with a  
 544 dot indicates that the  $\tau$ -interval search is significantly better with regard to the pairwise  
 545 Student's  $t$  test with a level of 0.1. As shown in Tables 3, 4, 5, 6, 7 and 8, the evaluation  
 546 score obtained by the  $\tau$ -interval search is highlighted in bold and is underlined in most  
 547 of the comparisons. In many comparisons, the  $\tau$ -interval search is statistically better  
 548 than other methods.

549 To further analysis the statical performance of each method, for each method, we  
 550 calculate the gap between the times of the significant wins and the times of signifi-  
 551 cant loses. An algorithm  $a$  significantly wins  $b$  if its mean and standard deviation are  
 552 satisfied:

$$553 \quad \mu_a - 1.96 \frac{\sigma_a}{\sqrt{t}} > \mu_b + 1.96 \frac{\sigma_b}{\sqrt{t}}, \quad (44)$$

554  
 555 where  $t$  is the number of comparison times; otherwise,  $a$  significantly loses  $b$  (Please refer  
 556 to reference Li et al. (2016) for more details). Figure 5 shows the results of the statisti-  
 557 cal comparison. Each bar in Fig. 5 represents the gap between the times of the significant  
 558 wins and the times of significant loses. As shown in Fig. 5, we observe that the bar of the  
 559  $\tau$ -interval search method is the highest w.r.t PA under different noise level. With respect  
 560 to A, the bar of the  $\tau$ -interval search is the highest when the noise level is % and 5%; and  
 561 when the label is not polluted by noise, the bar of the  $\tau$ -interval search is the second high-  
 562 est. In general, we can conclude that the  $\tau$ -interval search method can optimize the PA  
 563 value better and also can obtain a satisfactory A value.

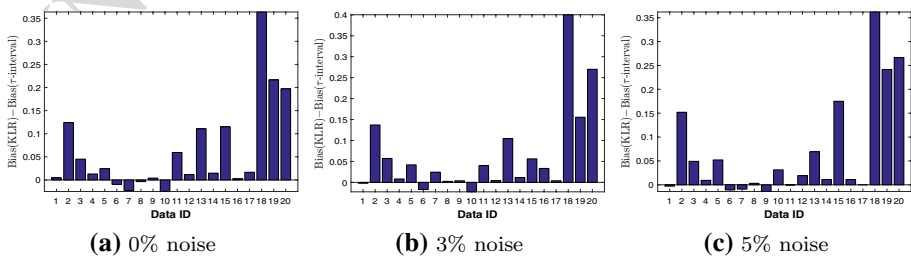


Fig. 4 Bias Gap of KLR and the  $\tau$ -interval method under different noise level. Each bar is the mean gap over 30 times on each data set

Table 3 Comparison of accuracy on data without noise

ID	KLR	p-cut	Bisection- $F_1$	Grid-search	$S_2$ -search	Bisection-PA	$\tau$ -interval
1	0.669 ± 0.012•	0.670 ± 0.012	0.647 ± 0.011•	0.672 ± 0.012	0.672 ± 0.012	0.670 ± 0.013•	<u>0.673</u> ± 0.012
2	0.671 ± 0.013	0.665 ± 0.014•	<u>0.672</u> ± 0.014	0.670 ± 0.014	0.671 ± 0.014	0.668 ± 0.013	0.670 ± 0.014
3	0.683 ± 0.013•	0.681 ± 0.012•	0.670 ± 0.013•	0.684 ± 0.014	0.684 ± 0.014•	0.683 ± 0.013•	<u>0.688</u> ± 0.013
4	0.663 ± 0.011•	0.662 ± 0.013•	0.654 ± 0.011•	0.661 ± 0.012•	0.661 ± 0.012•	0.662 ± 0.011•	<u>0.667</u> ± 0.011
5	0.683 ± 0.013	0.681 ± 0.014•	0.670 ± 0.012•	0.682 ± 0.016•	0.681 ± 0.016•	0.683 ± 0.014	<u>0.685</u> ± 0.012
6	0.859 ± 0.026•	0.864 ± 0.023	0.808 ± 0.030•	0.860 ± 0.021•	0.859 ± 0.022•	0.861 ± 0.025	<u>0.868</u> ± 0.021
7	0.789 ± 0.049•	0.786 ± 0.050•	0.792 ± 0.053•	0.770 ± 0.048•	0.776 ± 0.057•	0.788 ± 0.051•	<u>0.817</u> ± 0.038
8	0.850 ± 0.030	0.850 ± 0.032	0.821 ± 0.038•	0.848 ± 0.031•	0.847 ± 0.032•	0.850 ± 0.029	<u>0.857</u> ± 0.025
9	0.954 ± 0.022•	0.955 ± 0.022•	0.952 ± 0.022•	0.955 ± 0.023•	0.952 ± 0.023•	0.955 ± 0.022•	<u>0.965</u> ± 0.017
10	0.637 ± 0.059•	0.602 ± 0.069•	0.626 ± 0.060•	0.625 ± 0.063•	0.617 ± 0.074•	0.609 ± 0.067•	<u>0.654</u> ± 0.050
11	0.846 ± 0.042•	0.845 ± 0.042•	0.846 ± 0.042•	0.846 ± 0.046•	0.845 ± 0.043•	0.846 ± 0.041•	<u>0.867</u> ± 0.032
12	0.966 ± 0.013•	0.967 ± 0.013	0.946 ± 0.019•	0.970 ± 0.014	0.969 ± 0.014	0.966 ± 0.014	<u>0.971</u> ± 0.010
13	0.763 ± 0.028	0.750 ± 0.031•	0.758 ± 0.030•	0.755 ± 0.035•	0.756 ± 0.033	0.762 ± 0.029	<u>0.766</u> ± 0.029
14	0.779 ± 0.016•	0.767 ± 0.023•	0.777 ± 0.016•	0.779 ± 0.019	0.779 ± 0.019	0.780 ± 0.016	<u>0.783</u> ± 0.017
15	0.742 ± 0.022•	0.710 ± 0.027•	0.744 ± 0.024	0.725 ± 0.027•	0.726 ± 0.028•	0.732 ± 0.026•	<u>0.746</u> ± 0.023
16	0.993 ± 0.004•	0.992 ± 0.005•	0.992 ± 0.005•	0.994 ± 0.004•	0.995 ± 0.004•	0.993 ± 0.005•	<u>0.996</u> ± 0.003
17	0.999 ± 0.004	0.998 ± 0.005•	0.996 ± 0.006•	0.997 ± 0.006•	0.991 ± 0.012•	0.999 ± 0.004	<u>1.000</u> ± 0.003
18	<u>0.946</u> ± 0.000	0.716 ± 0.066•	0.878 ± 0.063	0.743 ± 0.220•	0.757 ± 0.207•	0.855 ± 0.035	0.843 ± 0.065
19	<u>0.958</u> ± 0.007	0.800 ± 0.054•	0.947 ± 0.013	0.934 ± 0.022•	0.935 ± 0.022•	0.927 ± 0.021•	0.942 ± 0.023
20	<u>0.968</u> ± 0.002	0.748 ± 0.063•	0.945 ± 0.030	0.932 ± 0.028•	0.935 ± 0.024•	0.924 ± 0.043•	0.940 ± 0.037
Rank	3.15	5.50	4.75	4.40	4.65	4.05	1.50

Table 4 Comparison of pure accuracy on data sets without noise

ID	KLR	p-cut	Bisection- $F_1$	Grid-search	$S_z$ -search	Bisection-PA	$\tau$ -interval
1	0.339 ± 0.025•	0.341 ± 0.025	0.293 ± 0.022•	0.345 ± 0.023	0.345 ± 0.024	0.340 ± 0.025•	<b>0.347</b> ± 0.024
2	0.318 ± 0.028•	0.325 ± 0.028	0.305 ± 0.028•	0.328 ± 0.029	<b>0.329</b> ± 0.028	0.326 ± 0.026	0.327 ± 0.029
3	0.361 ± 0.029•	0.363 ± 0.024•	0.320 ± 0.027•	0.366 ± 0.026	0.365 ± 0.027	0.364 ± 0.027•	<b>0.373</b> ± 0.026
4	0.324 ± 0.022•	0.325 ± 0.025•	0.295 ± 0.023•	0.325 ± 0.024•	0.325 ± 0.024•	0.324 ± 0.023•	<b>0.332</b> ± 0.023
5	0.356 ± 0.026•	0.365 ± 0.026	0.312 ± 0.024•	0.362 ± 0.029	0.361 ± 0.029•	0.365 ± 0.025	<b>0.369</b> ± 0.022
6	0.717 ± 0.053•	0.727 ± 0.045	0.605 ± 0.064•	0.719 ± 0.042•	0.718 ± 0.043•	0.722 ± 0.050	<b>0.736</b> ± 0.041
7	0.574 ± 0.099•	0.570 ± 0.101•	0.571 ± 0.111•	0.539 ± 0.097•	0.548 ± 0.116•	0.574 ± 0.104•	<b>0.629</b> ± 0.079
8	0.697 ± 0.060	0.698 ± 0.064	0.628 ± 0.081•	0.693 ± 0.063•	0.690 ± 0.065•	0.698 ± 0.059	<b>0.711</b> ± 0.051
9	0.902 ± 0.047•	0.903 ± 0.046•	0.896 ± 0.048•	0.902 ± 0.049•	0.896 ± 0.049•	0.903 ± 0.047•	<b>0.924</b> ± 0.037
10	0.196 ± 0.115•	0.181 ± 0.132•	0.188 ± 0.118•	0.186 ± 0.123•	0.179 ± 0.118•	0.185 ± 0.132•	<b>0.225</b> ± 0.108
11	0.647 ± 0.100•	0.644 ± 0.099•	0.644 ± 0.102•	0.649 ± 0.101•	0.650 ± 0.096•	0.646 ± 0.098•	<b>0.698</b> ± 0.073
12	0.924 ± 0.030•	0.927 ± 0.029	0.877 ± 0.044•	0.935 ± 0.031	0.931 ± 0.031	0.925 ± 0.031•	<b>0.935</b> ± 0.022
13	0.449 ± 0.071•	0.466 ± 0.069	0.415 ± 0.082•	0.459 ± 0.072	0.460 ± 0.072	0.470 ± 0.072	<b>0.474</b> ± 0.066
14	0.444 ± 0.040•	0.438 ± 0.042•	0.437 ± 0.042•	0.452 ± 0.044	0.451 ± 0.043	0.450 ± 0.042	<b>0.458</b> ± 0.042
15	0.345 ± 0.061•	0.372 ± 0.049	0.344 ± 0.070•	0.365 ± 0.054•	0.367 ± 0.054•	0.376 ± 0.059	<b>0.384</b> ± 0.067
16	0.973 ± 0.017•	0.968 ± 0.021•	0.969 ± 0.022•	0.975 ± 0.017•	0.978 ± 0.015•	0.973 ± 0.018•	<b>0.983</b> ± 0.014
17	0.986 ± 0.043•	0.983 ± 0.044•	0.960 ± 0.068•	0.972 ± 0.051•	0.885 ± 0.195•	0.986 ± 0.043•	<b>0.996</b> ± 0.022
18	0.000 ± 0.000•	0.063 ± 0.041	0.033 ± 0.050•	0.057 ± 0.037	0.057 ± 0.037	0.053 ± 0.047	<b>0.063</b> ± 0.037
19	0.147 ± 0.140•	0.163 ± 0.060•	0.208 ± 0.148•	0.241 ± 0.125	0.238 ± 0.133	0.231 ± 0.107	<b>0.250</b> ± 0.136
20	0.004 ± 0.026•	0.051 ± 0.038•	0.110 ± 0.112	0.109 ± 0.107	<b>0.112</b> ± 0.113	0.101 ± 0.097	0.110 ± 0.106
Rank	5.15	4.15	6.20	3.60	3.90	3.85	1.15

Table 5 Comparison of accuracy on data sets with 3% noise

ID	KLR	p-cut	Bisection- $F_1$	Grid-search	$S_2$ -search	Bisection-PA	$\tau$ -interval
1	0.660 ± 0.011•	0.659 ± 0.011•	0.642 ± 0.013•	0.659 ± 0.011•	0.660 ± 0.010•	0.659 ± 0.011•	<b>0.663</b> ± 0.012
2	<b>0.657</b> ± 0.012	0.650 ± 0.012•	0.654 ± 0.011	0.656 ± 0.014	0.655 ± 0.013	0.654 ± 0.012	0.656 ± 0.013
3	0.674 ± 0.015•	0.675 ± 0.012•	0.662 ± 0.016•	0.678 ± 0.013	0.678 ± 0.013	0.674 ± 0.013•	<b>0.679</b> ± 0.013
4	0.659 ± 0.012•	0.657 ± 0.013•	0.645 ± 0.012•	0.659 ± 0.012•	0.659 ± 0.012•	0.659 ± 0.013•	<b>0.663</b> ± 0.010
5	0.667 ± 0.015	0.663 ± 0.012•	0.658 ± 0.011•	0.665 ± 0.015	0.665 ± 0.014	0.666 ± 0.013	<b>0.669</b> ± 0.013
6	0.826 ± 0.029	0.827 ± 0.027	0.782 ± 0.035•	0.830 ± 0.026	0.827 ± 0.028•	0.827 ± 0.026	<b>0.837</b> ± 0.028
7	0.770 ± 0.055•	0.765 ± 0.049•	0.767 ± 0.045•	0.764 ± 0.055•	0.760 ± 0.055•	0.766 ± 0.050•	<b>0.801</b> ± 0.043
8	0.829 ± 0.033•	0.830 ± 0.033	0.800 ± 0.035•	0.831 ± 0.028•	0.829 ± 0.028•	0.829 ± 0.032•	<b>0.839</b> ± 0.026
9	0.904 ± 0.053•	0.898 ± 0.050•	0.893 ± 0.047•	0.900 ± 0.050•	0.904 ± 0.051•	0.903 ± 0.053•	<b>0.930</b> ± 0.035
10	0.594 ± 0.054•	0.572 ± 0.067•	0.589 ± 0.061•	0.596 ± 0.052•	0.600 ± 0.051•	0.577 ± 0.067•	<b>0.637</b> ± 0.046
11	0.811 ± 0.050•	0.809 ± 0.055•	0.816 ± 0.044•	0.810 ± 0.043•	0.808 ± 0.037•	0.808 ± 0.054•	<b>0.836</b> ± 0.044
12	0.935 ± 0.020	<b>0.939</b> ± 0.019	0.911 ± 0.019•	0.937 ± 0.020	0.939 ± 0.018	0.936 ± 0.020	0.939 ± 0.019
13	0.744 ± 0.030	0.721 ± 0.029•	0.741 ± 0.029	0.719 ± 0.042•	0.721 ± 0.043•	0.736 ± 0.031•	<b>0.746</b> ± 0.034
14	0.765 ± 0.015	0.744 ± 0.020•	0.763 ± 0.015•	0.762 ± 0.019	0.762 ± 0.019•	0.764 ± 0.015•	<b>0.768</b> ± 0.014
15	0.727 ± 0.026•	0.689 ± 0.036•	0.723 ± 0.027•	0.717 ± 0.028•	0.717 ± 0.026•	0.710 ± 0.026•	<b>0.740</b> ± 0.026
16	0.944 ± 0.019	0.884 ± 0.036•	0.929 ± 0.020•	0.941 ± 0.023•	0.942 ± 0.023•	0.942 ± 0.021•	<b>0.947</b> ± 0.020
17	0.914 ± 0.042•	0.914 ± 0.042•	0.914 ± 0.043•	0.917 ± 0.040	0.920 ± 0.037	0.914 ± 0.043•	<b>0.930</b> ± 0.030
18	<b>0.919</b> ± 0.000	0.658 ± 0.058•	0.833 ± 0.056	0.682 ± 0.192•	0.686 ± 0.197•	0.794 ± 0.051	0.788 ± 0.078
19	<b>0.929</b> ± 0.007	0.708 ± 0.052•	0.893 ± 0.032•	0.875 ± 0.046•	0.877 ± 0.047•	0.855 ± 0.062•	0.908 ± 0.039
20	<b>0.940</b> ± 0.001	0.673 ± 0.041•	0.886 ± 0.048	0.847 ± 0.136•	0.847 ± 0.139	0.821 ± 0.079•	0.893 ± 0.028
Rank	2.95	5.70	5.00	4.30	4.00	4.70	1.35

Table 6 Comparison of pure accuracy on data sets with 3% noise

ID	KLR	p-cut	Bisection- $F_1$	Grid-search	$S_2$ -search	Bisection-PA	$\tau$ -interval
1	0.320 ± 0.022•	0.319 ± 0.022•	0.282 ± 0.025•	0.318 ± 0.022•	0.320 ± 0.019•	0.319 ± 0.021•	<b>0.326</b> ± 0.024
2	0.288 ± 0.025•	0.296 ± 0.023	0.272 ± 0.023•	<b>0.301</b> ± 0.027	0.300 ± 0.027	0.299 ± 0.024	0.299 ± 0.028
3	0.343 ± 0.031•	0.350 ± 0.025	0.304 ± 0.033•	0.353 ± 0.025	0.353 ± 0.025	0.347 ± 0.026•	<b>0.355</b> ± 0.027
4	0.316 ± 0.024•	0.316 ± 0.026•	0.277 ± 0.025•	0.319 ± 0.023	0.319 ± 0.023	0.317 ± 0.026	<b>0.325</b> ± 0.019
5	0.325 ± 0.032•	0.329 ± 0.024	0.288 ± 0.024•	0.329 ± 0.029	0.328 ± 0.028	0.331 ± 0.025	<b>0.335</b> ± 0.026
6	0.652 ± 0.059•	0.654 ± 0.056	0.551 ± 0.074•	0.659 ± 0.051	0.654 ± 0.056•	0.653 ± 0.053	<b>0.673</b> ± 0.056
7	0.536 ± 0.112•	0.528 ± 0.100•	0.522 ± 0.095•	0.522 ± 0.112•	0.513 ± 0.114•	0.528 ± 0.102•	<b>0.597</b> ± 0.085
8	0.654 ± 0.066•	0.659 ± 0.066	0.584 ± 0.077•	0.659 ± 0.056•	0.656 ± 0.055•	0.657 ± 0.065	<b>0.676</b> ± 0.053
9	0.795 ± 0.115•	0.784 ± 0.108•	0.764 ± 0.104•	0.785 ± 0.110•	0.790 ± 0.115•	0.793 ± 0.114•	<b>0.851</b> ± 0.074
10	0.108 ± 0.119•	0.122 ± 0.143•	0.114 ± 0.130•	0.124 ± 0.091•	0.130 ± 0.091•	0.126 ± 0.142•	<b>0.193</b> ± 0.095
11	0.575 ± 0.108•	0.572 ± 0.117•	0.582 ± 0.098•	0.569 ± 0.095•	0.564 ± 0.083•	0.569 ± 0.115•	<b>0.632</b> ± 0.103
12	0.859 ± 0.042	<b>0.869</b> ± 0.040	0.800 ± 0.045•	0.864 ± 0.043	0.867 ± 0.039	0.861 ± 0.043	0.868 ± 0.041
13	0.410 ± 0.069•	0.404 ± 0.065•	0.387 ± 0.073•	0.394 ± 0.072•	0.397 ± 0.073•	0.415 ± 0.071•	<b>0.432</b> ± 0.074
14	0.418 ± 0.038	0.400 ± 0.040•	0.413 ± 0.037•	0.421 ± 0.037	0.420 ± 0.037	0.420 ± 0.035	<b>0.428</b> ± 0.034
15	0.320 ± 0.062•	0.335 ± 0.069•	0.303 ± 0.064•	0.340 ± 0.053•	0.339 ± 0.050•	0.345 ± 0.054•	<b>0.365</b> ± 0.057
16	0.782 ± 0.079	0.641 ± 0.079•	0.695 ± 0.100•	0.781 ± 0.079	0.783 ± 0.079	0.783 ± 0.076	<b>0.800</b> ± 0.072
17	0.498 ± 0.169	0.500 ± 0.169	0.495 ± 0.175•	0.508 ± 0.168	0.509 ± 0.159	0.502 ± 0.169	<b>0.558</b> ± 0.141
18	0.000 ± 0.000•	0.047 ± 0.041	0.025 ± 0.048•	<b>0.047</b> ± 0.033	0.046 ± 0.032	0.035 ± 0.043	0.044 ± 0.036
19	0.046 ± 0.087•	0.100 ± 0.051•	0.131 ± 0.096	0.128 ± 0.083	0.128 ± 0.081	<b>0.135</b> ± 0.076	0.133 ± 0.108
20	0.005 ± 0.022•	0.063 ± 0.043•	0.103 ± 0.088	0.103 ± 0.080	0.096 ± 0.078•	0.100 ± 0.069	<b>0.121</b> ± 0.077
Rank	5.10	4.50	6.10	3.45	3.80	3.70	1.35

Table 7 Comparison of accuracy on data sets with 5% noise

ID	KLR	p-cut	Bisection- $F_1$	Grid-search	$S_2$ -search	Bisection-PA	$\tau$ -interval
1	0.654 ± 0.012•	0.654 ± 0.012•	0.639 ± 0.014•	0.656 ± 0.012•	0.657 ± 0.012•	0.654 ± 0.012•	<b>0.661</b> ± 0.011
2	0.650 ± 0.013	0.646 ± 0.014•	0.650 ± 0.013	0.650 ± 0.014	<b>0.651</b> ± 0.014	0.648 ± 0.015	0.651 ± 0.014
3	0.668 ± 0.017•	0.666 ± 0.017•	0.659 ± 0.016•	0.670 ± 0.018•	0.669 ± 0.018•	0.668 ± 0.017•	<b>0.673</b> ± 0.016
4	0.641 ± 0.013•	0.642 ± 0.013•	0.632 ± 0.013•	0.643 ± 0.011	0.643 ± 0.011	0.643 ± 0.012	<b>0.646</b> ± 0.011
5	0.657 ± 0.011•	0.654 ± 0.014•	0.651 ± 0.013•	0.658 ± 0.012	0.658 ± 0.012	0.657 ± 0.013•	<b>0.661</b> ± 0.013
6	0.824 ± 0.028•	0.824 ± 0.026•	0.776 ± 0.036•	0.822 ± 0.027•	0.822 ± 0.027•	0.825 ± 0.028•	<b>0.835</b> ± 0.025
7	0.743 ± 0.051•	0.741 ± 0.058•	0.734 ± 0.055•	0.728 ± 0.059•	0.728 ± 0.057•	0.740 ± 0.055•	<b>0.770</b> ± 0.048
8	0.812 ± 0.029•	0.812 ± 0.027•	0.787 ± 0.033•	0.810 ± 0.032•	0.809 ± 0.035•	0.813 ± 0.030•	<b>0.822</b> ± 0.032
9	0.888 ± 0.034•	0.885 ± 0.031•	0.853 ± 0.029•	0.886 ± 0.032•	0.888 ± 0.032•	0.889 ± 0.035•	<b>0.911</b> ± 0.024
10	0.571 ± 0.059•	0.539 ± 0.055•	0.554 ± 0.059•	0.569 ± 0.066•	0.555 ± 0.079•	0.545 ± 0.058•	<b>0.595</b> ± 0.050
11	0.763 ± 0.062•	0.759 ± 0.060•	0.766 ± 0.063•	0.762 ± 0.067•	0.765 ± 0.062•	0.760 ± 0.063•	<b>0.802</b> ± 0.048
12	0.912 ± 0.023•	0.915 ± 0.025	0.880 ± 0.027•	0.914 ± 0.034	0.914 ± 0.031	0.914 ± 0.023	<b>0.919</b> ± 0.022
13	<b>0.731</b> ± 0.026	0.703 ± 0.031•	0.731 ± 0.030	0.712 ± 0.030•	0.712 ± 0.031•	0.717 ± 0.025•	0.726 ± 0.031
14	0.750 ± 0.017	0.732 ± 0.022•	0.748 ± 0.017•	0.750 ± 0.018	0.750 ± 0.018	0.750 ± 0.017•	<b>0.754</b> ± 0.017
15	0.719 ± 0.022	0.686 ± 0.027•	0.717 ± 0.023	0.702 ± 0.030•	0.702 ± 0.030•	0.704 ± 0.024•	<b>0.721</b> ± 0.028
16	<b>0.935</b> ± 0.012	0.868 ± 0.025•	0.915 ± 0.009•	0.933 ± 0.015	0.934 ± 0.015	0.932 ± 0.016	0.934 ± 0.013
17	0.885 ± 0.037	0.880 ± 0.043	0.885 ± 0.037	<b>0.888</b> ± 0.039	0.876 ± 0.049	0.885 ± 0.037	0.885 ± 0.037
18	<b>0.902</b> ± 0.000	0.650 ± 0.054•	0.803 ± 0.076	0.631 ± 0.189•	0.624 ± 0.189•	0.763 ± 0.049•	0.794 ± 0.080
19	<b>0.912</b> ± 0.008	0.695 ± 0.036•	0.890 ± 0.028	0.870 ± 0.040•	0.869 ± 0.047•	0.845 ± 0.062•	0.884 ± 0.027
20	<b>0.921</b> ± 0.002	0.629 ± 0.042•	0.843 ± 0.054	0.817 ± 0.110	0.819 ± 0.109	0.766 ± 0.086•	0.836 ± 0.079
Rank	3.05	5.70	4.90	4.25	4.10	4.35	1.65

Table 8 Comparison of pure accuracy on data sets with 5% noise

ID	KLR	p-cut	Bisection- $F_1$	Grid-search	$S_2$ -search	Bisection-PA	$\tau$ -interval
1	0.308 ± 0.024•	0.309 ± 0.024•	0.275 ± 0.029•	0.313 ± 0.023•	0.314 ± 0.023•	0.308 ± 0.024•	<b>0.322</b> ± 0.023
2	0.276 ± 0.028•	0.288 ± 0.027	0.267 ± 0.027•	0.290 ± 0.028	<b>0.290</b> ± 0.028	0.286 ± 0.030	0.290 ± 0.028
3	0.332 ± 0.035•	0.332 ± 0.033•	0.300 ± 0.032•	0.337 ± 0.035•	0.336 ± 0.036•	0.335 ± 0.035•	<b>0.345</b> ± 0.033
4	0.279 ± 0.026•	0.286 ± 0.026	0.252 ± 0.025•	0.289 ± 0.021	0.289 ± 0.021	0.287 ± 0.024	<b>0.291</b> ± 0.022
5	0.305 ± 0.023•	0.311 ± 0.027•	0.274 ± 0.028•	0.314 ± 0.022	0.314 ± 0.021	0.312 ± 0.026•	<b>0.319</b> ± 0.027
6	0.648 ± 0.057•	0.647 ± 0.053•	0.539 ± 0.075•	0.643 ± 0.054•	0.644 ± 0.054•	0.649 ± 0.058•	<b>0.670</b> ± 0.049
7	0.478 ± 0.105•	0.476 ± 0.117•	0.451 ± 0.118•	0.449 ± 0.119•	0.446 ± 0.117•	0.474 ± 0.112•	<b>0.534</b> ± 0.097
8	0.620 ± 0.058•	0.621 ± 0.053•	0.559 ± 0.069•	0.617 ± 0.067•	0.614 ± 0.072•	0.622 ± 0.060•	<b>0.642</b> ± 0.064
9	0.761 ± 0.071•	0.758 ± 0.065•	0.672 ± 0.068•	0.759 ± 0.068•	0.763 ± 0.068•	0.764 ± 0.073•	<b>0.807</b> ± 0.054
10	0.051 ± 0.127•	0.044 ± 0.113•	0.049 ± 0.117•	0.064 ± 0.094•	0.059 ± 0.096•	0.050 ± 0.117•	<b>0.111</b> ± 0.094
11	0.471 ± 0.136•	0.464 ± 0.132•	0.474 ± 0.139•	0.469 ± 0.145•	0.468 ± 0.139•	0.466 ± 0.138•	<b>0.552</b> ± 0.112
12	0.807 ± 0.052•	0.816 ± 0.053	0.728 ± 0.064•	0.816 ± 0.071	0.815 ± 0.066	0.812 ± 0.050•	<b>0.824</b> ± 0.049
13	0.391 ± 0.062	0.379 ± 0.063	0.376 ± 0.079	0.382 ± 0.060	0.381 ± 0.063	0.387 ± 0.052	<b>0.394</b> ± 0.066
14	0.389 ± 0.042	0.373 ± 0.042•	0.385 ± 0.041•	0.393 ± 0.044	0.393 ± 0.044	0.393 ± 0.040	<b>0.400</b> ± 0.042
15	0.306 ± 0.057•	0.331 ± 0.055•	0.303 ± 0.066•	0.320 ± 0.066•	0.319 ± 0.071•	0.339 ± 0.051	<b>0.355</b> ± 0.060
16	<b>0.759</b> ± 0.049	0.598 ± 0.055•	0.655 ± 0.042•	0.754 ± 0.054	0.755 ± 0.055	0.752 ± 0.053	0.758 ± 0.047
17	0.408 ± 0.148	0.402 ± 0.151	0.405 ± 0.146	<b>0.411</b> ± 0.147	0.392 ± 0.155	0.410 ± 0.143	0.408 ± 0.148
18	-0.000 ± 0.000•	0.043 ± 0.053	0.043 ± 0.055	0.049 ± 0.032	0.047 ± 0.030	0.040 ± 0.055	<b>0.053</b> ± 0.035
19	0.069 ± 0.092•	0.108 ± 0.063•	0.161 ± 0.094	0.159 ± 0.098	0.150 ± 0.105•	0.163 ± 0.104•	<b>0.195</b> ± 0.110
20	0.003 ± 0.015•	0.051 ± 0.036	<b>0.061</b> ± 0.062	0.043 ± 0.057	0.041 ± 0.058	0.058 ± 0.048	0.046 ± 0.054
Rank	4.70	4.95	5.75	3.45	4.10	3.60	1.45

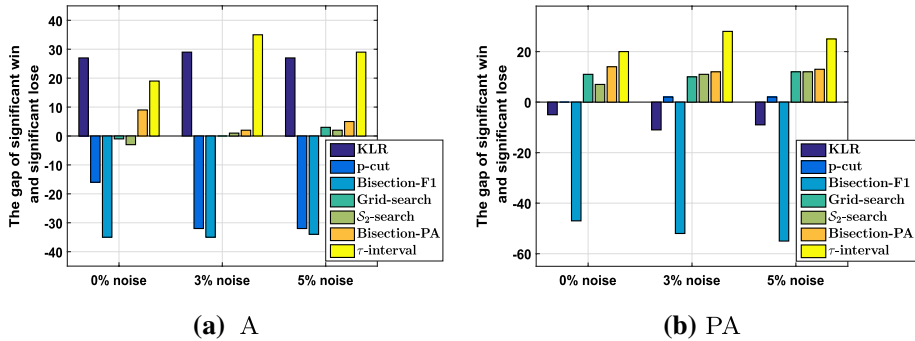


Fig. 5 Statistical comparison under different noise level. Each bar is the gap between the times of the significant wins and the times of significant loses. The significant lose and win is defined according to Eq. (44)

564 **7 Conclusion**

565 With an increase in the complexity of the data, eliminating random consistency from learn-  
 566 ing algorithms has great potential to improve the generalization ability. In this paper, first,  
 567 we have shown that the PA is insensitive to the class distribution of classifiers in evaluation  
 568 and is more fairer than the A in learning classifiers through two vivid examples. Second,  
 569 we have given some novel bounds to show that learning by PA can approach to the optimal  
 570 A and have shown that the empirical risk minimization process of the PA is Bayes-risk  
 571 consistent. Based on these theoretical guarantees, we have proposed a plug-in rule model  
 572 that optimizes the PA. The experimental results have shown the fairness and effectiveness  
 573 of learning by PA. An interesting future work is to establish the other strategies to define  
 574 the random consistency. An analysis of the random consistency for each instance maybe a  
 575 promising direction.

576 **Acknowledgements** This work is supported by National Key R&D Program of China (No.  
 577 2018YFB1004300), the National Natural Science Foundation of China (Nos. 61672332, 61872226,  
 578 61976129), the Overseas Returnee Research Program of Shanxi Province (No. 2017023), the Natural Sci-  
 579 ence Foundation of Shanxi Province (No. 201701D121052), and Program for the San Jin Young Scholars of  
 580 Shanxi.

581 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License,  
 582 which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long  
 583 as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Com-  
 584 mons licence, and indicate if changes were made. The images or other third party material in this article  
 585 are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the  
 586 material. If material is not included in the article's Creative Commons licence and your intended use is not  
 587 permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly  
 588 from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

589 **Appendix: Proofs**

590 **Lemma 1** *When the partitions in  $\mathcal{H}^{q(h)}$  are distributed uniformly, the expectation accu-*  
 591 *racy of partitions in  $\mathcal{H}^{q(h)}$  is:*



$$\mathbb{E}_{h' \in \mathcal{H}^{q(h)}} A(h') = pq(h) + (1-p)(1-q(h)). \tag{45}$$

**Proof** Without loss of generality, we assume that  $q(h) < p$ . Assuming that the size of data is  $N$ , we have

$$\mathbb{P}_{h' \in \mathcal{H}^{q(h)}} \left( TP(h') = \frac{j}{N} \right) = \frac{C_{Np}^j C_{N-Np}^{Nq(h)-j}}{C_N^{Nq(h)}}, \tag{46}$$

where  $j = 0, \dots, Nq(h)$  and  $C_n^m$  is the number of combinations of  $n$  items taken  $m$  at a time. From (46), we know that  $N \cdot TP(h')$  follows the hypergeometric distribution with the size of the population selected from be  $N$ ,  $Np$  elements of the population belonging to one group and  $N - Np$  belonging to the other group, and the number of samples drawn from the population be  $Nq(h)$ . Thus,

$$\mathbb{E}_{h' \in \mathcal{H}^{q(h)}} TP(h') = pq(h). \tag{47}$$

Then, according to  $TN(h') = 1 - p - (q(h) - TP(h'))$ , we have:

$$\begin{aligned} \mathbb{E}_{h' \in \mathcal{H}^{q(h)}} A(h') &= \mathbb{E}_{h' \in \mathcal{H}^{q(h)}} 1 - p - q(h) + 2TP(h') \\ &= 1 - p - q(h) + 2pq(h). \end{aligned} \tag{48}$$

□

**Example 2** Assume that two-class data are generated from two Gaussian distributions with uncommon means  $\mu_1, \mu_2$ , but a common covariance  $\Sigma$ :

$$m(\mathbf{x}|y = +1) = \mathcal{N}(\mu_1, \Sigma), \tag{49}$$

$$m(\mathbf{x}|y = -1) = \mathcal{N}(\mu_2, \Sigma) \tag{50}$$

and the probability of the positive class is  $p = \mathbb{P}(Y = +1)$ . The label of the minority class is corrupted by the instance-independent noise at the level  $s_1$ :  $\mathbb{P}(\tilde{Y} = -1|Y = +1) = s_1$ . For this learning task, the bias of  $h_A^*$  is:

$$\begin{aligned} \text{Bias}(h_A^*) &= \left| \mathbb{P}(d(X) < d_0|Y = -1) - \mathbb{P}(d(X) > d_0|Y = +1) \right| \end{aligned} \tag{51}$$

$$= \left| \Phi\left(\frac{d_0 + \Delta/2}{\sqrt{\Delta}}\right) - 1 + \Phi\left(\frac{d_0 - \Delta/2}{\sqrt{\Delta}}\right) \right|, \tag{52}$$

where  $\Phi(\bullet)$  is the cumulative distribution function of the standard normal distribution,  $\Delta = (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)$  and  $d_0 = \ln \frac{1-p}{p} \frac{1}{1-2s_1}$ .

**Proof** According to Lemma 2, the corrupted conditional class probability  $\mathbb{P}(\tilde{Y} = +1|X = \mathbf{x})$  is needed. Based on the Bayes' theorem:

$$\mathbb{P}(\tilde{Y} = +1|X = \mathbf{x}) \tag{53}$$

Author Proof

$$= \frac{\mathbf{m}(\mathbf{x}|\tilde{y} = +1)\mathbb{P}(\tilde{Y} = +1)}{\sum_{l \in \{-1, +1\}} \mathbf{m}(\mathbf{x}|\tilde{y} = l)\mathbb{P}(\tilde{Y} = l)}. \quad (54)$$

629  
630 Because  $\mathbb{P}(\tilde{Y} = +1) = p(1 - s_1)$ ,  $\mathbf{m}(\mathbf{x}|\tilde{y} = +1) = \mathcal{N}(\boldsymbol{\mu}_1, \Sigma)$  and

$$\mathbf{m}(\mathbf{x}|\tilde{y} = -1) = \sum_{l \in \{+1, -1\}} \mathbf{m}(\mathbf{x}|y = l, \tilde{y} = -1)\mathbb{P}(Y = l|\tilde{Y} = -1) \quad (55)$$

$$= \sum_{l \in \{+1, -1\}} \mathbf{m}(\mathbf{x}|y = l) \frac{\mathbb{P}(\tilde{Y} = +1|Y = l)\mathbb{P}(Y = l)}{\mathbb{P}(\tilde{Y} = -1)} \quad (56)$$

$$= \frac{(1-p)}{(1-p) + ps_1} \mathcal{N}(\boldsymbol{\mu}_2, \Sigma) + \frac{ps_1}{(1-p) + ps_1} \mathcal{N}(\boldsymbol{\mu}_1, \Sigma), \quad (57)$$

637 where  $\mathbf{m}(\mathbf{x}|y = l, \tilde{y} = -1) = \mathbf{m}(\mathbf{x}|y = l)$  is satisfied because the label noise is independent  
638 on instance:  $\mathbb{P}(\tilde{Y} = -1|Y = l, X = \mathbf{x}) = \mathbb{P}(\tilde{Y} = -1|Y = l)$ , we have:

$$\mathbb{P}(\tilde{Y} = +1|X = \mathbf{x}) = \frac{1 - s_1}{1 + \exp(\mathbf{w}^T \mathbf{x} + b)} \quad (58)$$

640 with  $\mathbf{w}^T = (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \Sigma^{-1}$  and  $b = \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 - \ln p - \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln(1-p)$ .

642 Further, according to Lemma 2, the optimal classifier in the sense of accuracy is

$$h_A^*(\mathbf{x}) = \begin{cases} +1, & d(\mathbf{x}) > d_0, \\ -1, & \text{otherwise.} \end{cases} \quad (59)$$

644 where  $d(\mathbf{x}) = \mathbf{x}' \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)' \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$  and  $d_0 = \ln \frac{1-p}{p} \frac{1}{1-2s_1}$ .

646 According to the additivity of the Gaussian distribution, we obtain the probability mass  
647 function of  $d(\mathbf{x})$ :

$$\mathbf{m}(d(\mathbf{x})|y = +1) = \mathcal{N}(\Delta/2, \Delta), \quad (60)$$

$$\mathbf{m}(d(\mathbf{x})|y = -1) = \mathcal{N}(-\Delta/2, \Delta) \quad (61)$$

652 where  $\Delta = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ . Then

$$\text{Bias}(h_A^*) = \left| \mathbb{P}(d(X) < d_0|Y = -1) - \mathbb{P}(d(X) > d_0|Y = +1) \right| \quad (62)$$

$$= \left| \Phi\left(\frac{d_0 + \Delta/2}{\sqrt{\Delta}}\right) - 1 + \Phi\left(\frac{d_0 - \Delta/2}{\sqrt{\Delta}}\right) \right|, \quad (63)$$

657 where  $\Phi(\bullet)$  is the cumulative distribution function of the standard normal distribution.

658 □

659 **Theorem 1** *The classifier that maximizes the PA is*

$$h_{PA}^*(\mathbf{x}) = \arg \max_h PA(h) \quad (64)$$

661

$$= \begin{cases} +1, & \eta(\mathbf{x}) > (\frac{1}{2} - p)PA^* + p, \\ -1, & \text{otherwise.} \end{cases} \quad (65)$$

662

663 where  $PA^* = PA(h_{PA}^*)$  and  $p = \mathbb{P}(Y = +1)$ .

665 **Proof** The formulation of the pure accuracy measure is fractional, which hinders obtaining  
666 the optimal classifier. Here, we resort to the cost-sensitive loss to obtain a non-closed-form  
667 solution. We begin this proof with two existing definitions and two lemmas:

668 **Definition 5** (Kotlowski and Dembczynski 2017) We refer to a measure as a linear-fractional  
669 performance measure if it is non-increasing with  $FP$ ,  $FN$  and formalized as

$$\Psi(FP, FN) = \frac{a_0 + a_1 FP + a_2 FN}{b_0 + b_1 FP + b_2 FN}, \quad (66)$$

670

671 where  $a_0, a_1, a_2, b_0, b_1, b_2 \in \mathcal{R}$  and  $b_0 + b_1 FP + b_2 FN \geq C_1 > 0$ .

673 **Definition 6** (Elkan 2001) The cost-sensitive loss is defined as  
674  $L_\rho(h) = \rho FP(h) + (1 - \rho)FN(h)$ , where  $\rho \in (0, 1)$ .

675 **Lemma 7** (Kotlowski and Dembczynski 2017) The regret w.r.t the linear-fractional per-  
676 formance measure  $\Psi(FP, FN)$  can be bounded by that w.r.t.  $L_\rho(h)$  when  $\rho = \frac{\Psi^* b_1 - a_1}{\Psi^*(b_1 + b_2) - (a_1 + a_2)}$

$$\Psi^* - \Psi(h) \leq C_2(L_\rho(h) - L_\rho^*), \quad (67)$$

677

678 where  $\Psi^* = \max_h \Psi(h)$ ,  $L_\rho^* = \min_h L_\rho(h)$  and  $C_2 = \frac{1}{C_1}(\Psi^*(b_1 + b_2) - (a_1 + a_2))$ .

680 **Lemma 8** (Elkan 2001) The classifier that minimizes  $L_\rho$  is

$$h_{L_\rho}^*(\mathbf{x}) = \begin{cases} +1, & \eta(\mathbf{x}) > \rho, \\ -1, & \text{otherwise.} \end{cases} \quad (68)$$

681

682 where  $\eta(\mathbf{x}) = \mathbb{P}(Y = 1|X = \mathbf{x})$ .

684 Because  $A = 1 - FN - FP$ ,  $RA = 1 - p - q(h) + 2pq$  and  $q(h) = p + FP - FN$ , we have

$$PA = \frac{A - RA}{1 - RA} = \frac{p(1 - p) - pFP - (1 - p)FN}{p(1 - p) + (\frac{1}{2} - p)(FP - FN)}. \quad (69)$$

685

686 According to Lemma 7, the regret of the PA can be bounded by that of  $L_\rho$  with  
687  $\rho = (\frac{1}{2} - p)PA^* + p$ . Then by Lemma 8, we obtain the formulation.  $\square$

689 **Lemma 3** For all distributions, the plug-in rule with  $\rho$  as the decision threshold

$$h_\rho(\mathbf{x}) = \begin{cases} +1, & \eta(\mathbf{x}) > \rho, \quad \text{where } \rho \in (0, \frac{1}{2}], \\ -1, & \text{otherwise,} \end{cases} \quad (70)$$

690

691 satisfies:

$$L(h_\rho) \leq \frac{1-\rho}{\rho} L^*, \quad (71)$$

694  
695 when  $\rho = 1/2$ , the equality holds.

$$L(h_\rho) = \mathbb{P}(h_\rho(X) = -1, Y = +1) + \mathbb{P}(h_\rho(X) = +1, Y = -1) \quad (72)$$

697 **Proof**

$$= \mathbb{E}_{X:\eta(X)<\rho} \eta(X) + \mathbb{E}_{X:\eta(X)\geq\rho} (1 - \eta(X)) \quad (73)$$

$$\leq \mathbb{E}_{X:\eta(X)<\rho} (1/\rho - 1)\eta(X) + \mathbb{E}_{X:\eta(X)\geq\rho} (1 - \eta(X)) \quad (74)$$

$$= \mathbb{E}_X \min\{(1/\rho - 1)\eta(X), 1 - \eta(X)\} \quad (75)$$

$$\leq (1/\rho - 1)\mathbb{E}_X \min\{\eta(X), 1 - \eta(X)\}. \quad (76)$$

706 □

707 **Lemma 4** For all distributions, suppose that  $p = \mathbb{P}(Y = +1) \leq \frac{1}{2}$ , the pure loss of  $h_A^*$   
708 satisfies:

$$PL(h_A^*) \leq \frac{L^*}{p\left(\frac{3}{2} - p\right) - L^*\left(\frac{1}{2} - p\right)}. \quad (77)$$

710  
711 **Proof** Let  $q_A^* = \mathbb{P}(h_A^* = +1)$ ,  $FP_A^* = \mathbb{P}(h_A^* = +1, Y = -1)$  and  $FN_A^* = \mathbb{P}(h_A^* = -1, Y = +1)$ .  
712 By definition,

$$PL(h_A^*) = \frac{L^*}{p + (1 - 2p)q_A^*}. \quad (78)$$

714 To obtain the upper bound of  $PL(h_A^*)$ , we derive the lower bound of  $q_A^*$ . Because:

$$L^* = \mathbb{E}_{X:\eta(X)\leq 1/2} \eta(X) + \mathbb{E}_{X:\eta(X)> 1/2} (1 - \eta(X)) \quad (79)$$

$$= \mathbb{E}_X \eta(X) - \mathbb{E}_{X:\eta(X)> 1/2} \eta(X) + \mathbb{E}_{X:\eta(X)> 1/2} (1 - \eta(X)) \quad (80)$$

$$= p - \mathbb{E}_X \max\{2\eta(X) - 1, 0\}, \quad (81)$$

722 and then, we have

$$q_A^* = \mathbb{E}_X \mathbf{I}\{\eta(X) - 1/2 \geq 0\} \quad (82)$$

$$\geq \mathbb{E}_X \max\{\eta(X) - 1/2, 0\} \quad (83)$$

$$= \frac{1}{2}(p - L^*), \quad (84)$$

728

729 where  $\mathbf{I}\{\bullet\}$  is the indicator function. Putting the lower bound of  $q_A^*$  into the formulation of  
 730  $PL(h_A^*)$ , we obtain the upper bound of  $PL(h_A^*)$ .  $\square$

731 **Theorem 3** For all distributions, suppose  $p \leq \frac{1}{2}$ , the pure loss of  $h_A^*$  satisfies:

$$732 \quad PL(h_{pA}^*) \leq PL(h_A^*) \quad (85)$$

733

$$734 \quad \leq \frac{2(1-p)}{p(3-2p) - L^*(1-2p)} PL(h_{pA}^*). \quad (86)$$

735

736 **Proof** For any  $q(h)$ , we have

$$737 \quad 1 - RA = p + (1 - 2p)q(h) \leq 1 - p, \quad (87)$$

738

739 hence

$$740 \quad L = (1 - RA)PL \leq (1 - p)PL. \quad (88)$$

741

742 Further amplifying the upper bound in Lemma 4:

$$743 \quad L^* \leq L(h_{pA}^*) \leq (1 - p)PL(h_{pA}^*), \quad (89)$$

744

745 we obtain the result.  $\square$

746 **Lemma 5** For two random variables  $Z_1, Z_2 \in [0, 1]$ , any  $\varepsilon \in (0, 1]$ , let  
 747  $\alpha = \mathbb{E}Z_1\mathbb{E}Z_2 / (2\mathbb{E}Z_1 + \mathbb{E}Z_2)$ , we have

$$748 \quad \mathbb{P}\left(\left|\frac{Z_1}{Z_2} - \frac{\mathbb{E}Z_1}{\mathbb{E}Z_2}\right| > \varepsilon\right) \quad (90)$$

$$\leq \mathbb{P}(|Z_1 - \mathbb{E}Z_1| > \alpha\varepsilon) + 3\mathbb{P}(|Z_2 - \mathbb{E}Z_2| > \alpha\varepsilon).$$

749

750 **Proof** For  $\beta \in [0, 1]$  and  $\gamma > 0$ , we have

$$751 \quad \mathbb{P}\left(\left|\frac{Z_1}{Z_2} - \frac{\mathbb{E}Z_1}{\mathbb{E}Z_2}\right| > \varepsilon\right) \quad (91)$$

$$= \mathbb{P}\left(\left|\frac{Z_1 - \mathbb{E}Z_1}{(Z_2 - \mathbb{E}Z_2) + \mathbb{E}Z_2} + \frac{(\mathbb{E}Z_2 - Z_2)\mathbb{E}Z_1}{(Z_2 - \mathbb{E}Z_2)\mathbb{E}Z_2 + (\mathbb{E}Z_2)^2}\right| > \varepsilon\right)$$

752

$$753 \quad \leq \mathbb{P}\left(\left|\frac{Z_1 - \mathbb{E}Z_1}{(Z_2 - \mathbb{E}Z_2) + \mathbb{E}Z_2}\right| > \beta\varepsilon\right) + \mathbb{P}\left(\left|\frac{(\mathbb{E}Z_2 - Z_2)\mathbb{E}Z_1}{(Z_2 - \mathbb{E}Z_2)\mathbb{E}Z_2 + (\mathbb{E}Z_2)^2}\right| > (1 - \beta)\varepsilon\right) \quad (92)$$

754

$$755 \quad \leq \mathbb{P}(|Z_1 - \mathbb{E}Z_1| > \beta|\mathbb{E}Z_2 - \gamma\varepsilon|\varepsilon) + 2\mathbb{P}(|Z_2 - \mathbb{E}Z_2| > \gamma\varepsilon)$$

$$+ \mathbb{P}\left(|Z_2 - \mathbb{E}Z_2| > \frac{(1 - \beta)\mathbb{E}Z_1}{\mathbb{E}Z_2}|\mathbb{E}Z_2 - \gamma\varepsilon|\varepsilon\right), \quad (93)$$

756

757 where the first inequality is obtained by

$$758 \quad \mathbb{P}(|a + b| > \varepsilon) \leq \mathbb{P}(|a| > \beta\varepsilon) + \mathbb{P}(|b| > (1 - \beta)\varepsilon), \quad (94)$$

759

760 and the second inequality is obtained by

$$\mathbb{P}(B_1) = \mathbb{P}(B_1|B_2)\mathbb{P}(B_2) + \mathbb{P}(B_1|B_2^c)\mathbb{P}(B_2^c) \quad (95)$$

$$\leq \mathbb{P}(B_1|B_2) + \mathbb{P}(B_2^c), \quad (96)$$

for any events  $B_1, B_2$ , where  $B_2^c$  complementary set of  $B_2$ . Here, we take the event  $|Z_2 - \mathbb{E}Z_2| \leq \gamma\epsilon$  as  $B_2$  to divide the two terms of the first inequality.

Let

$$\beta = \mathbb{E}Z_1 / (\mathbb{E}Z_1 + \mathbb{E}Z_2), \quad (97)$$

$$\gamma = \mathbb{E}Z_1\mathbb{E}Z_2 / (\mathbb{E}Z_1 + \epsilon\mathbb{E}Z_1 + \mathbb{E}Z_2), \quad (98)$$

we have

$$\begin{aligned} & \beta|\mathbb{E}Z_2 - \gamma\epsilon| \\ &= (1 - \beta)\mathbb{E}Z_1|\mathbb{E}Z_2 - \gamma\epsilon|/\mathbb{E}Z_2 \end{aligned} \quad (99)$$

$$= \gamma. \quad (100)$$

Finally, by the assumption  $\epsilon \leq 1$ , we have  $\gamma \geq \alpha$  and then get the result.  $\square$

**Theorem 4** Suppose the cardinality of  $\mathcal{H}$  is finite:  $|\mathcal{H}| < \infty$ , then for every  $h \in \mathcal{H}$ , any  $\epsilon \in (0, 1]$ , we have

$$\mathbb{P}\left\{\sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon\right\} \leq 8|\mathcal{H}| \exp\left\{-2N\left(\alpha\epsilon - \frac{Rc(\mathcal{H})}{2}\right)^2\right\}, \quad (101)$$

where  $\alpha = \min_{h \in \mathcal{H}} \frac{L(h)}{2PL(h)+1}$  and  $Rc(\mathcal{H})$  is the Rademacher complexity of  $\mathcal{H}$ .

**Proof** First, we process the superior limit in probability according to the union bound:

$$\begin{aligned} & \mathbb{P}\left\{\sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon\right\} \\ &= \mathbb{P}\left\{\exists h \in \mathcal{H} : \left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon\right\} \\ &\leq \sum_{h \in \mathcal{H}} \mathbb{P}\left\{\left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon\right\} \\ &\leq |\mathcal{H}| \sup_{h \in \mathcal{H}} \mathbb{P}\left\{\left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon\right\}. \end{aligned} \quad (102)$$

Second, we transform the gap in the sense of PL into that of L by Lemma 5. For every  $h \in \mathcal{H}$ , let  $\alpha = \min_{h \in \mathcal{H}} \frac{L(h)}{2PL(h)+1}$ , we have:

$$\begin{aligned} & \mathbb{P}\left\{\left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon\right\} \\ &\leq \mathbb{P}\left\{\left| \widehat{L}_N - L \right| > \alpha\epsilon\right\} + 3\mathbb{P}\left\{\left| \widehat{RA}_N - RA \right| > \alpha\epsilon\right\}. \end{aligned} \quad (103)$$

790 Third, applying Theorem 8 in Bartlett and Mendelson (2003), for every  $h \in \mathcal{H}$ , with  
791 probability at least  $1 - \delta/4$ , we obtain that:

$$792 \quad \left| \widehat{L}_N - L \right| \leq \frac{Rc(\mathcal{H})}{2} + \sqrt{\frac{\ln(8/\delta)}{2N}}. \quad (104)$$

793  
794 Let  $\delta = 8 \exp \{-2N(\alpha\epsilon - Rc(\mathcal{H})/2)^2\}$ , and then:

$$795 \quad \mathbb{P} \left\{ \left| \widehat{L}_N - L \right| > \alpha\epsilon \right\} \leq \delta/4. \quad (105)$$

796  
797 For the second term in (103), by  $|\mathcal{H}^{q(h)}| = C_N^{Ng(h)}$  and the triangle inequality, we have:

$$798 \quad \left| \widehat{RA}_N(h) - RA(h) \right| \leq \frac{1}{C_N^{Ng(h)}} \sum_{j=1}^{C_N^{Ng(h)}} \left| \widehat{L}_N(h_j) - L(h_j) \right|. \quad (106)$$

799  
800 According to Theorem 8 in Bartlett and Mendelson (2003), for every function  
801  $h_j \in \mathcal{H}^{q(h)}$ , with probability at least  $1 - \delta/4$ , holds that:

$$802 \quad \left| \widehat{L}_N(h_j) - L(h_j) \right| \leq \frac{Rc(\mathcal{H}^{q(h)})}{2} + \sqrt{\frac{\ln(8/\delta)}{2N}}, \quad (107)$$

803 because  $\mathcal{H}^{q(h)} \subseteq \mathcal{H}$ , we have  $Rc(\mathcal{H}^{q(h)}) \leq Rc(\mathcal{H})$ , and then for every function  $h_j \in \mathcal{H}^{q(h)}$ ,  
805 with probability at least  $1 - \delta/4$ , holds that:

$$806 \quad \left| \widehat{RA}_N(h) - RA(h) \right| \leq \frac{Rc(\mathcal{H})}{2} + \sqrt{\frac{\ln(8/\delta)}{2N}}. \quad (108)$$

807  
808 Putting  $\delta$  into inequality (108), we obtain for every  $h \in \mathcal{H}$ :

$$809 \quad \mathbb{P} \left\{ \left| \widehat{RA}_N(h) - RA(h) \right| > \alpha\epsilon \right\} \leq \delta/4. \quad (109)$$

810  
811 Thus, combining (102), (103), (105) and (109), we obtain the final result.  $\square$

812 **Lemma 6** Let  $\mathcal{S}'_N = \{(\mathbf{x}'_1, y'_1), \dots, (\mathbf{x}'_N, y'_N)\}$  be an independent and identically distributed  
813 collection as  $\mathcal{S}_N$  and  $PL'_N(h)$  is the corresponding empirical pure loss. Suppose  
814  $N \geq 5(6 + 4\alpha\epsilon)\alpha^{-2}\epsilon^{-2}$ , where  $\alpha = \min_{h \in \mathcal{H}} \frac{L(h)}{2PL(h)+1}$ ,  $\epsilon \in (0, 1]$ , then we have

$$815 \quad \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon \right\} \\ \leq 2\mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - \widehat{PL}'_N(h) \right| > \frac{\epsilon}{2} \right\}. \quad (110)$$

816  
817 **Proof** There exists at least one function  $h_0 \in \mathcal{H}$  satisfies  $\left| \widehat{PL}_N(h_0) - PL(h_0) \right| \geq \epsilon$ . For  $h_0$ ,

$$818 \quad \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - \widehat{PL}'_N(h) \right| > \frac{\epsilon}{2} \right\} \\ \geq \mathbb{E}_{\mathcal{S}_N} \left[ \mathbf{I} \left( \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \epsilon \right) \mathbb{P} \left\{ \left| \widehat{PL}'_N(h_0) - PL(h_0) \right| < \frac{\epsilon}{2} \mid \mathcal{S}_N \right\} \right]. \quad (111)$$

819

820 Here, we omit the detail proof of this inequality because the technique is the same as  
 821 Lemma 2 in Vapnik and Chervonenkis (1971) on accuracy.

822 According to Lemma 5, let  $\alpha = \min_{h \in \mathcal{H}} \frac{L(h)}{2PL(h)+1}$ , we have:

$$\begin{aligned}
 & \mathbb{P} \left\{ \left| \widehat{PL}'_N(h_0) - PL(h_0) \right| > \frac{\varepsilon}{2} \middle| \mathcal{S}_N \right\} \\
 & \leq \mathbb{P} \left\{ \left| \widehat{A}'_N(h_0) - A(h_0) \right| > \frac{\alpha\varepsilon}{2} \middle| \mathcal{S}_N \right\} + 3\mathbb{P} \left\{ \left| \widehat{RA}'_N(h_0) - RA(h_0) \right| > \frac{\alpha\varepsilon}{2} \middle| \mathcal{S}_N \right\}.
 \end{aligned}
 \tag{112}$$

824 For the first term of (112), according to the Bernstein's inequality, we have

$$\begin{aligned}
 & \mathbb{P} \left\{ \left| \widehat{A}'_N(h_0) - A(h_0) \right| > \frac{\alpha\varepsilon}{2} \middle| \mathcal{S}_N \right\} \\
 & \leq 2 \exp \left\{ - \frac{\frac{\alpha^2 \varepsilon^2 N}{4}}{2 \left( A(h_0)(1 - A(h_0)) + \frac{\alpha\varepsilon}{6} \right)} \right\} \\
 & \leq 2 \exp \left\{ - \frac{3\alpha^2 \varepsilon^2 N}{6 + 4\alpha\varepsilon} \right\} \\
 & \leq 2 \left\{ 1 + \frac{3\alpha^2 \varepsilon^2 N}{6 + 4\alpha\varepsilon} \right\}^{-1} \leq \frac{1}{8},
 \end{aligned}
 \tag{113}$$

827 where the second inequality is because for any  $\rho \in [0, 1]$ , it is satisfied that  $\rho(1 - \rho) \leq 1/4$ ,  
 828 the third inequality is obtained by  $e^{-x} \leq (1 + x)^{-1}$  for  $x > 0$  and the last inequality is  
 830 obtained by the assumption  $N \geq 5(6 + 4\alpha\varepsilon)\alpha^{-2}\varepsilon^{-2}$ .

831 For the second term of (112), by the definition of  $\widehat{RA}'_N(h_0)$ , the only difference in the  
 832 proof of the two terms in (112) is the number of terms for summation. Under the assump-  
 833 tion on  $N$ , we have  $NC_N^{Ng(h)} \geq 5(6 + 4\alpha\varepsilon)\alpha^{-2}\varepsilon^{-2}$ , and then using the same technique as  
 834 (113), we have

$$\mathbb{P} \left\{ \left| \widehat{RA}'_N(h_0) - RA(h_0) \right| > \frac{\alpha\varepsilon}{2} \middle| \mathcal{S}_N \right\} \leq \frac{1}{8}.
 \tag{114}$$

836 Combing the inequalities (112), (113), (114), we have

$$\mathbb{P} \left\{ \left| \widehat{PL}'_N(h_0) - \widehat{PL}(h_0) \right| > \frac{\varepsilon}{2} \middle| \mathcal{S}_N \right\} \leq \frac{1}{2}.
 \tag{115}$$

839 Thus, according to (111) and (115), we obtain the final result. □

841 **Theorem 5** As the same condition as Lemma 6 and suppose the VC dimension of  $\mathcal{H}$  is  
 842 finite:  $d_{vc}(\mathcal{H}) < \infty$ , we have

$$\begin{aligned}
 & \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{PL}_N(h) - PL(h) \right| > \varepsilon \right\} \\
 & \leq 4(N + 1) \exp \left\{ - \left( \frac{\varepsilon^2(1 - |\widehat{p}_N - 1|)^2}{16} - \frac{d_{vc}(\mathcal{H}) \ln(2eN/d_{vc}(\mathcal{H}))}{N} \right) N \right\}.
 \end{aligned}
 \tag{116}$$

844 **Proof** By Lemma 6,  
 845



$$\begin{aligned} & \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{P}L_N(h) - PL(h) \right| > \varepsilon \right\} \\ & \leq 2 \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{P}L_N(h) - \widehat{P}L'_N(h) \right| > \frac{\varepsilon}{2} \right\}. \end{aligned} \quad (117)$$

We divide the hypothesis space  $\mathcal{H}$  into  $N + 1$  subspaces according to the class distribution of hypothesis function:  $\mathcal{H} = \bigcup_{\hat{q}_N \in \{0, \frac{1}{N}, \dots, 1\}} \mathcal{H}^{\hat{q}_N}$ , where  $\mathcal{H}^{\hat{q}_N} = \{h : \frac{1}{N} \sum_{i=1}^N \mathbf{I}[h(X_i) = +1] = \hat{q}_N, h \in \mathcal{H}\}$ . Thus, according to the definition of pure loss and the union bound, we have

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{h \in \mathcal{H}} \left| \widehat{P}L_N(h) - \widehat{P}L'_N(h) \right| > \frac{\varepsilon}{2} \right\} \\ & = \mathbb{P} \left\{ \sup_{\hat{q}_N} \sup_{h \in \mathcal{H}^{\hat{q}_N}} \left| \widehat{P}L_N(h) - \widehat{P}L'_N(h) \right| > \frac{\varepsilon}{2} \right\} \\ & \leq (N + 1) \sup_{\hat{q}_N} \mathbb{P} \left\{ \sup_{h \in \mathcal{H}^{\hat{q}_N}} \left| \widehat{P}L_N(h) - \widehat{P}L'_N(h) \right| > \frac{\varepsilon}{2} \right\} \\ & = (N + 1) \sup_{\hat{q}_N} \mathbb{P} \left\{ \sup_{h \in \mathcal{H}^{\hat{q}_N}} \left| \widehat{L}_N(h) - \widehat{L}'_N(h) \right| > \frac{\varepsilon(1 - \widehat{R}A_N)}{2} \right\}. \end{aligned} \quad (118)$$

We employ Theorem 3.1 in Vapnik and Chervonenkis (1971) for the error terms. Besides, for any  $\hat{q}_N$ , it satisfies that  $1 - \widehat{R}A_N \geq \frac{1 - |2\hat{p}_N - 1|}{2}$  and  $d_{vc}(\mathcal{H}^{\hat{q}_N}) \leq d_{vc}(\mathcal{H})$  for  $\mathcal{H}^{\hat{q}_N} \subseteq \mathcal{H}$ . Then we have:

$$\begin{aligned} & \sup_{\hat{q}_N} \mathbb{P} \left\{ \sup_{h \in \mathcal{H}^{\hat{q}_N}} \left| \widehat{L}_N(h) - \widehat{L}'_N(h) \right| > \frac{\varepsilon(1 - \widehat{R}A_N)}{2} \right\} \\ & \leq 2 \sup_{\hat{q}_N} \exp \left\{ - \left( \frac{\varepsilon^2(1 - \widehat{R}A_N)^2}{4} - \frac{d_{vc}(\mathcal{H}^{\hat{q}_N})[\ln(2eN/d_{vc}(\mathcal{H}^{\hat{q}_N}))]}{N} \right) N \right\} \\ & \leq 2 \exp \left\{ - \left( \frac{\varepsilon^2(1 - |2\hat{p}_N - 1|)^2}{16} - \frac{d_{vc}(\mathcal{H}) \ln(2eN/d_{vc}(\mathcal{H}))}{N} \right) N \right\}. \end{aligned} \quad (119)$$

Combining (117), (118) and (119), we obtain the final result.  $\square$

## References

- Agarwal, S., Graepel, T., Herbrich, R., Har-Peled, S., & Roth, D. (2005a). Generalization bounds for the area under the ROC curve. *Journal of Machine Learning Research*, 6(2), 393–425.
- Agarwal, S., Harpeled, S., & Roth, D. (2005b). A uniform convergence bound for the area under the ROC curve. In *Proceedings of the international conference on artificial intelligence and statistics* (pp. 1–8).
- Albatineh, A. N., & Niewiadomska-Bugaj, M. (2011). Correcting Jaccard and other similarity indices for chance agreement in cluster analysis. *Advances in Data Analysis and Classification*, 5(3), 179–200.
- Albatineh, A. N., Niewiadomska-Bugaj, M., & Mihalko, D. (2006). On similarity indices and correction for chance agreement. *Journal of Classification*, 23(2), 301–313.

- 869 Alcalá-fdez, J., Sanchez, L., Garcia, S., Jesus, M. J. D., Ventura, S., Garrell, J. M., et al. (2008). KEEL:  
870 A software tool to assess evolutionary algorithms for data mining problems. *Soft Computing*, 13(3),  
871 307–318.
- 872 Bartlett, P. L., & Mendelson, S. (2003). Rademacher and Gaussian complexities: Risk bounds and structural  
873 results. *Journal of Machine Learning Research*, 3(3), 463–482.
- 874 Bartlett, P. L., Jordan, M. I., & McAuliffe, J. (2006). Convexity, classification, and risk bounds. *Journal of*  
875 *the American Statistical Association*, 101(473), 138–156.
- 876 Blair, E., & Stanley, F. J. (2008). Interobserver agreement in the classification of cerebral palsy. *Develop-*  
877 *mental Medicine & Child Neurology*, 27(5), 615–622.
- 878 Budescu, D., & Bar-Hillel, M. (1993). To guess or not to guess: A decision-theoretic view of formula scor-
- 879 ing. *Journal of Educational Measurement*, 30(4), 277–291.
- 880 Cameron, M. L., Briggs, K. K., & Steadman, J. R. (2003). Reproducibility and reliability of the outerbridge  
881 classification for grading chondral lesions of the knee arthroscopically. *American Journal of Sports*  
882 *Medicine*, 31(1), 83–86.
- 883 Chong, E. K. P., & Żak, S. H. (2011). *An introduction to optimization* (3rd ed.). New York: Wiley.
- 884 Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measure-*  
885 *ment*, 20(1), 37–46.
- 886 Devroye, L., Györfi, L., & Lugosi, G. (1996). *A probabilistic theory of pattern recognition*. New York:  
887 Springer.
- 888 Diamond, J., & Evans, W. (1973). The correction for guessing. *Review of Educational Research*, 43(2),  
889 181–191.
- 890 Dua, D., & Graff, C. (2017). UCI machine learning repository. <http://archive.ics.uci.edu/ml>. **AQ1**
- 891 Elkan, C. (2001). The foundations of cost-sensitive learning. *Proceedings of the International Joint Confer-*  
892 *ence on Artificial Intelligence*, 17, 973–978.
- 893 Espinosa, M. P., & Gardezabal, J. (2010). Optimal correction for guessing in multiple-choice tests. *Journal*  
894 *of Mathematical psychology*, 54(5), 415–425.
- 895 Ferri, C., Hernandezorrallo, J., & Modroui, R. (2009). An experimental comparison of performance meas-
- 896 ures for classification. *Pattern Recognition Letters*, 30(1), 27–38.
- 897 Gao, W., Wang, L., Jin, R., Zhu, S., & Zhou, Z. (2016). One-pass AUC optimization. *Artificial Intelligence*,  
898 236, 1–29.
- 899 Ghahramani, Z. (2015). Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553),  
900 452–459.
- 901 Goodman, L. A., & Kruskal, W. H. (1963). Measures of association for cross classifications. *Publications of*  
902 *the American Statistical Association*, 49(268), 732–764.
- 903 Hazan, T., Keshet, J., & McAllester, D. A. (2010). Direct loss minimization for structured prediction. In *Pro-*  
904 *ceedings of the advances in neural information processing systems* (pp. 1594–1602).
- 905 He, H., & Garcia, E. (2009). Learning from imbalanced data. *IEEE Transactions on Knowledge and Data*  
906 *Engineering*, 21(9), 1263–1284.
- 907 Hubert, L., & Arabie, P. (1985). Comparing partitions. *Journal of Classification*, 2(1), 193–218.
- 908 Ingo, S. (2005). Consistency of support vector machines and other regularized kernel classifiers. *IEEE*  
909 *Transactions on Information Theory*, 51(1), 128–142.
- 910 Joachims, T. (2005). A support vector method for multivariate performance measures. In *Proceedings of the*  
911 *International Conference on Machine Learning* (pp. 377–384).
- 912 Kotlowski, W., & Dembczynski, K. (2017). Surrogate regret bounds for generalized classification perfor-
- 913 mance metrics. *Machine Learning*, 106(4), 549–572.
- 914 Koyejo, O.O., Natarajan, N., Ravikumar, P.K., & Dhillon, I. S. (2014). Consistent binary classification with  
915 generalized performance metrics. In *Proceedings of the advances in neural information processing*  
916 *systems* (pp. 2744–2752).
- 917 Kuncheva, L. I. (2013). A bound on kappa-error diagrams for analysis of classifier ensembles. *IEEE Trans-*  
918 *actions on Knowledge and Data Engineering*, 25(3), 494–501.
- 919 Li, F., Qian, Y., Wang, J., & Liang, J. (2016). Multigranulation information fusion: A dempster-shafer evi-
- 920 dence theory-based clustering ensemble method. *Information Sciences*, 378(1), 58–63.
- 921 Li, F., Qian, Y., Wang, J., Dang, C., & Liu, B. (2018). Cluster's quality evaluation and selective clustering  
922 ensemble. *ACM Transactions on Knowledge Discovery from Data*, 12(5), 60.
- 923 Li, F., Qian, Y., Wang, J., Dang, C., & Jing, L. (2019). Clustering ensemble based on sample's stability.  
924 *Artificial Intelligence*, 273, 37–55.
- 925 Margineantu, D. D., & Dietterich, T. G. (1997). Pruning adaptive boosting. In *Proceedings of the fourteenth*  
926 *international conference on machine learning* (pp. 211–218).
- 927 Martinezmunoz, G., & Suarez, A. (2006). Pruning in ordered bagging ensembles. In *International confer-*  
928 *ence on machine learning* (pp. 609–616).

- 929 Menon, A. K., Narasimhan, H., Agarwal, S., & Chawla, S. (2013). On the statistical consistency of algo-  
930 rithms for binary classification under class imbalance. In *Proceedings of the international conference*  
931 *on machine learning* (pp. 603–611).
- 932 Musicant, D. R., Kumar, V., & Ozgur, A. (2003). Optimizing F-measure with support vector machines. In  
933 *Proceedings of the Florida AI Research Society* (pp. 356–360).
- 934 Narasimhan, H., & Agarwal, S. (2013). A new support vector method for optimizing partial AUC based on a  
935 tight convex upper bound. In *Proceedings of the conference on knowledge discovery and data mining*.
- 936 Narasimhan, H., Vaish, R., & Agarwal, S. (2014). On the statistical consistency of plug-in classifiers for  
937 non-decomposable performance measures. In *Advances in neural information processing systems* (pp.  
938 1493–1501).
- 939 Narasimhan, H., Ramaswamy, H. G., Saha, A., & Agarwal, S. (2015). Consistent multiclass algorithms for  
940 complex performance measures. In *Proceedings of the international conference on machine learning*  
941 (pp. 2398–2407).
- 942 Qian, Y., Li, F., Liang, J., Liu, B., & Dang, C. (2016). Space structure and clustering of categorical data.  
943 *IEEE Transactions on Neural Networks and Learning Systems*, 27(10), 2047–2059.
- 944 Sabers, D. L., & Feldt, L. S. (1968). An empirical study of the effect of the correction for chance success  
945 on the reliability and validity of an aptitude test. *Journal of Educational Measurement*, 5(3), 251–258.
- 946 Sanyal, A., Kumar, P., Kar, P., Chawla, S., & Sebastiani, F. (2018). Optimizing non-decomposable measures **AQ2**  
947 with deep networks. *Machine Learning* 1–24.
- 948 Scott, W. A. (1955). Reliability of content analysis: The case of nominal scale coding. *Public Opinion*  
949 *Quarterly*, 19(3), 321–325.
- 950 Song, Y., Schwing, A. G., Zemel, R. S., & Urtasun, R. (2016). Training deep neural networks via direct loss  
951 minimization. *Computer Science* 2169–2177.
- 952 Valiant, L. G. (1984). A theory of the learnable. *Communications of ACM*, 27(11), 1134–1142.
- 953 Vapnik, V., & Chervonenkis, A. Y. (1971). On the uniform convergence of relative frequencies of events to  
954 their probabilities. *Theory of Probability and Its Applications*, 16(2), 264–280.
- 955 Veira, S. M., Kaymak, U., & Sousa, J. (2010). Cohen's kappa coefficient as a performance measure for fea-  
956 ture selection. In *Proceedings of the international conference on fuzzy systems* (pp. 1–8).
- 957 Vinh, N.X., Epps, J., & Bailey, J. (2009). Information theoretic measures for clusterings comparison: is a  
958 correction for chance necessary? In *Proceedings of the international conference on machine learning*  
959 (pp. 1073–1080).
- 960 Vinh, N. X., Epps, J., & Bailey, J. (2010). Information theoretic measures for clusterings comparison: Vari-  
961 ants, properties, normalization and correction for chance. *Journal of Machine Learning Research*, 11,  
962 2837–2854.
- 963 Waegeman, W., Dembczyński, K., Jachnik, A., Cheng, W., & Hüllermeier, E. (2014). On the Bayes-optimal-  
964 ity of F-measure maximizers. *Journal of Machine Learning Research*, 15(1), 3333–3388.
- 965 Wu, Q., Laet, T.D., & Janssen, R. (2017). Elimination scoring versus correction for guessing: A simulation  
966 study. In *Proceedings of the meeting of the psychometric society*.
- 967 Zhang, T. (2003). Statistical behavior and consistency of classification methods based on convex risk mini-  
968 mization. *Annals of Statistics*, 32(1), 56–134.
- 969 Zhao, M., Edakunni, N. U., Pocock, A. C., & Brown, G. (2013). Beyond Fano's inequality: Bounds on the  
970 optimal F-score, BER, and cost-sensitive risk and their implications. *Journal of Machine Learning*  
971 *Research*, 14(1), 1033–1090.
- 972 Zhou, Z., Wu, J., & Tang, W. (2002). Ensembling neural networks: Many could be better than all. *Artificial*  
973 *Intelligence*, 137, 239–263.

974 **Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and  
975 institutional affiliations.

976