

Multi-Class Imbalance Problem: A Multi-Objective Solution

Yi-Xiao He^{a,b,1}, Dan-Xuan Liu^{a,b,1}, Shen-Huan Lyu^{c,d,a,2}, Chao Qian^{a,b,1,*},
Zhi-Hua Zhou^{a,b,1}

^a*National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, 210023, China*

^b*School of Artificial Intelligence, Nanjing University, Nanjing, 210023, China*

^c*Key Laboratory of Water Big Data Technology of Ministry of Water Resources, Hohai University, Nanjing, 211100, China*

^d*College of Computer Science and Software Engineering, Hohai University, Nanjing, 211100, China*

Abstract

Multi-class imbalance problems are frequently encountered in real-world applications of machine learning. They have fundamentally complex trade-offs between classes. Existing literature tends to use a predetermined rebalancing strategy and mainly focuses on overall performance measures. However, in many real-world problems, the true level of imbalance and the relative importance between classes are unknown, making it difficult to predetermine the rebalancing strategy and the evaluation criterion. In this paper, we explicitly consider the between-class trade-off issue in the multi-class imbalance problem. We consider all the classes to be important and find a set of optimal trade-offs for the decision-maker to choose from. To reduce the computational cost of this process and make it a practical method, we seek the help of selective ensemble and multiple undersampling rates, and propose the Multi-class Multi-objective Selective Ensemble (MMSE) framework. We further equip the objective modeling with margins to reduce the number of objectives when the task has many classes. Experimental results show that our proposed methods successfully obtain diverse and highly competitive solutions within an acceptable running time.

*Corresponding author.

¹{heyx,liudx,qianc,zhouzh}@lamda.nju.edu.cn

²lvsh@hhu.edu.cn

Keywords: class-imbalanced learning, ensemble method, multi-objective optimization

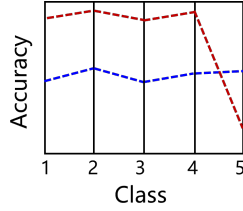
1. Introduction

Class imbalance is a problem frequently encountered in classification tasks [18]. The data collected can be naturally imbalanced, such as the number of patients with different diseases [22]. Abundant imbalanced learning methods have been developed to enhance the relative impact of the minority class in binary classification problems and achieved good results [17, 4, 5, 26]. However, multi-class imbalance problems are fundamentally more complex [36, 25].

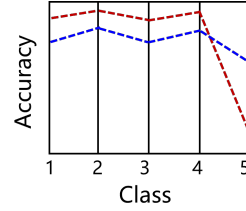
Firstly, in binary classification, even random guesses can achieve an accuracy of 50%, making the problem relatively easy. In contrast, in multi-class cases, vulnerable classes can perform extremely poorly. Secondly, in binary classification, the trade-offs are only between one small class and one big class, while in multi-class imbalance problems, the trade-offs are not only between small and big classes, but also between different small classes and between different big classes. Therefore, designing a rebalancing strategy for multi-class imbalance problems is more challenging. Finally, when it comes to model evaluation, it is hard to describe a multi-class classifier in one overall performance score.

In addition to multi-class classification being more complex than binary classification, another challenge we often face in real-world applications is that the ground-truth level of imbalance and the ground-truth relative importance of the classes are often unknown [48]. Note that under the traditional close-environment assumptions, we always know the targeted performance measure beforehand [49]. Nevertheless, in an open environment, it is not always possible to determine the relative importance of each class *a priori*. If we can provide the decision-maker with all the possible best trade-off performances of the model, it will greatly help her make decisions in an open environment.

Taking disease classification as an example, misdiagnosis of certain rare diseases (classes with a small number of samples) may cause serious problems, but meanwhile it is impossible to quantify the importance of each class. Figure 1 gives two examples of different trade-offs. In each example we assume that there are only two optimal trade-offs, in fact, there may be



(a) Assuming there are only two optimal trade-offs as shown in the figure, the decision-maker chooses the classifier shown in red.



(b) Assuming there are only two optimal trade-offs as shown in the figure, the decision-maker chooses the classifier shown in blue.

Figure 1: Different trade-offs of per-class accuracy. Different optimal trade-offs result in different choices by the decision-maker.

many more trade-offs in real applications. If the only two optimal trade-offs are as shown in Figure 1(a), the decision maker may choose the classifier shown in red because it *can distinguish at least the first four classes*. If the fifth class is indeed important, a separate inspection can be designed. If the only two optimal trade-offs are as shown in Figure 1(b), the decision-maker may choose the classifier shown in blue because it *achieves satisfactory performance on all classes*. The fundamental factor that affects the decision-maker's choice here is that the improvement of the fifth class has different effects on other classes. Only by presenting different optimal trade-offs to the decision-maker can she make better choices.

Therefore, when we cannot determine the importance of each class in advance, we hope to obtain diverse optimal trade-offs among classes for the decision-maker to choose from. To achieve this goal, we propose to model the multi-class imbalance problem as a multi-objective problem

$$\text{maximize } (M_1, M_2, \dots, M_l) , \quad (1)$$

where l denotes the number of classes, M_i is the model's performance on the i -th class. Given that solutions excelling in different objectives are incomparable, multi-objective problems usually have multiple optimal solutions [50, 30, 45]. These optimal trade-off solutions are referred to as *Pareto optimal solutions* (or the *Pareto front* in the objective space). It is assumed that revealing the Pareto front will better equip the decision-maker to make the final choice among these trade-offs.

In the process of searching for multiple optimal solutions on the Pareto front, we need to generate a large number of solutions, each emphasizing

different classes. This process can lead to significant model training overhead. Therefore, reducing this overhead is essential for transforming our goal into a practical learning algorithm. To address this issue, we propose the Multi-class Multi-objective Selective Ensemble (MMSE) framework. It encompasses three fundamental points. 1) We incorporate **selective ensemble** into the multi-objective modeling. In this way, we don't have to repeatedly train the entire model, but instead obtain different models through different combinations of base learners. 2) We use **undersampled** datasets to train base learners, which improves training efficiency. Meanwhile, the model obtained by ensembling multiple base learners can cover more training samples, which avoids the problem of information loss. 3) We undersample the dataset with **different undersampling ratios**. Different undersampling ratios for each class represent different rebalancing strategies. By combining base learners that have heterogeneous emphases over classes, we can obtain a variety of ensemble models with more diverse choices in performance across different classes.

With straightforward objective modeling where the performance of each class is modeled as an objective, we propose $\text{MMSE}_{\text{class}}$. However, scalability is another issue that must be taken into consideration. When the number of classes increases, the optimization problem becomes difficult because most of the generated solutions are incomparable. Considering this, we further propose a margin-based version called $\text{MMSE}_{\text{margin}}$. It optimizes common performance measures by optimizing label-wise and instance-wise margins. It not only reduces the number of objectives to 3 but also proves to be able to optimize common performance measures.

Our contributions are summarized as follows:

- We explore the multi-class imbalance problems from a new perspective, specifically when it is difficult to determine trade-offs between classes *a priori*.
- We model the problem as a multi-objective problem, where the performance of each class is optimized as a separate objective. But more importantly, in order to improve efficiency and make the method practical, we incorporate undersampling and selective ensemble, and develop the MMSE framework.
- Considering the scalability issue when the number of classes increases, we further propose a variant of objective modeling that equips with

margins, and analyze its optimization ability.

- We show in the experiments that both $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ not only achieve better performance on common performance measures, but also provide a variety of trade-offs between classes, and within an acceptable running time.

The rest of this paper is organized as follows. We start by introducing the related work in Section 2. In Section 3, we first demonstrate the problem settings, then introduce the proposed MMSE framework in detail. Theoretical analysis is provided in Section 4. In Section 5, experimental results are reported. Finally, we conclude our work in Section 6.

2. Related work

The most fundamental idea for solving class-imbalanced learning problems is rebalancing. The methods can be roughly categorized into the following three types. a) Sampling methods. These methods include random sampling, synthetic sampling [4, 16], and evolutionary-based sampling methods [11, 33]. b) Re-weighting methods. They are closely related to cost-sensitive learning where instances in small classes have higher misclassification costs [27]. c) Hybrid methods. They combine multiple techniques, such as integrating sampling methods in each boosting round [5, 34], ensembling multiple base learners trained on different balanced training sets [6, 26, 10].

Ensemble methods naturally have applications in solving class imbalance problems, because they can combine the strengths of multiple learners to achieve better performance [43, 42, 44, 39]. A highly representative approach is EasyEnsemble [26]. It combines undersampling with ensemble to achieve effective rebalancing while avoiding information loss. In addition, selective ensemble methods aim to use some base learners to achieve better results than a complete ensemble [19], and can also be applied to handle class imbalance problems [9].

It is worth noting that, many of the imbalanced-learning methods were originally proposed for binary problems, and the binary imbalanced classification has been studied more thoroughly [36, 9]. Although many learning methods are applicable to multi-class imbalance problems, they are generally direct extensions of the binary rebalancing strategies, without considering the complex trade-offs among multiple classes [38, 19]. Usually, a learner

Table 1: Definition of popular multi-class performance measures

Measure	Formulation	Note
Average Accuracy	$\text{Avg. Acc}(h) = \frac{1}{l} \sum_{y=1}^l \frac{1}{ D_y } \sum_{i \in D_y} \mathbb{I}[h(\mathbf{x}_i) = y]$	The average of per-class accuracy.
G-mean	$\text{G-mean}(h) = \left\{ \prod_{y=1}^l \left(\frac{1}{ D_y } \sum_{i \in D_y} \mathbb{I}[h(\mathbf{x}_i) = y] \right) \right\}^{\frac{1}{l}}$	The geometric mean of per-class accuracy.
macro-F1	$\text{macro-F1}(h) = \frac{1}{l} \sum_{y=1}^l \frac{2 \sum_{i \in D_y} \mathbb{I}[h(\mathbf{x}_i) = y]}{ D_y + \sum_{i \in D_y} \mathbb{I}[h(\mathbf{x}_i) = y]}$	F-measure averaging on each class.
micro-F1	$\text{micro-F1}(h) = \frac{2 \sum_{j=1}^l \sum_{i \in D_j} \mathbb{I}[h(\mathbf{x}_i) = y]}{ D + \sum_{j=1}^l \sum_{i \in D_j} \mathbb{I}[h(\mathbf{x}_i) = y]}$	F-measure averaging on the prediction matrix.
macro-AUC	$\text{macro-AUC}(f) = \frac{1}{l} \sum_{y=1}^l \frac{\mathcal{S}_{\text{macro}}^y}{ D_y D \setminus D_y }$ $\mathcal{S}_{\text{macro}}^y = \{(a, b) \in D_y \times \{D \setminus D_y\} \mid f^{(y)}(\mathbf{x}_a) \geq f^{(y)}(\mathbf{x}_b)\}$	AUC averaging on each class. $\mathcal{S}_{\text{macro}}$ is the set of correctly ordered instance pairs considering whether the instance belongs to the corresponding class.
MAUC [14]	$\text{MAUC}(f) = \frac{2}{i(i-1)} \sum_{i < j} \hat{A}(i, j)$ $\hat{A}(i, j) = [\hat{A}(i \mid j) + \hat{A}(j \mid i)]/2$ $\hat{A}(i \mid j) = \frac{1}{ D_i D_j } \{(a, b) \in D_i \times D_j \mid f^{(j)}(\mathbf{x}_a) \geq f^{(j)}(\mathbf{x}_b)\}_{\text{the } i\text{-th class.}}$	AUC averaging on each pair of classes. $\hat{A}(i \mid j)$ is the correctly ordered instance pairs of the i -th and j -th class based on the predicted probabilities on the i -th class.

is trained based on a pre-determined rebalancing strategy, and then the results on a series of evaluation criteria, such as F1, G-mean, and MAUC, are reported [44]. Table 1 summarizes six performance measures commonly used in multi-class imbalance studies. However, few studies have been conducted when the evaluation criteria and the relative importance of classes are unknown beforehand.

In this paper, we consider the performance of different classes as multiple objectives. Recently, many methods have been proposed to optimize multiple objectives simultaneously while training models [41, 40, 23, 24], such as simultaneously optimizing accuracy and regularization, or considering objectives related to specific tasks such as feature selection. However, they did not consider the trade-offs among classes. Instead, we directly model the performance of each class as an objective, and our goal is to provide different trade-offs between classes for the decision-maker to make choices. This is a clear difference that makes this paper a different study from existing literature. Although the idea of modeling each class as an objective is simple, making its optimization practical requires exquisite design, which is the focus of our work.

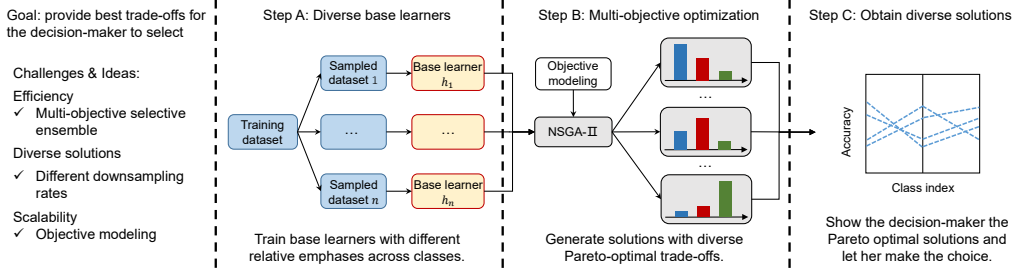


Figure 2: An illustration of our proposed MMSE framework.

3. The proposed approach

3.1. Problem description

Given the multi-class predictor $f : \mathbb{R}^d \rightarrow \mathbb{R}^l$, where $f^{(j)}(\mathbf{x})$ denotes the predicted probability of instance \mathbf{x} on the j -th class. Let $h(\mathbf{x}) = \arg \max_j f^{(j)}(\mathbf{x})$ denote the predicted class. Let D denote a dataset sampled i.i.d. from distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} = \mathbb{R}^d$ is the feature space and $\mathcal{Y} \in \{1, 2, \dots, l\}$ is the label space. Let D_y denote the set of sample indices with label y . $\mathbb{1}_{[\cdot]}$ is the indicator function, which returns 1 if \cdot is true and 0 otherwise.

In this paper, we consider the problem where the decision-maker’s evaluation criterion is not revealed until she sees the best possible trade-off solutions. We consider the following two scenarios of the evaluation process.

Scenario I: After the Pareto front is revealed, the decision-maker decides on a certain overall performance measure. The solution that has the best validation performance on this measure is chosen and the corresponding test performance is reported. We consider the measures in Table 1 to be the possible preferences of the decision-maker.

Scenario II: This scenario covers a broader context, in which the decision-maker may choose any solution on the Pareto front presented to her. Unlike scenario I, the decision-making process here may be a high-level consideration of the trade-offs between classes, which is hard to represent explicitly.

In this paper, we propose a multi-objective selective ensemble method that can deal with the two scenarios simultaneously. Our method not only achieves good performance in common overall performance measures, but also generates diverse trade-off solutions between classes.

169 *3.2. The multi-class multi-objective selective ensemble framework*

170 We present our framework MMSE, as illustrated in Figure 2. It incor-
 171 porates selective ensemble in the multi-objective optimization to enhance
 172 training and storage efficiency, and employs undersampling with different
 173 ratios to help generate diverse solutions.

174 *Multi-objective optimization.* To explicitly consider different trade-offs be-
 175 tween classes, we use the validation accuracy of each class as an objective,
 176 and the multi-objective problem is formulated as

$$\mathbf{g}(h, V) = \left(\frac{1}{|V_1|} \sum_{i \in V_1} \mathbb{1}_{[h(\mathbf{x}_i)=1]}, \dots, \frac{1}{|V_l|} \sum_{i \in V_l} \mathbb{1}_{[h(\mathbf{x}_i)=l]} \right), \quad (2)$$

177 where V denotes the validation set, and V_i denotes the subset of samples
 178 belonging to the i -th class. Usually, the solution to this multi-objective
 179 optimization problem contains many optimal classifiers h , which have their
 180 different advantages in different classes.

181 *Selective ensemble.* Let $F_{\mathbf{s}}$ denote a selective ensemble with selector vector
 182 $\mathbf{s} \in \{0, 1\}^n$, where $s_t = 1$ means that the base learner f_t is incorporated in
 183 the ensemble. If we consider soft voting to combine the base learners, the
 184 predicted probability of ensemble $F_{\mathbf{s}}$ on an instance \mathbf{x} is

$$F_{\mathbf{s}}(\mathbf{x}) = \frac{1}{|\mathbf{s}|} \sum_{t=1}^n s_t f_t(\mathbf{x}),$$

185 where $|\mathbf{s}| = \sum_{t=1}^n s_t$ represents the ensemble size. And let

$$H_{\mathbf{s}}(\mathbf{x}) = \arg \max_j F_{\mathbf{s}}^{(j)}(\mathbf{x}),$$

186 denote the predicted class. In this way, the multi-objective optimization
 187 problem becomes a search on the selector vector, i.e.,

$$\mathbf{g}(\mathbf{s}, V) = \left(\frac{1}{|V_1|} \sum_{i \in V_1} \mathbb{1}_{[H_{\mathbf{s}}(\mathbf{x}_i)=1]}, \dots, \frac{1}{|V_l|} \sum_{i \in V_l} \mathbb{1}_{[H_{\mathbf{s}}(\mathbf{x}_i)=l]}, -|\mathbf{s}| \right). \quad (3)$$

188 Combining selective ensemble with multi-objective optimization leads to
 189 greatly reduced time and storage consumption. Without the design of incor-
 190 porating selective ensemble, to find these solutions using a multi-objective

191 evolutionary algorithm, we have to search many (usually thousands of) rebal-
 192 ancing settings. Since for each setting we have to train a classifier, in total,
 193 we need to train thousands of classifiers from scratch. In contrast, using the
 194 framework we proposed, we only need to search thousands of combinations.

195 *Generating base learners.* When generating the base learners, we construct
 196 multiple undersampled subsets from the training set Tr . Undersampling is
 197 an efficient way to obtain rebalanced datasets with low training overhead.
 198 Compared to it, oversampling has a higher training cost and may also cause
 199 overfitting. The only weakness of undersampling is the possibility of discard-
 200 ing useful samples. But this disadvantage can be compensated for by ensem-
 201 bling multiple undersampled datasets, which avoids information loss [26].
 202 Based on this idea, we use each subset to train a separate classifier, and the
 203 final prediction is made by combining the predictions of all the classifiers.
 204 But there is a novel design in this step of our method, i.e., we undersample
 205 the dataset with *different undersampling ratios*. From EasyEnsemble, we
 206 know that when ensembling base learners trained on balanced subsets, the
 207 ensemble performance will vary depending on the number of base learners.
 208 Obviously, if the sampling ratios on different classes change for different data
 209 subsets, the performance of the obtained ensemble will also exhibit more
 210 diversity. As our goal is to obtain heterogeneous trade-offs among classes,
 211 combining base learners with heterogeneous emphases over classes will help.

212 3.3. Objective modeling for many-class cases

213 In the previous subsection, we use the Eq. (3) version of objective mod-
 214 eling, where the validation accuracy of each class is modeled as objective.
 215 Therefore, we name this method as $MMSE_{\text{class}}$. This type of objective mod-
 216 eling is flexible, and if the optimization problem is well solved, any opti-
 217 mal trade-off between classes can be obtained. However, when the num-
 218 ber of classes is large, the multi-objective problem becomes difficult to op-
 219 timize because most of the generated solutions are incomparable. In such
 220 cases, we propose a margin-based version of objective modeling, and we
 221 name the MMSE method equipped with margin-based objective modeling
 222 as $MMSE_{\text{margin}}$.

223 The concept of margin has long been used in evaluating a model’s train-
 224 ing performance [13], showing its effectiveness in both generalization ability
 225 and robustness. There have been some new research results recently, such

as applying it to multi-label problems [37], or using its distribution to characterize classifier performance more precisely [28]. Inspired by the fact that optimizing label-wise and instance-wise margins can optimize various commonly used multi-label performance measures [37], we decided to optimize the multi-class version of label-wise and instance-wise margins to address our Scenario I. And we apply different methods to aggregate per-class margins so that our method can retain certain advantages in Scenario II. Here we introduce the multi-class version of label-wise and instance-wise margins.

The *label-wise margin* on instance \mathbf{x}_i is defined to be

$$\gamma_i^{\text{label}}(f, \mathbf{x}_i) = \min_{y'} \left\{ f^{(y)}(\mathbf{x}_i) - f^{(y' \neq y)}(\mathbf{x}_i) \right\}, \quad (4)$$

where y is the ground-truth label of instance \mathbf{x}_i . We group the label-wise margin on the instances from the y -th class

$$\bar{\gamma}_y^{\text{label}}(f, V) = \frac{1}{|V_y|} \sum_{i \in V_y} \gamma_i^{\text{label}}(f, \mathbf{x}_i). \quad (5)$$

The *instance-wise margin* on label y is defined to be

$$\gamma_y^{\text{inst}}(f, V) = \min_{a, b} \left\{ f^{(y)}(\mathbf{x}_a) - f^{(y)}(\mathbf{x}_b) \mid a \in V_y, b \in V \setminus V_y \right\}. \quad (6)$$

Instance-wise margin is already defined on each class. But in practice, using the minimum margin of all pairs of instances is not robust, since noise or difficult instances may easily cause a meaningless value of γ_y^{inst} . Therefore we modify Eq. (6) into a more robust mean version

$$\bar{\gamma}_y^{\text{inst}}(f, V) = \left\{ \frac{1}{|V_y|} \sum_{a \in V_y} f^{(y)}(\mathbf{x}_a) - \frac{1}{|V \setminus V_y|} \sum_{b \in V \setminus V_y} f^{(y)}(\mathbf{x}_b) \right\}. \quad (7)$$

The objective vector for $\text{MMSE}_{\text{margin}}$ is defined as

$$\mathbf{g}(\mathbf{s}, V) = \left(\gamma^{\text{label}}(F_{\mathbf{s}}, V), \gamma^{\text{inst}}(F_{\mathbf{s}}, V), -|\mathbf{s}| \right), \quad (8)$$

where

$$\gamma^{\text{label}}(F_{\mathbf{s}}, V) = \frac{1}{l} \sum_y \bar{\gamma}_y^{\text{label}}(F_{\mathbf{s}}, V), \quad (9)$$

$$\gamma^{\text{inst}}(F_{\mathbf{s}}, V) = \min_y \bar{\gamma}_y^{\text{inst}}(F_{\mathbf{s}}, V). \quad (10)$$

Algorithm 1 MMSE

Input: Training data Tr , validation data V , objective modeling \mathbf{g} , evaluation criterion $eval$ denoting the decision making process.

Output: An ensemble.

- 1: Train base learners $\{h_i\}_{i=1}^n$ using different training samples obtained by different sampling strategies.
 - 2: Use NSGA-II to solve the problem $\arg \max_{\mathbf{s}} \mathbf{g}(\mathbf{s}, V)$ obtain a set of Pareto optimal solutions.
 - 3: Present the optimal ensembles to the decision-maker and she selects an ensemble according to $eval$.
-

244 We use the average and minimum for $\bar{\gamma}_y^{\text{label}}$ and $\bar{\gamma}_y^{\text{inst}}$ respectively to em-
245 phasize different aspects of performance across classes. With the objective
246 modeling in Eq. (8), the number of objectives is limited to 3, no matter how
247 many classes there are. Meanwhile, the label-wise and instance-wise margins
248 are related to common performance measures, and the third objective $-|\mathbf{s}|$
249 benefits the theoretical analysis. An analysis for $\text{MMSE}_{\text{margin}}$ is provided in
250 Section 4.

251 The pseudocode of MMSE is shown in Algorithm 1. It applies to both
252 $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$, the only difference is the objective modeling
253 \mathbf{g} . NSGA-II [8] is adopted as the multi-objective optimization algorithm.
254 It is a well-established multi-objective evolutionary algorithm suitable for
255 such combinatorial multi-objective problems. It is suitable for $\text{MMSE}_{\text{margin}}$
256 with only three objectives and can achieve a theoretical guarantee of opti-
257 mization time complexity as will be shown in Section 4. For consistency,
258 we also use NSGA-II for $\text{MMSE}_{\text{class}}$. The evaluation criterion $eval$ denotes
259 the decision-maker’s decision process after obtaining a set of Pareto-optimal
260 solutions. When presenting the obtained solutions to the decision-maker,
261 we can use multi-dimensional data visualization methods, such as parallel
262 coordinates [20, 47], where Figure 1 and Step C in Figure 2 is an example.

263 4. Theoretical analysis

264 4.1. Theoretical results

265 In this section, we prove that $\text{MMSE}_{\text{margin}}$ can optimize common multi-
266 class performance measures with approximation guarantee. Detailed proofs
267 for theorems will be given in Section 4.2.

268 As we have reduced the number of objectives in $\text{MMSE}_{\text{margin}}$, we need
 269 to analyze the expressiveness of the objective modeling. We now show that
 270 if the multi-class version of label-wise margin and instance-wise margins are
 271 optimized, then common multi-class imbalance measures can be optimized.

272 **Proposition 1.** *If all the label-wise margins on dataset D are positive, then*
 273 *Average Accuracy, G-mean, macro-F1, micro-F1 are optimized.*

274 **Proposition 2.** *If all the instance-wise margins on dataset D are positive,*
 275 *then macro-AUC and MAUC are optimized.*

276 Then we analyze the approximation guarantee of $\text{MMSE}_{\text{margin}}$, with NSGA-
 277 II being its multi-objective optimization algorithm. This analysis ensures
 278 that the two objectives of $\text{MMSE}_{\text{margin}}$ can be optimized and have a time
 279 complexity guarantee. Let the selector vector \mathbf{s} represent a subset S of V by
 280 assigning $s_i = 1$ if the i -th base learner of V is in S and $s_i = 0$ otherwise.
 281 Obviously, γ^{label} and γ^{inst} are two set functions that are both non-monotone³
 282 and non-submodular⁴. Therefore, we introduce the ϵ -approximate mono-
 283 tonicity in Definition 1 and β -approximate submodularity in Definition 2 to
 284 characterize how close a set function g is to monotonicity and submodularity,
 285 respectively.

286 **Definition 1** (ϵ -Approximate Monotonicity [21]). Let $\epsilon \geq 0$. A set function
 287 $g : 2^V \rightarrow \mathbb{R}$ is ϵ -approximately monotone if for any $S \subseteq V$ and $v \notin S$.

$$g(S \cup \{v\}) \geq g(S) - \epsilon.$$

288 It is easy to see that g is monotone iff $\epsilon = 0$.

289 **Definition 2** (β -Approximate Submodularity [7]). Let $0 \leq \beta \leq 1$. A set
 290 function $g : 2^V \rightarrow \mathbb{R}$ is β -approximately submodular if for any $S, T \subseteq V$ and
 291 $v \in V$,

$$\sum_{v \in T \setminus S} (g(S \cup \{v\}) - g(S)) \geq \beta(g(S \cup T) - g(S)).$$

292 It is easy to see that g is submodular iff $\beta = 1$.

³A set function $g : 2^n \rightarrow \mathbb{R}$ is monotone if $\forall X \subseteq Y \subseteq V, g(X) \leq g(Y)$.

⁴A set function g is submodular if it satisfies the “diminishing returns” property, i.e.,
 $\forall X \subseteq Y \subseteq V, \sum_{v \in Y \setminus X} (g(X \cup \{v\}) - g(X)) \geq g(X \cup Y) - g(X)$.

293 Assume the solutions in the first nondominated front will not be excluded
 294 from the population by NSGA-II. Let ϵ_1 and β_1 be the approximate mono-
 295 tonicity and approximate submodularity parameter of γ^{label} , respectively, ϵ_2
 296 and β_2 be the approximate monotonicity and approximate submodularity
 297 parameter of γ^{inst} , respectively. Proposition 3 gives the approximation guar-
 298 antee of $\text{MMSE}_{\text{margin}}$ on γ^{label} and γ^{inst} .

299 **Proposition 3.** *For the selective ensemble problem defined in Eq. (8) for*
 300 *$\text{MMSE}_{\text{margin}}$, the expected number of iterations of NSGA-II until finding a*
 301 *solution \mathbf{s} with $|\mathbf{s}| \leq m$ and $\gamma^{\text{label}} \geq (1 - e^{-\beta_1}) \cdot (\text{OPT}^{\text{label}} - m\epsilon_1)$, and a*
 302 *solution \mathbf{t} with $|\mathbf{t}| \leq m$ and $\gamma^{\text{inst}} \geq (1 - e^{-\beta_2}) \cdot (\text{OPT}^{\text{inst}} - m\epsilon_2)$ is $O(n(\log n +$
 303 $m))$, where $\text{OPT}^{\text{label}}$ and OPT^{inst} denote the optimal value of γ^{label} and the
 304 optimal value of γ^{inst} , respectively.*

305 *Proof sketch.* We first prove that with the approximate monotonicity and
 306 approximate submodularity assumption, we can always find an element to
 307 add to a set with certain improvements. Then by tracking the probability
 308 that such an improvement happens on the best solution in the population,
 309 we count the expected number of iterations required by NSGA-II to achieve
 310 the desired approximate guarantee. \square

311 **Remark 1.** As Proposition 3 demonstrates, the multi-objective selective en-
 312 semble procedure of $\text{MMSE}_{\text{margin}}$ can achieve the approximate optimal values
 313 of average label-wise margin γ^{label} and minimum instance-wise margin γ^{inst} .
 314 These two margins are statistics of label-wise margin γ_i^{label} and instance-wise
 315 margin γ_y^{inst} . And from Proposition 1 and Proposition 2 we know that, if γ_i^{label}
 316 and γ_y^{inst} are optimized on all instances and all classes, common multi-class
 317 performance measures are optimized.

318 4.2. Proofs

319 4.2.1. Proof of Proposition 1

Proof. If label-wise margin is positive on an instance \mathbf{x}_i , we have $f^{(y)}(\mathbf{x}_i) >$
 $f^{(y' \neq y)}(\mathbf{x}_i)$. Therefore,

$$\forall \mathbf{x}_i, h(\mathbf{x}_i) = \arg \max_j f^{(j)}(\mathbf{x}_i) = y.$$

320 Then we have $\forall y, \frac{1}{|D_y|} \sum_{i \in D_y} \mathbb{1}_{[h(\mathbf{x}_i)=y]} = 1$. Hence, Avg. $\text{Acc}(h) = 1$, G-mean(h) =
 321 1.

We also have $\sum_{i \in D_y} \mathbb{1}_{[h(\mathbf{x}_i)=y]} = |D_y|$, therefore

$$\text{macro-F1}(h) = \frac{1}{l} \sum_{y=1}^l \frac{2|D_y|}{|D_y| + |D_y|} = 1,$$

$$\text{micro-F1}(h) = \frac{2 \sum_{j=1}^l |D_j|}{|D| + \sum_{j=1}^l |D_j|} = \frac{2|D|}{|D| + |D|} = 1.$$

322

□

323 4.2.2. Proof of Proposition 2

Proof. If instance-wise margin on label y is positive, then

$$f^{(y)}(\mathbf{x}_a) > f^{(y)}(\mathbf{x}_b), \forall a \in D_y, b \in D \setminus D_y.$$

324 Hence,

$$\begin{aligned} \mathcal{S}_{\text{macro}}^y &= \{(a, b) \in D_y \times \{D \setminus D_y\} \mid f^{(y)}(\mathbf{x}_a) \geq f^{(y)}(\mathbf{x}_b)\} \\ &= |D_y| |D \setminus D_y|. \end{aligned}$$

325 If it holds for all y , then

$$\text{macro-AUC}(f) = \frac{1}{l} \sum_{y=1}^l \frac{\mathcal{S}_{\text{macro}}^y}{|D_y| |D \setminus D_y|} = 1.$$

326 We also have

$$\begin{aligned} \hat{A}(i \mid j) &= \frac{1}{|D_i| |D_j|} \{(a, b) \in D_i \times D_j \mid f^{(j)}(\mathbf{x}_a) \geq f^{(j)}(\mathbf{x}_b)\} \\ &= 1, \end{aligned}$$

327 and

$$\hat{A}(i, j) = [\hat{A}(i \mid j) + \hat{A}(j \mid i)]/2 = 1.$$

328 Therefore, $\text{MAUC}(f) = \frac{2}{l(l-1)} \sum_{i < j} \hat{A}(i, j) = 1.$

□

329 *4.2.3. Proof of Proposition 3*

330 *Proof.* The proof relies on Lemma 1 and Lemma 2, which are inspired by
 331 [32]. The detailed proofs of these lemmas are presented later.

332 **Lemma 1.** *Assume that a set function g is ϵ -approximately monotone as in*
 333 *Definition 1 and β -approximately submodular as in Definition 2. For any*
 334 *$\mathbf{s} \in \{0, 1\}^n$ with $|\mathbf{s}| < m$, there exists one element $v \notin \mathbf{s}$ such that*

$$g(\mathbf{s} \cup \{v\}) - g(\mathbf{s}) \geq \beta/m \cdot (\text{OPT} - g(\mathbf{s})) - \beta \cdot \epsilon ,$$

335 *where m is the size constraint.*

336 Assume that the number of selected base learners does not exceed m ,
 337 Lemma 1 proves that for any $\mathbf{s} \in \{0, 1\}^n$ with $|\mathbf{s}| < m$, there always ex-
 338 exists another element, the inclusion of which can bring an improvement on g
 339 roughly proportional to the current distance to the optimum.

340 **Lemma 2.** *To maximize an ϵ -approximately monotone and β -approximately*
 341 *submodular set function g , the expected number of iterations of the NSGA-II*
 342 *until finding a solution \mathbf{s} with $|\mathbf{s}| \leq m$ and $g(\mathbf{s}) \geq (1 - e^{-\beta}) \cdot (\text{OPT} - m\epsilon)$ is*
 343 *$O(n(\log n + m))$, where OPT denotes the optimal value.*

344 Lemma 2 proves the approximation guarantee of NSGA-II on any ϵ -
 345 approximately monotone and β -approximately submodular set function g .
 346 As in previous analyses (e.g., [2, 12]), we may assume that there is a set S_d of
 347 m "dummy" elements whose marginal contribution to any set is 0, i.e., for
 348 any $S \subseteq V$, $g(S) = g(S \setminus S_d)$.

349 By substituting the parameters ϵ_1 and β_1 of γ^{label} as well as ϵ_2 and β_2 of
 350 γ^{inst} into Lemma 2, the theorem can be directly obtained. \square

351 *Proof of Lemma 1.* Let \mathbf{s}^* be an optimal solution containing at most m
 352 items, i.e., $\mathbf{s}^* = \arg \max_{\mathbf{s} \in \{0, 1\}^n, |\mathbf{s}| \leq m} g(\mathbf{s})$, and OPT denote the optimal value,
 353 i.e., $g(\mathbf{s}^*) = \text{OPT}$. We denote the elements in $\mathbf{s} \setminus \mathbf{s}^*$ by $u_1^*, u_2^*, \dots, u_t^*$, where
 354 $t = |\mathbf{s} \setminus \mathbf{s}^*|$. Note that $t < m$ as $|\mathbf{s}| < m$. Because g is ϵ -approximately
 355 monotone, we have

$$\begin{aligned} g(\mathbf{s}^* \cup \mathbf{s}) &= g(\mathbf{s}^* \cup \{u_1^*, u_2^*, \dots, u_t^*\}) \\ &\geq g(\mathbf{s}^* \cup \{u_1^*, u_2^*, \dots, u_{t-1}^*\}) - \epsilon \\ &\geq \dots \geq g(\mathbf{s}^*) - t\epsilon \\ &\geq g(\mathbf{s}^*) - m\epsilon, \end{aligned} \tag{11}$$

where the first three inequalities hold by Definition 1. We denote the elements in $\mathbf{s}^* \setminus \mathbf{s}$ by $v_1^*, v_2^*, \dots, v_l^*$, where $l = |\mathbf{s}^* \setminus \mathbf{s}| \leq m$. Then, we have

$$\begin{aligned}
g(\mathbf{s}^*) - g(\mathbf{s}) - m\epsilon &\leq g(\mathbf{s} \cup \mathbf{s}^*) - g(\mathbf{s}) \\
&= g(\mathbf{s} \cup \{v_1^*, v_2^*, \dots, v_l^*\}) - g(\mathbf{s}) \\
&\leq \frac{1}{\beta} \sum_{j=1}^l (g(\mathbf{s} \cup \{v_j^*\}) - g(\mathbf{s})), \tag{12}
\end{aligned}$$

where the first inequality holds by Eq. (11), the first equality holds by the definition of $\mathbf{s}^* \setminus \mathbf{s}$, and the last inequality holds by Definition 2. Let $v^* = \arg \max_{v \in V/\mathbf{s}} g(\mathbf{s} \cup \{v\})$. Eq. (12) implies that

$$g(\mathbf{s}^*) - g(\mathbf{s}) - m\epsilon \leq l/\beta \cdot (g(\mathbf{s} \cup \{v^*\}) - g(\mathbf{s})).$$

Due to the existence of m dummy elements and $|\mathbf{s}| < m$, there must exist one dummy element $v \notin \mathbf{s}$ satisfying $g(\mathbf{s} \cup \{v\}) - g(\mathbf{s}) = 0$; this implies that $g(\mathbf{s} \cup \{v^*\}) - g(\mathbf{s}) \geq 0$. As $l \leq m$, we have $g(\mathbf{s}^*) - g(\mathbf{s}) - m\epsilon \leq m/\beta \cdot (g(\mathbf{s} \cup \{v^*\}) - g(\mathbf{s}))$, leading to $g(\mathbf{s} \cup \{v^*\}) - g(\mathbf{s}) \geq \beta/m \cdot (\text{OPT} - g(\mathbf{s})) - \beta \cdot \epsilon$. \square

Proof of Lemma 2. We divide the optimization process into two phases: (1) starts from an initial population P with constant size N and finishes after including the special solution $\mathbf{0}$ (i.e., empty set) in population; (2) starts after phase (1) and finishes after finding a solution with the desired approximation guarantee.

For phase (1), we consider the minimum number of 1-bits of the solutions in the population P , denoted by J_{\min} . That is, $J_{\min} = \min\{|\mathbf{s}| \mid \mathbf{s} \in P\}$. Assume that currently $J_{\min} = i > 0$, and let \mathbf{s} be a corresponding solution, i.e., $|\mathbf{s}| = i$. It is easy to see that J_{\min} cannot increase because \mathbf{s} cannot be weakly dominated by a solution with more 1-bits. In each iteration of NSGA-II, to decrease J_{\min} , it is sufficient to select \mathbf{s} and flip only one 1-bit of \mathbf{s} by the bit-wise mutation operator. This is because the newly generated solution \mathbf{s}' now has the smallest number of 1-bits (i.e., $|\mathbf{s}'| = i - 1$) and no solution in P can dominate it; thus it will be included into P . In our setting, the bit-wise mutation is performed with a probability of $1/2$, randomly selecting a parent solution and independently flipping each bit with a probability of $1/n$. Thus, the probability of selecting \mathbf{s} from the population and flipping only one 1-bit of \mathbf{s} by bit-wise mutation is $\frac{1}{2} \cdot \frac{1}{N} \cdot \frac{i}{n} (1 - 1/n)^{n-1} \geq \frac{i}{2enN}$,

383 since the probability of operating bit-wise mutation is $\frac{1}{2}$, the probability of
 384 selecting \mathbf{s} is $\frac{1}{N}$ due to uniform selection and \mathbf{s} has i 1-bits.

385 In each iteration of NSGA-II, there are N offspring solutions to be gener-
 386 ated. Thus, the probability of decreasing J_{min} by at least 1 in each iteration
 387 of NSGA-II is at least $N \cdot \frac{i}{2enN} = \frac{i}{2en}$. Note that $J_{min} \leq n$. We can then
 388 get that the expected number of iterations of phase (1) (i.e., J_{min} reaches $\mathbf{0}$)
 389 is at most $\sum_{i=1}^n \frac{2en}{i} = O(n \log n)$. Note that the solution $\mathbf{0}$ will always be
 390 kept in P once generated, since it has the smallest subset size 0 and no other
 391 solution can weakly dominate it.

392 For phase (2), we consider a quantity J_{max} , which is defined as

$$J_{max} = \max\{j \in \{0, 1, \dots, m\} \mid \exists \mathbf{s} \in P : \\ |\mathbf{s}| \leq j \wedge g(\mathbf{s}) \geq \left(1 - \left(1 - \frac{\beta}{m}\right)^j\right) \cdot (\text{OPT} - m\epsilon)\}.$$

393 That is, J_{max} denotes the maximum value of $j \in \{0, 1, \dots, m\}$ such that in
 394 the population P , there exists a solution \mathbf{s} with $|\mathbf{s}| \leq j$ and $g(\mathbf{s}) \geq (1 - (1 - \beta/m)^j) \cdot (\text{OPT} - m\epsilon)$. The solution that satisfies this condition may not be
 395 unique in the population, but there must be one in the first front. We consider
 396 the solution \mathbf{s} in the first front of NSGA-II. We analyze the expected number
 397 of iterations until $J_{max} = m$, which implies that there exists one solution \mathbf{s}
 398 in P satisfying that $|\mathbf{s}| \leq m$ and $g(\mathbf{s}) \geq (1 - (1 - \beta/m)^m) \cdot (\text{OPT} - m\epsilon) \geq$
 399 $(1 - e^{-\beta}) \cdot (\text{OPT} - m\epsilon)$. That is, the desired approximation guarantee is
 400 reached.
 401

402 The current value of J_{max} is at least 0, since the population P contains
 403 the solution $\mathbf{0}$, which will always be kept in P once generated. Assume that
 404 currently $J_{max} = i < m$. Let \mathbf{s} be a corresponding solution with the value
 405 i , i.e., $|\mathbf{s}| \leq i$ and $g(\mathbf{s}) \geq (1 - (1 - \beta/m)^i) \cdot (\text{OPT} - m\epsilon)$. It is easy to
 406 see that J_{max} cannot decrease because cleaning \mathbf{s} from P implies that \mathbf{s} is
 407 weakly dominated by a newly generated solution $\hat{\mathbf{s}}$, which must satisfy that
 408 $|\hat{\mathbf{s}}| \leq |\mathbf{s}|$ and $g(\hat{\mathbf{s}}) \geq g(\mathbf{s})$. By Lemma 1, we know that flipping one specific
 409 0-bit of \mathbf{s} (i.e., adding a specific element) can generate a new solution \mathbf{s}' ,
 410 which satisfies $g(\mathbf{s}') - g(\mathbf{s}) \geq \frac{\beta}{m}(\text{OPT} - g(\mathbf{s})) - \beta\epsilon$. Then, we have

$$\begin{aligned} g(\mathbf{s}') &\geq \left(1 - \frac{\beta}{m}\right) g(\mathbf{s}) + \frac{\beta}{m} \text{OPT} - \beta\epsilon \\ &\geq \left(1 - \left(1 - \frac{\beta}{m}\right)^{i+1}\right) \cdot (\text{OPT} - m\epsilon), \end{aligned}$$

411 where the last inequality is derived by $g(\mathbf{s}) \geq (1 - (1 - \beta/m)^i) \cdot (\text{OPT} - m\epsilon)$.
 412 After generating \mathbf{s}' , it can be guaranteed that there must be a solution weakly
 413 dominant \mathbf{s}' in the first front, and $J_{max} \geq i + 1$. Thus, J_{max} can increase by
 414 at least 1 in one iteration with probability at least $N \cdot \frac{1}{N} \cdot \frac{1}{2} \cdot \frac{1}{n} (1 - \frac{1}{n})^{n-1} \geq \frac{1}{2en}$,
 415 where $N \cdot \frac{1}{N}$ is the expectation of selecting \mathbf{s} as a parent solution when the
 416 NSGA-II generates N offspring solutions in each iteration, $\frac{1}{2}$ is the probability
 417 of operating bit-wise mutation to the parent solution \mathbf{s} and $\frac{1}{n} (1 - \frac{1}{n})^{n-1}$ is the
 418 probability of flipping a specific bit of \mathbf{s} while keeping other bits unchanged.
 419 This implies that it needs at most $2en$ expected number of iterations to
 420 increase J_{max} . Thus, after at most $2emn = O(mn)$ iterations in expectation,
 421 J_{max} must have reached m .

422 Then, by summing up the expected number of iterations of two phases, we
 423 get that the expected number of iterations of NSGA-II for finding a solution
 424 \mathbf{s} with $|\mathbf{s}| \leq m$ and $g(\mathbf{s}) \geq (1 - e^{-\beta}) \cdot (\text{OPT} - m\epsilon)$ is $O(n(\log n + m))$. \square

425 5. Experiments

426 In this section, we show with experiments that our methods can efficiently
 427 generate many diverse and highly competitive classification models.

428 5.1. Experimental setup

429 5.1.1. Compared methods

430 Considering that our methods employ multiple rebalancing strategies
 431 (specifically, all are forms of undersampling) and decision tree ensembles,
 432 we select compared methods that share these key components. The com-
 433 pared methods must be capable of handling multi-class problems. Unlike
 434 our methods, which offer a wide range of choices for decision-makers, exist-
 435 ing methods can only use predetermined rebalancing strategies and offer only
 436 one solution.

437 We compare our proposed methods $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ to the
 438 following six state-of-the-art ensemble-based multi-class imbalanced learning
 439 methods.

- 440 • SMOTE [4]: It is a synthesized oversampling algorithm. We over-
 441 sample all the other classes to have the same training samples as the
 442 majority class. Then we use multi-class AdaBoost [15] classifier on the
 443 rebalanced dataset.

- 444 • EasyEnsemble [26]: It uses undersampling without replacement to gen-
445 erate multiple balanced training subsets and trains a multi-class Ad-
446 aBoost on each of them, then combines them.
- 447 • BalancedRF [6]: It uses undersampling with replacement to generate
448 multiple balanced subsets first, then trains a decision tree with random
449 feature selection on each of the subsets, then combines them.
- 450 • SMOTEBoost [5]: It adds a step of synthesized oversampling to make
451 a balanced training set in each round of boosting. We extend it to
452 multi-class cases in a way similar to multi-class AdaBoost.
- 453 • MDEP [38]: It is a multi-objective selective ensemble method that
454 simultaneously optimizes validation error, ensemble size, and margin
455 distribution. We use rebalanced base learners as input.
- 456 • DEP [19]: It is a two-stage selective ensemble method that first opti-
457 mizes the combination error and then solely optimizes the validation
458 error. We use rebalanced base learners as input.

459 5.1.2. Datasets

460 We conduct experiments on ten multi-class datasets, including seven LIB-
461 SVM datasets [3], one UCI dataset, and two real-world application datasets.
462 The number of classes varies from 3 (*dna*) to 26 (*letter*). The number of
463 features varies from 6 (*car*) to 2565 (*miRNA*). Table 2 records the number
464 of training instances of each class. In the last column we show the imbalance
465 rate of each dataset, which is calculated by dividing the number of samples
466 in the largest class by the number of samples in the smallest class.

467 Among the benchmark datasets, *car*, *dna* is naturally imbalanced, and
468 *vehicle*, *satimage*, *pendigits*, *usps*, *letter*, *segment* are artificially made imbal-
469 anced.

470 The real-world dataset *acoustic* is naturally imbalanced. The task aims at
471 predicting the function of an acoustic system. The dataset has 21 continuous
472 features, indicating the angle of the placements of 21 acoustic units that
473 determine the function of the system. The four classes are namely *amplify*,
474 *minify*, *cage*, *harvest*. The first two classes mean that the sound will decrease
475 or increase inside the acoustic system. *cage* means there is a sharp decrease in
476 the sound field that the system becomes a cage to shield from the sound [35].
477 *harvest* means the energy is greatly magnified in a small area that it can

Table 2: Information of the datasets

Dataset	Number of training instances in each class	Imbalance rate
car	[307 55 968 52]	18.6
vehicle	[170 140 110 80]	2.1
dna	[507 487 1074]	2.2
satimage	[993 486 956 414 425 809]	2.4
pendigits	[700 600 500 400 300 200 100 70 50 30]	23.3
usps	[800 600 400 400 300 200 100 80 60 50]	16
letter	[520 500 480 460 440 420 400 380 360 340 320 300 280 260 240 220 200 180 160 140 120 100 80 60 40 20]	26
segment	[264 210 160 110 80 50 30]	8.8
acoustic	[2477 723 2674 526]	5.1
miRNA	[2207 256 92 92 92 92 92 92 70 64]	34.5

be captured in the form of electricity [1, 29], meanwhile can be dangerous when the energy focusing is undesired. The extreme cases *cage* and *harvest* naturally happen less often.

The real-world dataset *miRNA* is naturally imbalanced. Circulating microRNAs (miRNAs) are promising biomarkers that could be applied to early detection of cancer. We experimented with data processed from serum miRNA profiles [46], which has 2565 features, each one of which denotes the expression level of certain miRNA⁵. The ten classes are *Healthy*, *Ovarian Cancer*, *Breast Cancer*, *Colorectal Cancer*, *Gastric Cancer*, *Lung Cancer*, *Pancreatic Cancer*, *Sarcoma*, *Esophageal Cancer* and *Hepatocellular Carcinoma*.

5.1.3. Configurations

Experiments were run on a Windows 10 machine with a 3.40 GHz Intel i7-13700KF CPU and 32 GB memory. Each dataset is randomly partitioned into training and test sets, and this partitioning process is repeated 10 times independently and the average result is reported. In the training process of all the methods, the training set is further partitioned into model training set and validation set with ratio 3:1 and with stratified sampling, where the validation set is used for selective ensemble and model selection.

⁵The miRNA data can be downloaded from <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE106817>

497 For the proposed methods $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$, 100 data sub-
 498 sets are generated each randomly using ‘not minority’ or ‘middle’ sampling
 499 strategies, with ‘not minority’ we undersample all the other classes to have
 500 the same training instances as the minority class, and with ‘middle’ we first
 501 randomly select a class, then undersample classes bigger than that class to
 502 have the same number of training instances as that class. Therefore, a class
 503 has different undersampling rates in different subsets. For each data sub-
 504 set, the base learner is randomly chosen from an Adaboost with 10 trees
 505 or a random forest with 5 trees. The population size of NSGA-II is set to
 506 100, and the maximum number of generations is 100. When generating new
 507 solutions, we randomly perform crossover or mutation with probability 0.5
 508 respectively. When doing crossover, we randomly select two parent solutions
 509 uniformly, and randomly select the position of encodings to combine them
 510 into a new solution. When performing mutation, we randomly select a parent
 511 solution and operate a bit-wise mutation that independently flips each bit of
 512 solution with probability $1/n$. Considering the estimation of performance on
 513 the validation set is not accurate, inspired by PONSS [31] that deals with
 514 noisy problems, we use a domination strategy with a threshold.

515 The hyperparameters of the compared methods are selected based on the
 516 performance on the validation set. Specifically, we rank the performance of
 517 each hyperparameter value, and then select the hyperparameter with the best
 518 average rank on the six performance measures. The number of neighbors in
 519 SMOTE is selected from $\{3, 5\}$. The number of base learners in EasyEnsem-
 520 ble is selected from $\{10, 20, 50\}$ and the number of trees in each Adaboost
 521 is set to 10. The number of decision trees in BalancedRF is selected from
 522 $\{10, 20, 50\}$. The maximum number of base learners in SMOTEBoost is se-
 523 lected from $\{10, 20, 50\}$. As one of the objectives of MMSE is to reduce the
 524 size of ensemble, the number of base learners output by MMSE is less than
 525 50. The above settings ensure that the obtained models contain roughly the
 526 same number of individual learners. For MDEP, the individual learners are
 527 the same as MMSE, the population size is 100 and the maximum number
 528 of generations is 100. For DEP, the data subsets are generated the same as
 529 MMSE, the base learners are decision trees. The maximum number of gener-
 530 ations in each stage is 50, and the population size is 100 for both stages. This
 531 setting ensures that the total number of fitness evaluations during MDEP and
 532 DEP is the same as that of MMSE.

533 5.2. Results and discussion

534 We show that our proposed methods are superior in both Scenario I and
535 Scenario II decision-making processes.

536 5.2.1. Scenario I

537 After $\text{MMSE}_{\text{class}}$ or $\text{MMSE}_{\text{margin}}$ obtains a collection of diverse optimal
538 solutions, we examine them with varied performance measures as described
539 in Section 3.1. In detail, we choose the best ensemble on the validation set
540 under each measure and report the corresponding result on the test set. And
541 for the compared methods that each generate a single model only, we simply
542 report the model performance on all six measures.

543 Table 3 and Table 4 show the results on six common performance mea-
544 sures, where the ranking of each method under each performance measure
545 is recorded in the parentheses. From the experimental results, our methods
546 $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ outperform other methods in all evaluation met-
547 rics on almost all the datasets, and obtain very competitive results on the
548 others.

549 Specifically, on *letter* dataset, $\text{MMSE}_{\text{margin}}$ has a better average score than
550 $\text{MMSE}_{\text{class}}$ on all the performance measures. This is because *letter* has 26
551 classes, which is a relatively large number. For $\text{MMSE}_{\text{class}}$, this means the
552 number of objectives is large and the optimization process becomes difficult.
553 At this point, $\text{MMSE}_{\text{margin}}$ is able to perform well because the number of ob-
554 jectives remains unchanged. This demonstrates that the objective modeling
555 in $\text{MMSE}_{\text{margin}}$, which incorporates margin to aggregate the performances of
556 the classes, is proved to be successful. On the other hand, $\text{MMSE}_{\text{class}}$ has its
557 own advantages. For example, on *acoustic* dataset, which has only 4 classes,
558 $\text{MMSE}_{\text{class}}$ outperforms $\text{MMSE}_{\text{margin}}$ on all the measures.

559 To show a summary of the compared methods on all datasets, Figure 3
560 plots the average rank of each method on each performance measure. Accord-
561 ing to the Friedman-Nemenyi test at significance level 0.1, we can observe
562 that 1) Our methods $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ achieve the best average
563 rank on all the performance measures, and they are roughly equally good. 2)
564 MDEP, BalancedRF and SMOTE are significantly worse than our methods
565 on all the performance measures. 3) Compared with DEP, SMOTEBoost,
566 and EasyEnsemble, our methods have no significant advantage, but have bet-
567 ter average rank on all the performance measures. This indicates the high
568 competitiveness of our method on these measures, and in Section 5.2.2, we
569 will further show the richness of the solutions we provide.

Table 3: Experimental results on benchmark datasets of common performance measures. The results are shown in mean \pm std.(rank) of 10 times of running. The smaller the rank, the better the performance. The best accuracy is highlighted in bold type. An entry is marked with a bullet ‘●’ if the method is significantly worse than MMSE_{class} or MMSE_{margin} based on the Wilcoxon rank-sum test with confidence level 0.1.

Dataset	Method	avg. acc	G-mean	F1-macro	F1-micro	macro-AUC	MAUC
car	SMOTE	0.912 \pm 0.018(6)●	0.908 \pm 0.018(6)●	0.909 \pm 0.011(4)●	0.953 \pm 0.013(3)	0.967 \pm 0.017(8)●	0.962 \pm 0.020(8)●
	EasyEnsemble	0.911 \pm 0.017(7)●	0.908 \pm 0.017(7)●	0.797 \pm 0.032(7)●	0.844 \pm 0.019(7)●	0.976 \pm 0.005(6)●	0.989 \pm 0.003(5)●
	BalancedRF	0.918 \pm 0.024(5)●	0.915 \pm 0.024(5)●	0.833 \pm 0.041(6)●	0.873 \pm 0.021(6)●	0.979 \pm 0.005(5)●	0.988 \pm 0.005(6)●
	SMOTEBoost	0.928 \pm 0.038(4)●	0.923 \pm 0.043(4)●	0.939 \pm 0.034(2)	0.968\pm0.017(1)	0.997\pm0.002(1)	0.995 \pm 0.003(4)●
	MDEP	0.884 \pm 0.027(8)●	0.880 \pm 0.029(8)●	0.796 \pm 0.068(8)●	0.843 \pm 0.055(8)●	0.967 \pm 0.014(7)●	0.980 \pm 0.008(7)●
	DEP	0.949 \pm 0.016(3)●	0.948 \pm 0.016(3)●	0.907 \pm 0.022(5)●	0.930 \pm 0.019(5)●	0.994 \pm 0.003(4)●	0.997 \pm 0.002(3)●
	MMSE _{class} (ours)	0.957 \pm 0.023(2)	0.956 \pm 0.023(2)	0.929 \pm 0.030(3)	0.953 \pm 0.016(4)	0.996 \pm 0.003(3)	0.998 \pm 0.002(2)
	MMSE _{margin} (ours)	0.964\pm0.020(1)	0.963\pm0.021(1)	0.945\pm0.024(1)	0.962 \pm 0.016(2)	0.997 \pm 0.002(2)	0.998\pm0.002(1)
vehicle	SMOTE	0.661 \pm 0.041(8)●	0.641 \pm 0.044(8)●	0.666 \pm 0.039(8)●	0.660 \pm 0.041(8)●	0.774 \pm 0.027(8)●	0.774 \pm 0.027(8)●
	EasyEnsemble	0.729 \pm 0.031(5)	0.689 \pm 0.041(6)	0.721 \pm 0.033(5)	0.726 \pm 0.031(5)	0.919 \pm 0.008(2)	0.920\pm0.008(1)
	BalancedRF	0.727 \pm 0.024(6)	0.692 \pm 0.034(5)	0.720 \pm 0.026(6)	0.725 \pm 0.024(6)	0.909 \pm 0.008(6)●	0.910 \pm 0.008(6)●
	SMOTEBoost	0.737 \pm 0.020(2)	0.704\pm0.032(1)	0.732\pm0.023(1)	0.735\pm0.021(1)	0.914 \pm 0.011(5)	0.915 \pm 0.011(5)
	MDEP	0.699 \pm 0.016(7)●	0.662 \pm 0.026(7)●	0.696 \pm 0.019(7)●	0.697 \pm 0.017(7)●	0.895 \pm 0.009(7)●	0.895 \pm 0.009(7)●
	DEP	0.730 \pm 0.021(4)	0.692 \pm 0.033(4)	0.724 \pm 0.023(4)	0.728 \pm 0.022(4)	0.917 \pm 0.007(4)	0.918 \pm 0.007(4)
	MMSE _{class} (ours)	0.738\pm0.032(1)	0.698 \pm 0.041(2)	0.731 \pm 0.031(2)	0.734 \pm 0.030(2)	0.919\pm0.009(1)	0.919 \pm 0.009(2)
	MMSE _{margin} (ours)	0.732 \pm 0.027(3)	0.696 \pm 0.037(3)	0.728 \pm 0.026(3)	0.732 \pm 0.025(3)	0.918 \pm 0.008(3)	0.919 \pm 0.007(3)
dna	SMOTE	0.882 \pm 0.015(8)●	0.881 \pm 0.016(8)●	0.879 \pm 0.015(8)●	0.893 \pm 0.014(8)●	0.918 \pm 0.022(8)●	0.917 \pm 0.022(8)●
	EasyEnsemble	0.938 \pm 0.006(4)●	0.937 \pm 0.006(4)●	0.923 \pm 0.007(6)●	0.928 \pm 0.007(6)●	0.991 \pm 0.002(4)●	0.992 \pm 0.002(4)●
	BalancedRF	0.933 \pm 0.011(6)●	0.932 \pm 0.011(5)●	0.925 \pm 0.011(5)●	0.933 \pm 0.009(5)●	0.989 \pm 0.002(6)●	0.989 \pm 0.002(5)●
	SMOTEBoost	0.933 \pm 0.012(5)●	0.932 \pm 0.012(6)●	0.931 \pm 0.011(4)	0.939 \pm 0.011(3)	0.989 \pm 0.004(6)●	0.989 \pm 0.003(6)●
	MDEP	0.920 \pm 0.008(7)●	0.919 \pm 0.008(7)●	0.906 \pm 0.013(7)●	0.913 \pm 0.013(7)●	0.982 \pm 0.003(7)●	0.982 \pm 0.003(7)●
	DEP	0.942 \pm 0.007(2)	0.942 \pm 0.007(2)	0.932 \pm 0.007(2)	0.938 \pm 0.007(4)	0.992 \pm 0.002(3)	0.992 \pm 0.002(3)
	MMSE _{class} (ours)	0.944\pm0.007(1)	0.943\pm0.007(1)	0.934\pm0.009(1)	0.941\pm0.009(1)	0.993 \pm 0.002(2)	0.993 \pm 0.002(2)
	MMSE _{margin} (ours)	0.940 \pm 0.009(3)	0.940 \pm 0.009(3)	0.932 \pm 0.012(3)	0.939 \pm 0.011(2)	0.993\pm0.001(1)	0.993\pm0.001(1)
satimage	SMOTE	0.825 \pm 0.011(8)●	0.814 \pm 0.015(8)●	0.824 \pm 0.010(8)●	0.847 \pm 0.008(8)●	0.897 \pm 0.006(8)●	0.895 \pm 0.007(8)●
	EasyEnsemble	0.892 \pm 0.010(4)	0.888 \pm 0.011(4)	0.887 \pm 0.009(4)●	0.901 \pm 0.008(5)●	0.988 \pm 0.002(4)	0.987 \pm 0.002(4)
	BalancedRF	0.885 \pm 0.010(5)●	0.880 \pm 0.012(5)●	0.884 \pm 0.009(6)●	0.900 \pm 0.008(6)●	0.986 \pm 0.002(6)●	0.986 \pm 0.002(6)●
	SMOTEBoost	0.878 \pm 0.008(6)●	0.861 \pm 0.011(7)●	0.886 \pm 0.008(5)●	0.909\pm0.007(2)	0.986 \pm 0.002(5)●	0.986 \pm 0.002(5)●
	MDEP	0.876 \pm 0.009(7)●	0.871 \pm 0.009(6)●	0.873 \pm 0.010(7)●	0.890 \pm 0.011(7)●	0.983 \pm 0.004(7)●	0.982 \pm 0.004(7)●
	DEP	0.895 \pm 0.010(2)	0.890 \pm 0.011(3)	0.893 \pm 0.008(2)	0.908 \pm 0.007(3)	0.988 \pm 0.002(3)	0.988 \pm 0.002(3)
	MMSE _{class} (ours)	0.894 \pm 0.011(3)	0.892 \pm 0.012(2)	0.893 \pm 0.011(3)	0.908 \pm 0.009(4)	0.989\pm0.002(1)	0.988\pm0.002(1)
	MMSE _{margin} (ours)	0.899\pm0.012(1)	0.894\pm0.014(1)	0.895\pm0.009(1)	0.909 \pm 0.008(1)	0.988 \pm 0.002(2)	0.988 \pm 0.002(2)
pendigits	SMOTE	0.875 \pm 0.012(8)●	0.864 \pm 0.016(8)●	0.872 \pm 0.013(8)●	0.877 \pm 0.012(8)●	0.931 \pm 0.007(8)●	0.931 \pm 0.007(8)●
	EasyEnsemble	0.949 \pm 0.008(4)●	0.948 \pm 0.008(4)●	0.949 \pm 0.008(4)●	0.949 \pm 0.008(4)●	0.998 \pm 0.001(4)●	0.998 \pm 0.001(4)●
	BalancedRF	0.939 \pm 0.011(5)●	0.937 \pm 0.011(5)●	0.939 \pm 0.011(5)●	0.939 \pm 0.011(5)●	0.997 \pm 0.001(6)●	0.997 \pm 0.001(6)●
	SMOTEBoost	0.931 \pm 0.017(6)●	0.922 \pm 0.024(7)●	0.928 \pm 0.020(7)●	0.932 \pm 0.017(6)●	0.997 \pm 0.002(5)●	0.997 \pm 0.002(5)●
	MDEP	0.930 \pm 0.011(7)●	0.928 \pm 0.011(6)●	0.930 \pm 0.011(6)●	0.931 \pm 0.011(7)●	0.996 \pm 0.002(7)●	0.996 \pm 0.002(7)●
	DEP	0.963\pm0.005(1)	0.962\pm0.006(1)	0.963 \pm 0.005(3)	0.963 \pm 0.005(3)	0.999 \pm 0.000(3)	0.999 \pm 0.000(3)
	MMSE _{class} (ours)	0.959 \pm 0.006(3)	0.958 \pm 0.006(3)	0.963 \pm 0.006(2)	0.964\pm0.006(1)	0.999 \pm 0.000(2)	0.999\pm0.000(1)
	MMSE _{margin} (ours)	0.962 \pm 0.007(2)	0.962 \pm 0.007(2)	0.964\pm0.006(1)	0.963 \pm 0.006(2)	0.999\pm0.000(1)	0.999 \pm 0.000(2)
usps	SMOTE	0.798 \pm 0.012(8)●	0.789 \pm 0.014(8)●	0.801 \pm 0.012(8)●	0.821 \pm 0.009(8)●	0.889 \pm 0.006(8)●	0.888 \pm 0.007(8)●
	EasyEnsemble	0.916 \pm 0.006(4)●	0.916 \pm 0.006(4)●	0.916 \pm 0.005(4)●	0.924 \pm 0.004(4)●	0.995 \pm 0.001(4)●	0.995 \pm 0.001(4)●
	BalancedRF	0.896 \pm 0.009(5)●	0.895 \pm 0.010(5)●	0.897 \pm 0.008(6)●	0.907 \pm 0.007(6)●	0.992 \pm 0.001(6)●	0.991 \pm 0.001(6)●
	SMOTEBoost	0.893 \pm 0.008(6)●	0.887 \pm 0.008(6)●	0.899 \pm 0.007(5)●	0.910 \pm 0.006(5)●	0.993 \pm 0.001(5)●	0.992 \pm 0.001(5)●
	MDEP	0.887 \pm 0.014(7)●	0.884 \pm 0.015(7)●	0.889 \pm 0.015(7)●	0.900 \pm 0.014(7)●	0.989 \pm 0.005(7)●	0.988 \pm 0.005(7)●
	DEP	0.922\pm0.006(1)	0.920\pm0.007(1)	0.924 \pm 0.006(2)	0.931 \pm 0.005(2)	0.996 \pm 0.001(3)	0.995 \pm 0.001(3)
	MMSE _{class} (ours)	0.921 \pm 0.008(2)	0.920 \pm 0.008(2)	0.926\pm0.007(1)	0.934\pm0.006(1)	0.996\pm0.001(1)	0.995 \pm 0.001(2)
	MMSE _{margin} (ours)	0.920 \pm 0.007(3)	0.919 \pm 0.008(3)	0.924 \pm 0.007(3)	0.931 \pm 0.007(3)	0.996 \pm 0.001(2)	0.995\pm0.001(1)

Continued on Table 4

Table 4: Experimental results on benchmark datasets of common performance measures (continued). The results are shown in mean \pm std.(rank) of 10 times of running. The smaller the rank, the better the performance. The best accuracy is highlighted in bold type. An entry is marked with a bullet ‘•’ if the method is significantly worse than MMSE_{class} or MMSE_{margin} based on the Wilcoxon rank-sum test with confidence level 0.1.

Dataset	Method	avg. acc	G-mean	F1-macro	F1-micro	macro-AUC	MAUC
letter	SMOTE	0.769 \pm 0.006(8)•	0.761 \pm 0.007(8)•	0.768 \pm 0.006(8)•	0.769 \pm 0.006(8)•	0.880 \pm 0.003(8)•	0.880 \pm 0.003(8)•
	EasyEnsemble	0.836 \pm 0.007(5)•	0.834 \pm 0.007(5)•	0.838 \pm 0.006(5)•	0.837 \pm 0.007(5)•	0.993 \pm 0.001(4)•	0.993 \pm 0.001(4)•
	BalancedRF	0.804 \pm 0.006(7)•	0.801 \pm 0.007(7)•	0.806 \pm 0.006(7)•	0.805 \pm 0.006(7)•	0.983 \pm 0.001(6)•	0.983 \pm 0.001(6)•
	SMOTEBoost	0.875 \pm 0.006(4)•	0.866 \pm 0.007(4)•	0.873 \pm 0.006(4)•	0.875 \pm 0.006(4)•	0.990 \pm 0.001(5)•	0.990 \pm 0.001(5)•
	MDEP	0.810 \pm 0.058(6)•	0.802 \pm 0.058(6)•	0.809 \pm 0.056(6)•	0.810 \pm 0.057(6)•	0.978 \pm 0.020(7)•	0.978 \pm 0.020(7)•
	DEP	0.892 \pm 0.005(3)	0.885 \pm 0.006(3)	0.891 \pm 0.005(3)•	0.893 \pm 0.005(3)	0.996\pm0.001(1)	0.996\pm0.001(1)
	MMSE _{class} (ours)	0.892 \pm 0.004(2)	0.887 \pm 0.003(2)	0.892 \pm 0.002(2)	0.894 \pm 0.003(2)	0.996 \pm 0.001(3)	0.996 \pm 0.001(3)
	MMSE _{margin} (ours)	0.894\pm0.002(1)	0.888\pm0.002(1)	0.893\pm0.003(1)	0.895\pm0.003(1)	0.996 \pm 0.000(2)	0.996 \pm 0.000(2)
segment	SMOTE	0.934 \pm 0.009(7)•	0.929 \pm 0.010(7)•	0.932 \pm 0.009(7)•	0.934 \pm 0.009(7)•	0.961 \pm 0.005(8)•	0.961 \pm 0.005(8)•
	EasyEnsemble	0.953 \pm 0.009(4)•	0.952 \pm 0.010(4)•	0.953 \pm 0.009(4)•	0.953 \pm 0.009(4)•	0.997 \pm 0.001(4)•	0.997 \pm 0.001(4)•
	BalancedRF	0.931 \pm 0.008(8)•	0.927 \pm 0.009(8)•	0.930 \pm 0.008(8)•	0.931 \pm 0.008(8)•	0.994 \pm 0.002(7)•	0.994 \pm 0.002(7)•
	SMOTEBoost	0.949 \pm 0.007(5)•	0.945 \pm 0.009(5)•	0.948 \pm 0.008(5)•	0.949 \pm 0.007(5)•	0.996 \pm 0.001(5)•	0.996 \pm 0.001(5)•
	MDEP	0.942 \pm 0.012(6)•	0.939 \pm 0.013(6)•	0.942 \pm 0.012(6)•	0.942 \pm 0.012(6)•	0.995 \pm 0.002(6)•	0.995 \pm 0.002(6)•
	DEP	0.959 \pm 0.009(2)	0.957 \pm 0.010(2)	0.959 \pm 0.009(3)	0.959 \pm 0.009(2)	0.997 \pm 0.001(3)	0.997 \pm 0.001(3)
	MMSE _{class} (ours)	0.957 \pm 0.009(3)	0.955 \pm 0.010(3)	0.959 \pm 0.009(2)	0.959 \pm 0.008(3)	0.997 \pm 0.001(2)	0.997 \pm 0.001(2)
	MMSE _{margin} (ours)	0.960\pm0.008(1)	0.958\pm0.009(1)	0.959\pm0.008(1)	0.961\pm0.010(1)	0.998\pm0.001(1)	0.998\pm0.001(1)
acoustic	SMOTE	0.904 \pm 0.007(7)•	0.903 \pm 0.007(7)•	0.890 \pm 0.008(7)•	0.926 \pm 0.006(6)•	0.939 \pm 0.004(8)•	0.936 \pm 0.004(8)•
	EasyEnsemble	0.943 \pm 0.004(5)•	0.942 \pm 0.004(5)•	0.893 \pm 0.005(6)•	0.923 \pm 0.003(7)•	0.996 \pm 0.000(4)•	0.995 \pm 0.001(4)•
	BalancedRF	0.943 \pm 0.006(4)•	0.943 \pm 0.006(4)•	0.914 \pm 0.005(4)•	0.939 \pm 0.004(5)•	0.995 \pm 0.001(5)•	0.995 \pm 0.001(5)•
	SMOTEBoost	0.893 \pm 0.010(8)•	0.889 \pm 0.011(8)•	0.910 \pm 0.009(5)•	0.945 \pm 0.005(4)•	0.995 \pm 0.001(6)•	0.994 \pm 0.001(6)•
	MDEP	0.924 \pm 0.009(6)•	0.924 \pm 0.010(6)•	0.889 \pm 0.011(8)•	0.923 \pm 0.010(8)•	0.990 \pm 0.003(7)•	0.990 \pm 0.003(7)•
	DEP	0.947 \pm 0.007(2)	0.947 \pm 0.007(2)	0.922 \pm 0.007(3)•	0.946 \pm 0.005(3)•	0.996 \pm 0.000(3)•	0.995 \pm 0.000(3)•
	MMSE _{class} (ours)	0.948\pm0.005(1)	0.947\pm0.005(1)	0.932\pm0.005(1)	0.954\pm0.003(1)	0.996\pm0.000(1)	0.996\pm0.000(1)
	MMSE _{margin} (ours)	0.947 \pm 0.003(3)	0.947 \pm 0.003(3)	0.926 \pm 0.008(2)	0.949 \pm 0.006(2)	0.996 \pm 0.000(2)	0.996 \pm 0.000(2)
miRNA	SMOTE	0.583 \pm 0.033(8)•	0.548 \pm 0.043(8)•	0.566 \pm 0.031(8)•	0.836 \pm 0.009(7)•	0.781 \pm 0.017(8)•	0.768 \pm 0.018(8)•
	EasyEnsemble	0.791 \pm 0.025(4)•	0.773 \pm 0.034(3)	0.722 \pm 0.024(4)•	0.876 \pm 0.007(5)•	0.990 \pm 0.002(3)•	0.978 \pm 0.005(2)•
	BalancedRF	0.702 \pm 0.027(5)•	0.672 \pm 0.031(5)•	0.649 \pm 0.028(6)•	0.852 \pm 0.010(6)•	0.974 \pm 0.003(6)•	0.946 \pm 0.006(6)•
	SMOTEBoost	0.658 \pm 0.025(6)•	0.583 \pm 0.039(7)•	0.693 \pm 0.024(5)•	0.901\pm0.007(1)	0.988 \pm 0.002(5)•	0.967 \pm 0.005(5)•
	MDEP	0.638 \pm 0.051(7)•	0.599 \pm 0.056(6)•	0.593 \pm 0.051(7)•	0.834 \pm 0.032(8)•	0.958 \pm 0.019(7)•	0.924 \pm 0.020(7)•
	DEP	0.797 \pm 0.026(3)	0.771 \pm 0.037(4)	0.745 \pm 0.030(2)	0.888 \pm 0.009(3)•	0.989 \pm 0.002(4)•	0.977 \pm 0.005(4)•
	MMSE _{class} (ours)	0.804\pm0.014(1)	0.787\pm0.016(1)	0.747\pm0.025(1)	0.894 \pm 0.007(2)	0.991\pm0.002(1)	0.979\pm0.004(1)
	MMSE _{margin} (ours)	0.799 \pm 0.016(2)	0.781 \pm 0.022(2)	0.740 \pm 0.012(3)	0.888 \pm 0.009(4)	0.990 \pm 0.002(2)	0.977 \pm 0.004(3)

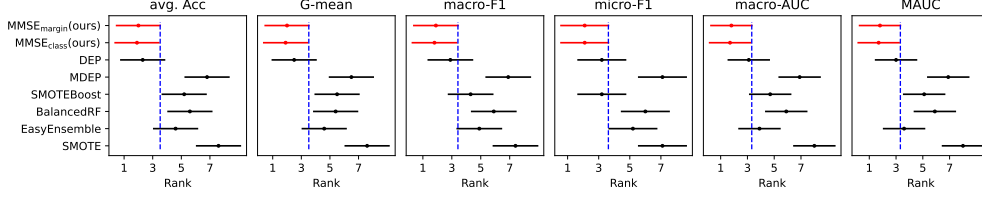


Figure 3: The result of the Friedman-Nemenyi test of the compared methods on different performance measures. The dots indicate the average ranks. The bars indicate the critical difference with the Nemenyi test at significance level 0.1, and two methods with non-overlapping bars are significantly different in performance.

570 In summary, our methods select different solutions based on the decision-
 571 maker’s preferred criterion, and achieve better results than the compared
 572 methods. This quantitatively demonstrates that our method provides highly
 573 competitive choices.

574 5.2.2. Scenario II

575
 576 In Scenario II, the decision-maker may choose any solution on the Pareto
 577 front presented to her. So in order to demonstrate the effectiveness of our
 578 approach, we need to show that we can provide decision-makers with diverse
 579 and good choices.

580 For ease of presentation, we select three out of six compared methods,
 581 namely DEP, EasyEnsemble, and SMOTEBoost. These three methods are
 582 better because they are not significantly inferior to our methods. We com-
 583 pare the solution sets generated by our methods with the single solution
 584 generated by each of the three selected methods separately. We take the
 585 *acoustic* dataset as an example and show the classifiers’ validation accuracy
 586 for each class in Figure 4. The solutions in red dominate the compared classi-
 587 fier, which means they perform better than the compared classifier in all the
 588 classes. The solutions in orange are incomparable with the compared classi-
 589 fier, which means they perform better than the compared classifier in at least
 590 one class. In other words, all solutions of our methods shown in Figure 4
 591 have their advantages. And we can observe that these solutions are also very
 592 diverse. This shows that our method can provide the decision-maker with
 593 rich choices, and these choices are no worse than the best three compared
 594 methods.

595 If we compare the performance of $MMSE_{class}$ and $MMSE_{margin}$ more care-

fully, we can observe that the performance of $\text{MMSE}_{\text{class}}$ is more widely spread in each class, which clearly reflects the waxing and waning relationship between the performance of each class. In contrast, the solution distribution of $\text{MMSE}_{\text{margin}}$ on each class has a relatively consistent trend. This is because $\text{MMSE}_{\text{margin}}$ does not directly optimize the accuracy of each class. But even so, it still provides many different trade-offs.

Figure 5 and Figure 6 show the performance of $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ respectively on the rest datasets. We can see that both $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ obtain diverse and highly competitive solutions on all the datasets.

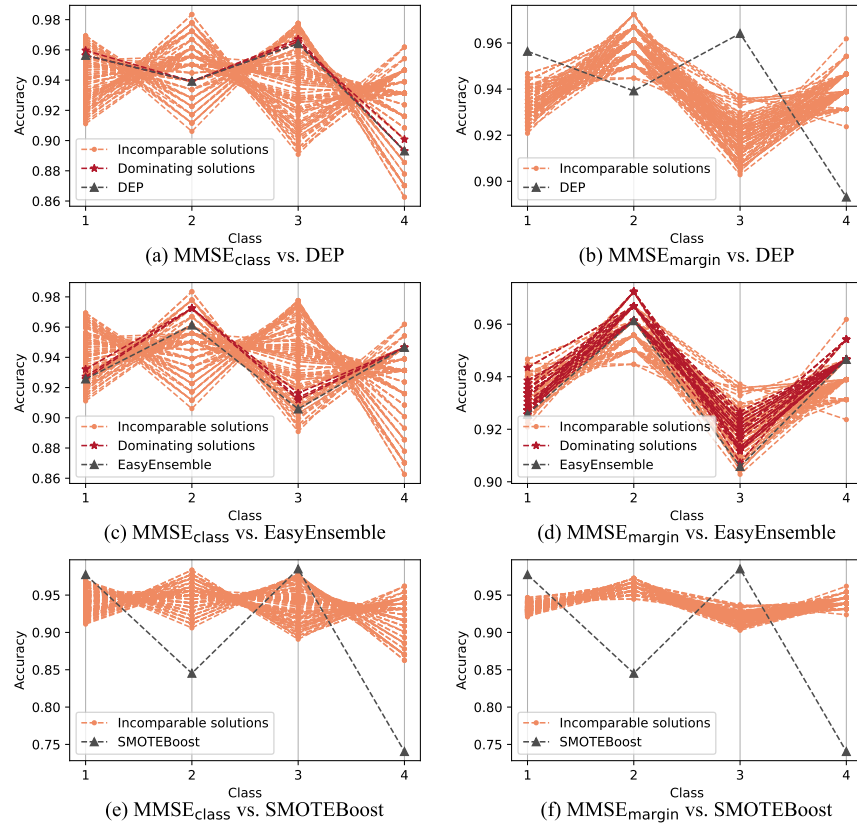


Figure 4: The solutions generated by $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ compared with the single classifier generated by DEP, EasyEnsemble, and SMOTEBoost on *acoustic* dataset. The red solutions dominate the compared classifier, and the orange solutions are incomparable with the compared classifier.

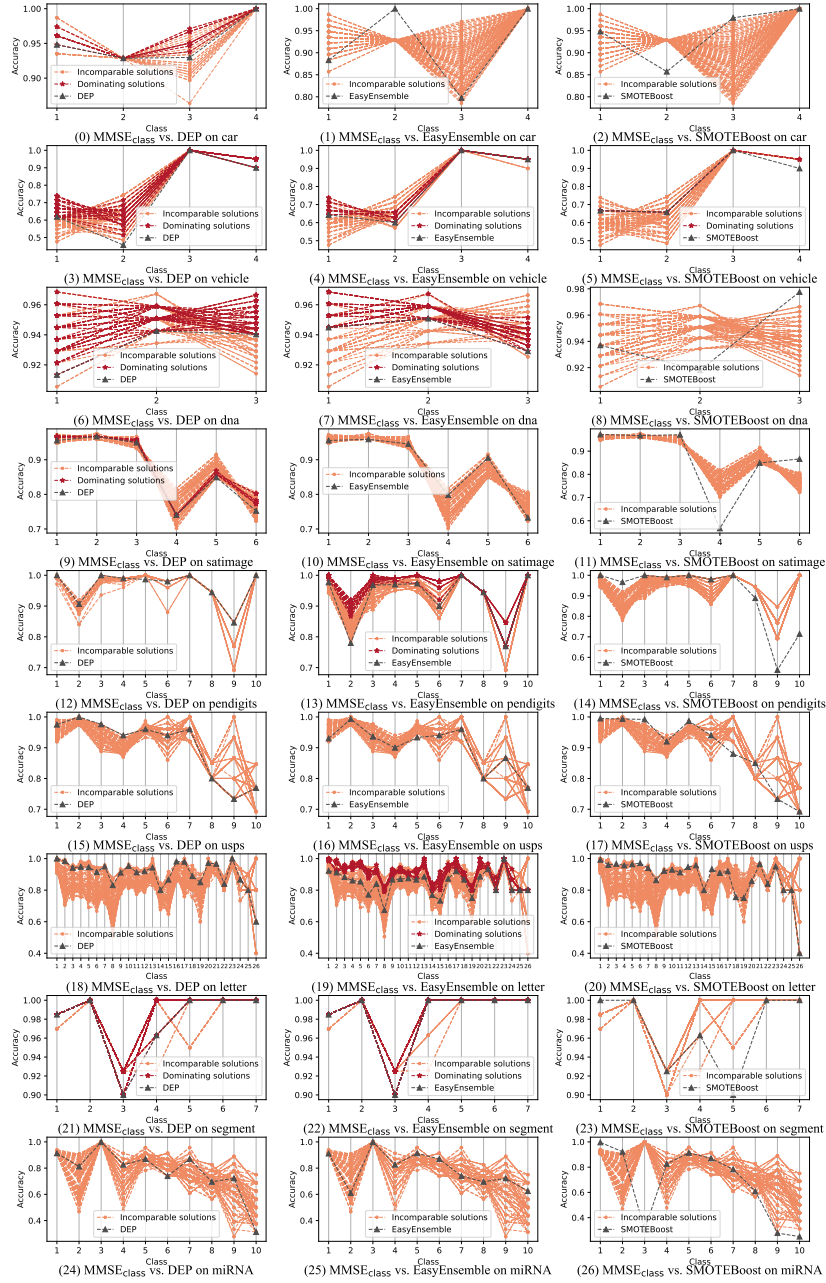


Figure 5: The solutions generated by $MMSE_{class}$ compared with the single classifier generated by DEP, EasyEnsemble, and SMOTEBoost on the other nine datasets. The red solutions dominate the compared classifier, and the orange solutions are incomparable with the compared classifier.

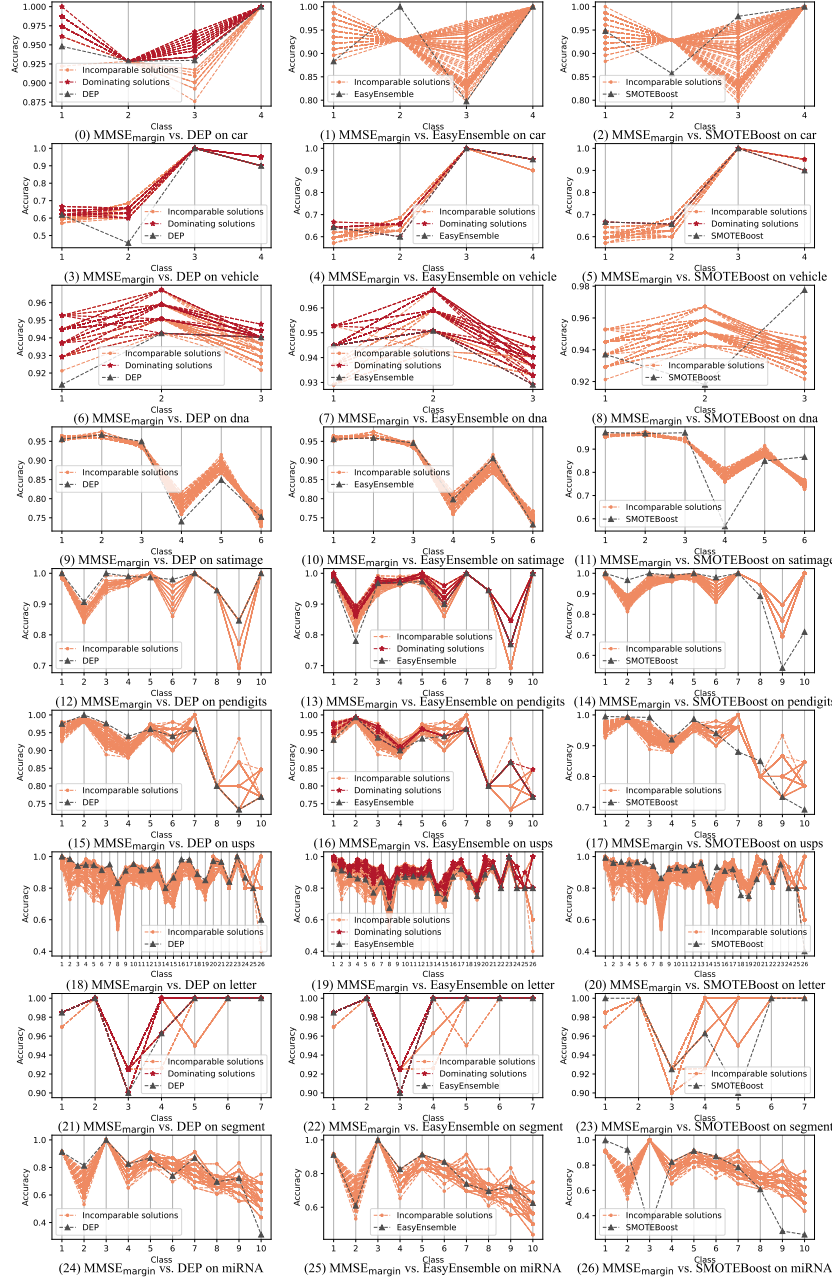


Figure 6: The solutions generated by $MMSE_{margin}$ compared with the single classifier generated by DEP, EasyEnsemble, and SMOTEBoost on the other nine datasets. The red solutions dominate the compared classifier, and the orange solutions are incomparable with the compared classifier.

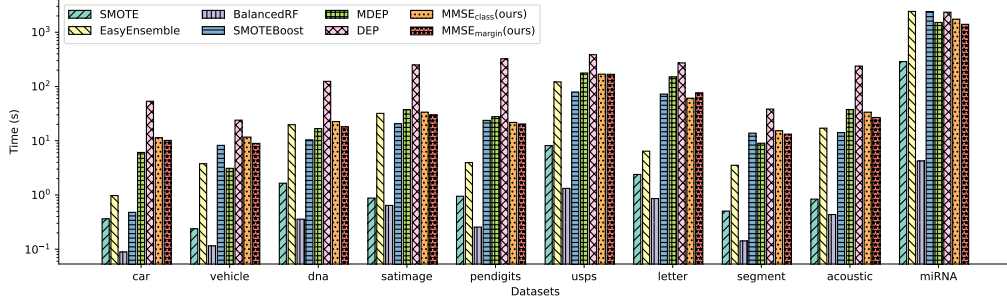


Figure 7: Running time comparison.

5.3. Running time comparison

In this subsection, we compare the running time of different methods. The running time of our methods $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ include training of base learners, multi-objective evolutionary optimization, and the evaluation of the obtained solution set on all the performance measures. Because our methods need to show the decision-maker the performance of all the obtained solutions in all the classes and different evaluation criteria, it is fair to include this part of the time. The running time of the compared methods includes the hyper-parameter tuning and the evaluation of the obtained single model on all the performance measures. As we can observe in Figure 7, the running time of $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ is comparable with EasyEnsemble and SMOTEBoost, the running time of MDEP is roughly the same, while DEP has even longer running time. That is to say, our methods successfully obtain diverse highly competitive solutions efficiently.

5.4. Effectiveness of optimizing margins

$\text{MMSE}_{\text{margin}}$ is a novel design of objective modeling proposed to reduce the number of objectives. In Section 4, we proved that optimizing label-wise margin can optimize *Average Accuracy*, *G-mean*, *macro-F1*, *micro-F1*, and optimizing the instance-wise margin can optimize *macro-AUC* and *MAUC*. Therefore, in this subsection, we experimentally verify it. We choose the *letter* dataset, which has a large number of classes that can best demonstrate the advantages of $\text{MMSE}_{\text{margin}}$. We record the objective values and performance measures of all solutions generated during the multi-objective optimization process. Figure 8 verifies the positive correlation between optimizing the label-wise margin and *Average Accuracy*, *G-mean*, *macro-F1*,

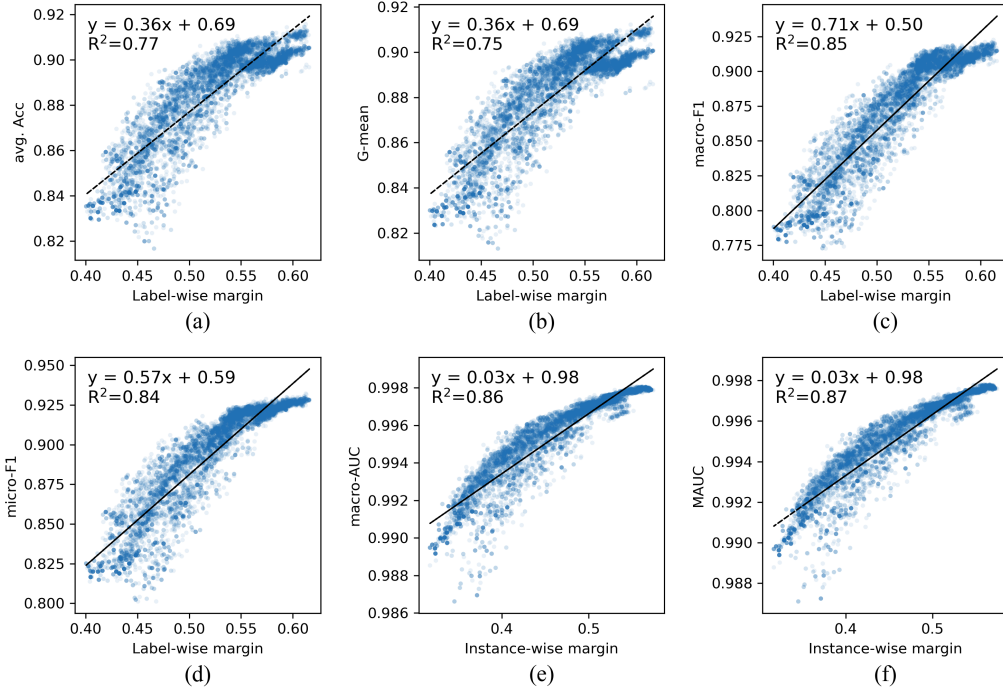


Figure 8: The relationship between the optimization objective and the performance measure that can be optimized in theory. The points are all the solutions generated during the multi-objective evolutionary optimization of applying $MMSE_{margin}$ on the *letter* dataset.

630 *micro-F1*, and the positive correlation between optimizing the instance-wise
631 margin and *macro-AUC* and *MAUC* through two-dimensional scatter plots
632 and the linear fit lines. The slopes of the fitted lines vary greatly because the
633 solutions have different ranges of values on different performance measures,
634 but all slopes are positive, indicating a positive correlation. The key point to
635 note is that the R^2 values in all the subplots are good, as an R^2 value close
636 to 1 indicates a good fit.

637 6. Conclusion

638 In this paper, we revisit the multi-class imbalance problem from the per-
639 spective of multi-objective optimization. Instead of using a predefined re-
640 balancing strategy and generating a single model, we propose the MMSE
641 framework to generate a set of ensembles with the best possible trade-offs

between classes. In real-world applications where it is difficult to choose between different trade-off strategies *a priori*, the decision-maker will be in a better position to make the final choice if the optimal trade-offs are given. Specifically, we propose $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$. The latter enjoys a theoretical guarantee. And experimental results verify that both $\text{MMSE}_{\text{class}}$ and $\text{MMSE}_{\text{margin}}$ can obtain diverse and highly competitive solutions within an acceptable running time.

Currently, we are dealing with class imbalance problems where there is a relative lack of samples in the small classes. An interesting future work is to explore how to use the small class information more effectively when the small class samples are extremely scarce. Another interesting direction for future work is to design specific optimization algorithms for this combinatorial multi-objective optimization problem.

Acknowledgments

This work was supported by the National Science Foundation of China (62276124, 62306104), Jiangsu Science Foundation (BK20230949), China Postdoctoral Science Foundation (2023TQ0104), Jiangsu Excellent Postdoctoral Program (2023ZB140). The authors would like to thank Professor Bin Liang and student Wei Wang for providing an acoustic model that generates the acoustic dataset.

References

- [1] Ahmed, R., Mir, F., Banerjee, S., 2017. A review on energy harvesting approaches for renewable energies from ambient vibrations and acoustic waves using piezoelectricity. *Smart Materials and Structures* 26, 085031.
- [2] Buchbinder, N., Feldman, M., Naor, J., Schwartz, R., 2014. Submodular maximization with cardinality constraints, in: *Proceedings of the 25th Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1433–1452.
- [3] Chang, C.C., Lin, C.J., 2011. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology* 2, 1–27.
- [4] Chawla, N.V., Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. Smote: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research* 16, 321–357.

- 675 [5] Chawla, N.V., Lazarevic, A., Hall, L.O., Bowyer, K.W., 2003. Smote-
676 boost: Improving prediction of the minority class in boosting, in: Pro-
677 ceedings of the 7th European Conference on Principles of Data Mining
678 and Knowledge Discovery, pp. 107–119.
- 679 [6] Chen, C., Liaw, A., Breiman, L., et al., 2004. Using random forest to
680 learn imbalanced data. University of California, Berkeley 110, 24.
- 681 [7] Das, A., Kempe, D., 2011. Submodular meets spectral: Greedy algo-
682 rithms for subset selection, sparse approximation and dictionary selec-
683 tion, in: Proceedings of the 28th International Conference on Machine
684 Learning, pp. 1057–1064.
- 685 [8] Deb, K., Pratap, A., Agarwal, S., Meyarivan, T., 2002. A fast and
686 elitist multiobjective genetic algorithm: NSGA-II. IEEE Transactions
687 on Evolutionary Computation 6, 182–197.
- 688 [9] Du, H., Zhang, Y., Zhang, L., Chen, Y.C., 2023. Selective ensemble
689 learning algorithm for imbalanced dataset. Computer Science and In-
690 formation Systems , 23–23.
- 691 [10] Estabrooks, A., Jo, T., Japkowicz, N., 2004. A multiple resampling
692 method for learning from imbalanced data sets. Computational Intelli-
693 gence 20, 18–36.
- 694 [11] Fernandes, E.R., de Carvalho, A.C., Yao, X., 2019. Ensemble of clas-
695 sifiers based on multiobjective genetic sampling for imbalanced data.
696 IEEE Transactions on Knowledge and Data Engineering 32, 1104–1115.
- 697 [12] Friedrich, T., Göbel, A., Neumann, F., Quinzan, F., Rothenberger, R.,
698 2019. Greedy maximization of functions with bounded curvature un-
699 der partition matroid constraints, in: Proceedings of the 33rd AAAI
700 Conference on Artificial Intelligence, pp. 2272–2279.
- 701 [13] Guo, Y., Zhang, C., 2021. Recent advances in large margin learning.
702 IEEE Transactions on Pattern Analysis and Machine Intelligence 44,
703 7167–7174.
- 704 [14] Hand, D.J., Till, R.J., 2001. A simple generalisation of the area un-
705 der the roc curve for multiple class classification problems. Machine
706 Learning 45, 171–186.

- [15] Hastie, T., Rosset, S., Zhu, J., Zou, H., 2009. Multi-class adaboost. *Statistics and its Interface* 2, 349–360.
- [16] He, H., Bai, Y., Garcia, E.A., Li, S., 2008. Adasyn: Adaptive synthetic sampling approach for imbalanced learning, in: *Proceedings of International Joint Conference on Neural Networks*, pp. 1322–1328.
- [17] He, H., Garcia, E.A., 2009. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering* 21, 1263–1284.
- [18] He, H., Ma, Y., 2013. *Imbalanced Learning: Foundations, Algorithms, and Applications*. John Wiley & Sons.
- [19] He, Y.X., Wu, Y.C., Qian, C., Zhou, Z.H., 2024. Margin distribution and structural diversity guided ensemble pruning. *Machine Learning* doi:10.1007/s10994-023-06429-3.
- [20] Inselberg, A., Dimsdale, B., 1990. Parallel coordinates: a tool for visualizing multi-dimensional geometry, in: *Proceedings of the 1st IEEE conference on visualization*, pp. 361–378.
- [21] Krause, A., Singh, A.P., Guestrin, C., 2008. Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning Research* 9, 235–284.
- [22] Krawczyk, B., Galar, M., Jeleń, Ł., Herrera, F., 2016. Evolutionary undersampling boosting for imbalanced classification of breast cancer malignancy. *Applied Soft Computing* 38, 714–726.
- [23] Li, Y., Zhang, J., Zhang, S., Xiao, W., Zhang, Z., 2022. Multi-objective optimization-based adaptive class-specific cost extreme learning machine for imbalanced classification. *Neurocomputing* 496, 107–120.
- [24] Liang, J., Zhang, Y., Chen, K., Qu, B., Yu, K., Yue, C., Suganthan, P.N., 2024. An evolutionary multiobjective method based on dominance and decomposition for feature selection in classification. *Science China Information Sciences* 67, 120101.
- [25] Liu, X.Y., Li, Q.Q., Zhou, Z.H., 2013. Learning imbalanced multi-class data with optimal dichotomy weights, in: *Proceedings of the 13th IEEE International Conference on Data Mining*, pp. 478–487.

- [26] Liu, X.Y., Wu, J., Zhou, Z.H., 2009. Exploratory undersampling for class-imbalance learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 2, 539–550.
- [27] Liu, X.Y., Zhou, Z.H., 2006. The influence of class imbalance on cost-sensitive learning: An empirical study, in: *Proceedings of the 6th IEEE International Conference on Data Mining*, pp. 970–974.
- [28] Lyu, S.H., Yang, L., Zhou, Z.H., 2019. A refined margin distribution analysis for forest representation learning, in: *Advances in Neural Information Processing Systems* 32, pp. 5531–5541.
- [29] Pillai, M.A., Deenadayalan, E., 2014. A review of acoustic energy harvesting. *International Journal of Precision Engineering and Manufacturing* 15, 949–965.
- [30] Prajapati, A., Parashar, A., Rathee, A., 2023. Multi-dimensional information-driven many-objective software remodularization approach. *Frontiers of Computer Science* 17, 173209.
- [31] Qian, C., Shi, J.C., Yu, Y., Tang, K., Zhou, Z.H., 2017. Subset selection under noise, in: *Advances in Neural Information Processing Systems* 30, pp. 3563–3573.
- [32] Qian, C., Yu, Y., Zhou, Z.H., 2015. Subset selection by pareto optimization, in: *Advances in Neural Information Processing Systems* 28, pp. 1765–1773.
- [33] Roshan, S.E., Asadi, S., 2020. Improvement of bagging performance for classification of imbalanced datasets using evolutionary multi-objective optimization. *Engineering Applications of Artificial Intelligence* 87, 103319.
- [34] Seiffert, C., Khoshgoftaar, T.M., Van Hulse, J., Napolitano, A., 2009. RUSBoost: A hybrid approach to alleviating class imbalance. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* 40, 185–197.
- [35] Shen, C., Xie, Y., Li, J., Cummer, S.A., Jing, Y., 2018. Acoustic metacages for sound shielding with steady air flow. *Journal of Applied Physics* 123, 124501.

- 770 [36] Wang, S., Yao, X., 2012. Multiclass imbalance problems: Analysis and
771 potential solutions. *IEEE Transactions on Systems, Man, and Cyber-*
772 *netics, Part B (Cybernetics)* 42, 1119–1130.
- 773 [37] Wu, X.Z., Zhou, Z.H., 2017. A unified view of multi-label performance
774 measures, in: *Proceedings of the 34th International Conference on Ma-*
775 *chine Learning*, pp. 3780–3788.
- 776 [38] Wu, Y.C., He, Y.X., Qian, C., Zhou, Z.H., 2022. Multi-objective evolu-
777 tionary ensemble pruning guided by margin distribution, in: *Proceedings*
778 *of the 17th International Conference on Parallel Problem Solving from*
779 *Nature*, pp. 427–441.
- 780 [39] Xu, Y., Yu, Z., Chen, C.P., 2024. Classifier ensemble based on multiview
781 optimization for high-dimensional imbalanced data classification. *IEEE*
782 *Transactions on Neural Networks and Learning Systems* 31, 870–883.
- 783 [40] Xue, Y., Cai, X., Neri, F., 2022. A multi-objective evolutionary algo-
784 rithm with interval based initialization and self-adaptive crossover op-
785 erator for large-scale feature selection in classification. *Applied Soft*
786 *Computing* 127, 109420.
- 787 [41] Xue, Y., Tang, Y., Xu, X., Liang, J., Neri, F., 2021. Multi-objective
788 feature selection with missing data in classification. *IEEE Transactions*
789 *on Emerging Topics in Computational Intelligence* 6, 355–364.
- 790 [42] Yang, K., Yu, Z., Chen, C.P., Cao, W., Wong, H.S., You, J., Han, G.,
791 2021. Progressive hybrid classifier ensemble for imbalanced data. *IEEE*
792 *Transactions on Systems, Man, and Cybernetics: Systems* 52, 2464–
793 2478.
- 794 [43] Yang, K., Yu, Z., Chen, C.P., Cao, W., You, J., Wong, H.S., 2022.
795 Incremental weighted ensemble broad learning system for imbalanced
796 data. *IEEE Transactions on Knowledge and Data Engineering* 34, 5809–
797 5824.
- 798 [44] Yang, K., Yu, Z., Wen, X., Cao, W., Chen, C.P., Wong, H.S., You, J.,
799 2020. Hybrid classifier ensemble for imbalanced data. *IEEE transactions*
800 *on neural networks and learning systems* 31, 1387–1400.

- 801 [45] Yang, P., Zhang, L., Liu, H., Li, G., 2024. Reducing idleness in financial
802 cloud services via multi-objective evolutionary reinforcement learning
803 based load balancer. *Science China Information Sciences* 67, 120102.
- 804 [46] Yokoi, A., Matsuzaki, J., Yamamoto, Y., Yoneoka, Y., Takahashi, K.,
805 Shimizu, H., Uehara, T., Ishikawa, M., Ikeda, S.i., Sonoda, T., et al.,
806 2018. Integrated extracellular microRNA profiling for ovarian cancer
807 screening. *Nature Communications* 9, 1–10.
- 808 [47] Zhen, L., Li, M., Peng, D., Yao, X., 2020. Objective reduction for
809 visualising many-objective solution sets. *Information Sciences* 512, 278–
810 294.
- 811 [48] Zhou, Z.H., 2012. *Ensemble Methods: Foundations and Algorithms*.
812 Chapman & Hall/CRC, Boca Raton, FL.
- 813 [49] Zhou, Z.H., 2022. Open-environment machine learning. *National Science*
814 *Review* 9, nwac123. doi:10.1093/nsr/nwac123.
- 815 [50] Zhou, Z.H., Yu, Y., Qian, C., 2019. *Evolutionary Learning: Advances*
816 *in Theories and Algorithms*. Springer, Singapore.