

A post-selection algorithm for improving dynamic ensemble selection methods

Paulo R.G. Cordeiro¹, George D.C. Cavalcanti² and Rafael M.O. Cruz³

Abstract—Dynamic Ensemble Selection (DES) is a Multiple Classifier Systems (MCS) approach that aims to select an ensemble for each query sample during the selection phase. Even with the proposal of several DES approaches, no particular DES technique is the best choice for different problems. Thus, we hypothesize that selecting the best DES approach per query instance can lead to better accuracy. To evaluate this idea, we introduce the Post-Selection Dynamic Ensemble Selection (PS-DES) approach, a post-selection scheme that evaluates ensembles selected by several DES techniques using different metrics. Experimental results show that using accuracy as a metric to select the ensembles, PS-DES performs better than individual DES techniques. PS-DES source code is available in a GitHub repository⁴.

I. INTRODUCTION

Multiple Classifier Systems (MCS) are often used to improve the accuracy and reliability of machine learning models [1]. MCS has three main phases: generation, selection, and combination. In the generation phase, a pool of classifiers is created by training base classifiers with techniques such as Bagging [2], different models, and variations of learning algorithms [3].

The selection phase selects the most competent classifiers in the pool to predict a given query sample. Two main approaches for selecting classifiers are static selection [4] and dynamic selection (DS) [1]. In the static approach, the selection is performed during training and works on the classifiers' overall performance on a validation set. In contrast, the DS approach selects the classifiers on the fly based on their competence in predicting a specific query sample. When only one classifier is selected, it is called Dynamic Classifier Selection (DCS), and when more than one classifier is selected, it is called Dynamic Ensemble Selection (DES). Examples of DES techniques include META-DES [5], Dynamic Ensemble Selection Performance (DESP) [6], and K-Nearest Oracles Union (KNORA-U) [4]. The last phase of an MCS is combination, also called integration. In this phase, the output of all the classifiers selected is combined to produce a final prediction. The combination or integration of the predictions can be done in various ways, including voting and weighting [7].

Currently, research efforts in DES are focused on proposing new methods to improve phases, such as generation [8], selection [9], and combination [10]. Additionally, there have been attempts to apply DES to other areas of knowledge [11]. Despite the progress that has been made, no DES technique is suitable for all problems. This is in line with the statistical rationale for MCS [12], which suggests that combining multiple classifiers increases the likelihood of finding the optimal result for any given problem. However, to the authors' knowledge, the field still lacks techniques that work on evaluating the ensembles selected by DES methods and explores the advantages of pre-selected ensembles to obtain better performance.

Aiming to evaluate this gap in DES' research field, we pose the following research question: "How to analyze ensembles selected by different DES techniques and choose the one having the highest correct prediction potential?" To investigate this question, we propose the Post-Selection Dynamic Ensemble Selection (PS-DES) approach. PS-DES is based on the assumption that different selection criteria may lead to different selected ensembles, and the best criteria used to select an ensemble may differ on an instance-to-instance basis. PS-DES aims to analyze and choose the best ensemble from a set of ensembles generated by various DES techniques to obtain more reliable predictions. Therefore, our proposal works as a post-selection scheme, i.e., it performs after the selection phase of different DES methods and before the combination phase.

Moreover, the best ensemble is selected based on a new concept of ensemble potential proposed in this work. In contrast to the selection criteria employed in many DES methods such as META-DES [5] that work by estimating the quality or competence of each model, the proposed ensemble potential evaluates whether the final selected ensemble of classifiers is reliable. We propose three approaches based on classical performance estimation metrics for measuring the ensemble potential: Accuracy, F-score, and Matthew's Correlation Coefficient.

Experiments over 20 classification datasets and considering three different performance evaluation metrics demonstrate that the post-selection scheme based on the ensemble potential leads to systematic improvement in classification performance over state-of-the-art DES methods. Thus, the evaluation of the pre-selected ensemble capabilities should not be neglected. The rest of the paper is organized as follows: Section II shows a literature review on DES. Section III presents our proposal. Section IV shows the experimental setup. The results are discussed in Section V, and Section VI

¹Paulo R.G. Cordeiro is with Centro de Informática, Universidade Federal de Pernambuco, Recife, Brazil (prgc@cin.ufpe.br)

²George D.C. Cavalcanti is with Centro de Informática, Universidade Federal de Pernambuco, Recife, Brazil (gdcc@cin.ufpe.br)

³Rafael M.O. Cruz is with Département de génie Logiciel et des TI, École de Technologie Supérieure, Montreal, Canada (rafael.menelau-cruz@etsmtl.ca)

⁴<https://github.com/prgc/ps-des>

presents the conclusions.

II. LITERATURE REVIEW

Typically, the development of a DES involves three stages. Firstly, in the generation phase, a pool of classifiers is generated. Secondly, in the selection phase, a subset of classifiers (ensemble) is chosen from the pool created in the generation phase. Finally, the classifiers in the ensemble are combined to classify a given query sample in the combination, or integration, phase.

During DES' generation phase, a pool of classifiers is created, denoted as $P = \{C_1, C_2, \dots, C_M\}$, where M represents the number of classifiers in the pool. The classifiers in the pool must exhibit both diversity and accuracy. Diversity [12] refers to the property that the classifiers should not make the same prediction mistakes, as this is crucial to cover the feature space adequately. Several approaches can be used to generate a pool. These approaches include using different distributions of the training set, such as Bagging [2], using different parameters for the same base classifier (e.g., variations in the number of neighbors in a k-Nearest Neighbors algorithm), or using different base classifiers altogether, which are called heterogeneous ensembles [3]. Heterogeneous ensembles tend to be more diverse than homogeneous ones due to their different mathematical formulations, which typically result in different classification results [13].

The second phase of developing a DES is selection, which aims to choose a subset of classifiers ($P' \subseteq P$), also known as an ensemble. There are two approaches to selection, namely static and dynamic [1]. A fixed subset of classifiers is chosen for all test samples in the static approach. In contrast, the dynamic approach, called Dynamic Ensemble Selection (DES), involves selecting a subset of the pool for each query sample \mathbf{x}_q . In dynamic selection, classifiers are chosen based on some criteria, given the pool created in the previous phase. Among the criteria found in the literature are the Oracle approach, as seen in KNORA-E, KNORA-U [4], and K-Nearest Output Profile (KNOP) [14], accuracy-based methods, such as DES Performance (DES-P) [6], and meta-learning, as in the case of META-DES [5]. These criteria are typically computed from the Region of Competence (RoC), a local region for a query sample (\mathbf{x}_q), denoted as $\theta_{\mathbf{x}_q}$, which is a fundamental concept in dynamic selection approaches. The RoC is usually obtained by applying k-NN or clustering methods to a validation set (DSEL) or the training set itself, such that $\theta_{\mathbf{x}_q} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$, where k is the size of the ROC.

The final phase of a DES is integration, also called aggregation or combination, which involves combining the classifiers selected in the selection phase when multiple classifiers are chosen. Techniques used in this phase include majority vote, product rule, and sum rule [12].

It is worth noting that research papers related to DES do not focus on assessing ensembles that DES techniques have already generated. Elmi and Eftekhari [9] presented a solution by utilizing the selection phase of DES approaches.

However, their proposed approach only allowed for layer-by-layer ensemble selection and did not allow the evaluation of the collective output of all ensembles generated by DES methods.

III. POST-SELECTION DYNAMIC ENSEMBLE SELECTION

The proposed Post-Selection Dynamic Ensemble Selection (PS-DES) is based on the notion of potential, the capability of an ensemble selected by a given DS technique to make a correct prediction. Consequently, it works as a post-processing scheme for ensembles chosen according to different criteria (e.g., meta-learning, Oracle, accuracy). In addition, this proposal aims to evaluate the quality or potential of a selected ensemble, which contrasts with the current DES methods that build an ensemble by selecting multiple competent classifiers individually without trying to characterize the selected dynamic ensemble. This approach consists of three phases: **(1)** pool generation and setup, **(2)** post-selection, and **(3)** combination.

A. Phase 1: Pool generation and DES' setup

The initial stage of PS-DES involves generating a pool of classifiers and configuring the DES techniques. First, Bagging generates multiple bootstraps (T^b) from the original dataset (T), where b is the number of bootstraps. Then, a pool of $b \times m$ classifiers, denoted as $P = \{C_1^1, C_1^2, \dots, C_m^b\}$, is constructed by training each of the m classifiers (C_1, \dots, C_m) on each of the b bootstraps generated by Bagging. Finally, any DES techniques specified by the user, including META-DES, KNORA-U, and DES-P, are initialized using the same pool P as input. These techniques can be used to select an intermediate ensemble later validated by our post-processing scheme to obtain the optimal one. All DES approaches, $DES_{set} = \{des_1, des_2, \dots, des_n\}$, are consolidated into the set DES_{set} .

B. Phase 2: Selection

This phase seeks to identify the optimal ensemble from a set of ensembles generated by different DES techniques, and Algorithm 1 shows its pseudo-code. Given a query sample (\mathbf{x}_q), the validation dataset $DSEL$, and a set of DES techniques (DES_{set}), this phase involves selecting several dynamic ensembles P' , each one created using a different DES method, and assessing their effectiveness in order to determine which ensemble is most likely to perform well for the given \mathbf{x}_q .

This phase begins by computing the Region on Competence (RoC) ($\theta_{\mathbf{x}_q}$) for the query sample \mathbf{x}_q using the k-NN algorithm. It is essential to note that all DES_n utilize the same RoC ($\theta_{\mathbf{x}_q}$) based on the k-NN, thereby reducing the computational burden of implementing multiple selection criteria.

Then, for each technique in DES_{set} , a set of classifiers is selected according to its competence estimation and selection criterion, forming the ensemble P' (lines 5 and 6). Subsequently, the potential of the generated ensemble P' is evaluated (line 7). As the class label of \mathbf{x}_q is unknown, the

Algorithm 1 PS-DES selection

```
1: procedure PS-DES_SELECTION( $\mathbf{x}_q, DSEL, DES_{set}$ )
2:    $pot_{max} \leftarrow 0$ 
3:    $P^{sel} \leftarrow \emptyset$ 
4:    $\theta_{\mathbf{x}_q} \leftarrow \text{calculate\_ROC}(\mathbf{x}_q, DSEL)$ 
5:   for  $des$  in  $DES_{set}$  do
6:      $P' \leftarrow \text{get\_ensemble}(des, \theta_{\mathbf{x}_q})$ 
7:      $pot_{des} \leftarrow \text{calculate\_pot}(P')$ 
8:     if  $pot_{des} \geq pot_{max}$  then
9:        $pot_{max} \leftarrow pot_{des}$ 
10:       $P^{sel} \leftarrow P'$ 
11:     end if
12:   end for
13:   return  $P^{sel}$ 
14: end procedure
```

potential assumes that the output class of \mathbf{x}_q corresponds to the majority vote of the ensemble. Consequently, it computes the potential of this ensemble by assessing the proportion of methods in it that contribute to this decision. For instance, given a binary classification problem and an ensemble with seven base classifiers, $P' = \{C_1, \dots, C_7\}$ selected by a given DES technique, and let $y_{P'} = \{0, 1, 0, 0, 1, 1, 1\}$ be the predictions of the base classifier for the given \mathbf{x}_q . The majority vote would give the class 1 as the answer. The potential is then estimated based on a classical performance metric comparing the ensemble majority vote and the votes of each classifier using a performance metric. If accuracy is used to calculate P' potential, the value would be $pot_{des} = 0.57$. If the F-score is chosen as the potential metric, the P' potential is $pot_{des} = 0.73$.

After evaluating the potential of all possible ensembles, the one that obtained the highest value, P^{sel} , is returned as the selected one for the combination step.

C. Phase 3: Combination

Once the P^{sel} selection is complete, Phase 3 begins, which is accountable for combining the classifiers into P^{sel} , using techniques such as majority vote or sum rule.

IV. EXPERIMENTAL SETUP

Datasets. The experiments were conducted using 20 datasets from the UCI Machine Learning Repository [15], which vary in sample size, dimensions, number of classes, and Imbalance Ratio (IR) (Table I). Each dataset T is split into three parts: training (50%), $DSEL$ (25%), and testing (25%). This split is stratified, meaning that the proportions of the classes between the three datasets are maintained. For each dataset, we run 30 replications, changing the distribution of the sets (holdout) to obtain the average values for the evaluated metrics. The data is scaled using the Standard Scaler (also known as Z-score normalization [16]).

Phase 1. First, Bagging, 100 bootstraps were used for all experiments, consistent with previous studies [1], [4], [5]. For the pool generation, three base classifiers (Perceptron, Logistic Regression, and Naive Bayes) were selected for

TABLE I

DATASETS MAIN CHARACTERISTICS. THE NUMBER OF SAMPLES, DIMENSIONS (DIM), CLASSES, AND IMBALANCE RATIO (IR).

Datasets	Examples	Dim	Classes	IR
appendicitis	106	7	2	2.52
australian	690	14	2	1.12
balance	625	4	3	2.63
cmc	1473	9	3	1.30
column_3C	310	6	3	2.50
diabetes	768	8	2	1.86
glass1	214	9	2	1.82
glass6	214	9	2	6.38
haberman	306	3	2	2.78
hayes	160	4	3	3.40
heart	270	13	2	1.25
led7digit	500	7	10	1.54
mammographic	830	5	2	1.15
musk	476	166	2	1.29
pima	768	8	2	1.90
sonar	208	60	2	1.14
vehicle	846	18	4	1.10
vehicle2	846	18	2	2.88
vowel	990	13	11	1.00
wdbc	683	9	2	1.85

the experiments. As they have different mathematical foundations and low computational costs [1], [10], [17] they are suitable for building a diverse and lightweight pool of classifiers. Thus, the classifier pool (P) consisted of 300 classifiers (3 base classifiers \times 100 bootstraps). Since the focus of the research was not on optimizing each base model's hyperparameters, the default hyperparameters values from scikit-learn were used.

Four DES approaches (KNORA-U, KNOP, DES-P, and META-DES) were chosen due to their application of various selection criteria (e.g., Oracle, accuracy, meta-learning). These approaches showed superior performance in a recent empirical study [1]. We applied these DES methods default hyperparameter configurations of the DESlib 0.3 library [18] to guarantee experiment consistency. Moreover, the same pool of classifiers was utilized to fairly compare all DES techniques.

Phase 2. The Region of Competence (RoC) was calculated applying k-Nearest Neighbors (k-NN) with $k = 7$, as suggested in [1]. To assess the performance of the ensembles, we employed a range of evaluation metrics, including accuracy, F-score, and Matthews Correlation Coefficient (MCC). Accuracy is a popular metric for DES techniques, although it may not be suitable for imbalanced datasets (i.e., high IR). Meanwhile, F-score and MCC are more suitable for such datasets. F-score is advantageous in scenarios where there is an appreciation for recall and precision, since these two metrics are used in its calculation. The MCC considers false-negative rates in its formulation, which can be of interest to specific problems.

The PS-DES variants are labeled according to the metric used to calculate the potential: accuracy (PS-DES-acc), F-score (PS-DES-F), and Matthews Correlation Coefficient (PS-DES-MCC). To assess whether the proposed metrics perform better than random selection, we also conducted an

experiment that randomly selected the best ensemble (PS-DES-Random).

Phase 3. Finally, majority voting was used as a combination approach since individual DES techniques usually apply it [1].

V. RESULTS AND DISCUSSIONS

The proposed method is evaluated based on three metrics: accuracy (Table II), F-score (Table III), and MCC (Table IV). Upon examining the tables, our results indicate that PS-DES-acc outperforms all the other approaches in all metrics. The PS-DES-acc obtained the best rank considering all performance metrics, followed by the variant using the F-score metric for computing the ensemble’s potential. These results are interesting since, even though the final proposal may be evaluated regarding a different performance metric (e.g., F-score or MCC), using accuracy as the metric for computing the ensemble potential is more advantageous.

Analogously, MCC obtained the lowest ranking among all PS-DES variants even when in the scenario that MCC is used as a performance evaluation metric to compute the overall method performance. This result indicates no relation between the metrics selected for calculating ensemble potential and the same metric applied to evaluate the approaches. For accuracy, F-score, and MCC, the chosen metric for calculating the potential does not interfere with the approach’s evaluation. Nevertheless, according to these tables, the average ranking of PS-DES approaches is systematically better when compared to individual DS techniques (e.g., META-DES). Thus, the proposed post-processing selection scheme indeed leads to more robust dynamic ensemble selection systems.

However, to see if such a difference in performance is significant, we need to go further and perform a more fine-grained analysis by comparing pair of techniques over multiple datasets. Hence, we also conducted one analysis based on the number of wins, ties, and losses (w/t/l) obtained by a control technique and the Wilcoxon signed rank test with a confidence level of 95%. Results of these pairwise comparisons are presented in Tables V, VI, and VII for the PS-DES-acc, PS-DES-F and PS-DES-MCC methods, respectively.

The pairwise statistical analysis of PS-DES-acc shows it outperforms KNORA-U, KNOP, META-DES, Random, and PS-DES-MCC regarding accuracy (Table V). No significant difference is observed between PS-DES-acc and DES-P or PS-DES-F. However, considering the presence of datasets with $IR > 1$, it is necessary to consider F-score and MCC. The F-score analysis reveals that PS-DES-acc outperforms all DES individual techniques and Random, with no significant difference to PS-DES-F and PS-DES-MCC. For MCC, PS-DES-acc performs exceptionally well and obtains significantly better results compared to all techniques apart from DES-P and PS-DES-F. Ultimately, this variant based on accuracy for computing the ensemble potential obtained more victories than all other models, regardless of the performance metric used.

The statistical analysis of PS-DES-F (Table VI) indicates that it performs better than KNORA-U, KNOP, and Random on all three metrics. However, no statistical difference is found for MCC when compared with DES-P. However, the win-tie-loss analysis demonstrates that the PS-DES-F systematically obtained more wins against the state-of-the-art DES techniques and the random selection scheme (between 13 to 15 wins over the 20 datasets). In contrast, the analysis of PS-DES-MCC (Table VII) presents the worst results compared to PS-DES-acc and PS-DES-F. Based on Wilcoxon’s test analysis, PS-DES-MCC scores better than KNORA-U and KNOP only in F-score. The hypothesis that PS-DES-MCC scores better cannot be refuted for all other metrics and comparisons.

In summary, the results indicate that PS-DES-acc and PS-DES-F yield comparable outcomes. Still, PS-DES-acc holds a slight advantage over its competitor, particularly when it is compared against the state-of-the-art DES methods.

VI. CONCLUSION

This work proposed a new Dynamic Ensemble Selection (DES) method: Post-Selection Dynamic Ensemble Selection (PS-DES). This method is based on the idea that the optimal criteria for ensemble selection may differ at the instance level leading to ensembles with different qualities or “potentials”. To this end, the approach evaluates the potential of ensembles chosen by various DES techniques to determine which is more suitable for labeling a given instance.

Experiments demonstrate no direct correlation between the metrics applied for calculating the ensemble potential and for evaluating the approaches, as the PS-DES-acc was found to achieve the best overall results in all cases. Additionally, PS-DES was consistently superior to the existing state-of-the-art DES techniques, which implies that evaluating the selected ensembles as a collective is more important than assessing and choosing each base classifier separately, as is the trend in most DES methods. Thus, post-processing approaches in DES are vital, and future works will explore new metrics for measuring the ensemble’s potential.

ACKNOWLEDGMENT

The authors would like to thank the Instituto Federal de Pernambuco, Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), Fundação de Amparo à Ciência e Tecnologia de Pernambuco (FACEPE), and the Natural Sciences and Engineering Research Council of Canada (NSERC).

REFERENCES

- [1] R. M. Cruz, R. Sabourin, and G. D. Cavalcanti, “Dynamic classifier selection: Recent advances and perspectives,” *Information Fusion*, vol. 41, pp. 195–216, 2018.
- [2] L. Breiman, “Bagging predictors,” *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [3] G. Tsoumakas, L. Angelis, and I. Vlahavas, “Selective fusion of heterogeneous classifiers,” *Intelligent Data Analysis*, vol. 9, no. 6, pp. 511–525, 2005.
- [4] A. H. Ko, R. Sabourin, and A. S. Brito Jr, “From dynamic classifier selection to dynamic ensemble selection,” *Pattern recognition*, vol. 41, no. 5, pp. 1718–1731, 2008.

TABLE II

AVERAGE ACCURACY AND RANKINGS OF THE EVALUATED METHODS FOR EACH DATASET. THE BEST RESULT PER DATASET IS PRESENTED IN BOLD.

Datasets	KNORA-U	KNOP	META-DES	DES-P	Random	PS-DES-MCC	PS-DES-F	PS-DES-acc
appendicitis	0.859	0.862	0.851	0.864	0.855	0.859	0.868	0.869
australian	0.847	0.846	0.848	0.846	0.846	0.849	0.849	0.849
balance	0.883	0.886	0.893	0.883	0.885	0.893	0.891	0.889
cmc	0.509	0.512	0.488	0.513	0.506	0.483	0.505	0.510
column_3C	0.845	0.844	0.841	0.848	0.846	0.839	0.845	0.847
diabetes	0.768	0.769	0.762	0.765	0.767	0.762	0.767	0.770
glass1	0.642	0.640	0.674	0.682	0.653	0.682	0.672	0.672
glass6	0.940	0.938	0.947	0.944	0.943	0.948	0.948	0.947
haberman	0.731	0.731	0.726	0.737	0.734	0.727	0.737	0.734
hayes	0.611	0.619	0.670	0.640	0.640	0.665	0.684	0.682
heart	0.836	0.838	0.833	0.841	0.836	0.837	0.838	0.837
led7digit	0.725	0.726	0.708	0.723	0.721	0.697	0.704	0.723
mammographic	0.830	0.831	0.822	0.828	0.830	0.827	0.824	0.825
musk	0.779	0.781	0.788	0.792	0.786	0.790	0.791	0.794
pima	0.770	0.768	0.761	0.766	0.766	0.762	0.763	0.766
sonar	0.771	0.771	0.796	0.784	0.787	0.787	0.796	0.798
vehicle	0.750	0.752	0.761	0.752	0.750	0.755	0.758	0.761
vehicle2	0.949	0.947	0.956	0.951	0.951	0.955	0.954	0.955
vowel	0.965	0.965	0.981	0.967	0.969	0.980	0.978	0.977
wdbc	0.970	0.970	0.969	0.970	0.970	0.968	0.969	0.968
ranking	5.45	5.03	4.90	3.88	5.23	4.73	3.70	3.10

TABLE III

AVERAGE F-SCORE AND RANKINGS OF THE EVALUATED METHODS FOR EACH DATASET. THE BEST RESULT PER DATASET IS PRESENTED IN BOLD.

Datasets	KNORA-U	KNOP	META-DES	DES-P	Random	PS-DES-MCC	PS-DES-F	PS-DES-acc
appendicitis	0.708	0.718	0.717	0.749	0.715	0.732	0.760	0.754
australian	0.845	0.843	0.845	0.844	0.843	0.846	0.847	0.846
balance	0.612	0.615	0.631	0.613	0.615	0.637	0.628	0.626
cmc	0.472	0.477	0.454	0.484	0.472	0.451	0.474	0.478
column_3C	0.799	0.796	0.792	0.803	0.800	0.790	0.798	0.801
diabetes	0.725	0.725	0.719	0.723	0.724	0.719	0.726	0.729
glass1	0.422	0.415	0.585	0.577	0.512	0.584	0.530	0.535
glass6	0.850	0.847	0.880	0.868	0.863	0.879	0.878	0.873
haberman	0.506	0.507	0.525	0.542	0.519	0.520	0.555	0.549
hayes	0.627	0.634	0.688	0.658	0.654	0.681	0.697	0.699
heart	0.832	0.834	0.829	0.837	0.832	0.833	0.834	0.833
led7digit	0.718	0.719	0.700	0.715	0.714	0.687	0.693	0.719
mammographic	0.829	0.830	0.821	0.826	0.829	0.825	0.822	0.824
musk	0.770	0.772	0.781	0.785	0.778	0.783	0.784	0.787
pima	0.727	0.724	0.718	0.724	0.722	0.718	0.721	0.725
sonar	0.768	0.767	0.792	0.780	0.783	0.784	0.793	0.795
vehicle	0.742	0.745	0.757	0.744	0.743	0.751	0.753	0.755
vehicle2	0.933	0.931	0.943	0.936	0.935	0.941	0.939	0.941
vowel	0.886	0.886	0.940	0.894	0.899	0.936	0.929	0.929
wdbc	0.966	0.966	0.965	0.967	0.967	0.965	0.966	0.964
ranking	5.80	5.45	4.45	3.95	5.45	4.60	3.45	2.85

- [5] R. M. Cruz, R. Sabourin, G. D. Cavalcanti, and T. Ing Ren, "META-DES: A dynamic ensemble selection framework using meta-learning," *Pattern Recognition*, vol. 48, no. 5, p. 1925–1935, 2015.
- [6] T. Woloszynski, M. Kurzynski, P. Podsiadlo, and G. W. Stachowiak, "A measure of competence based on random classification for dynamic ensemble selection," *Information Fusion*, vol. 13, no. 3, pp. 207–213, 2012.
- [7] Z.-H. Zhou, *Ensemble Methods: Foundations and Algorithms*, 1st ed. Chapman and Hall/CRC, 2012.
- [8] M. Monteiro, A. S. Britto, J. P. Barddal, L. S. Oliveira, and R. Sabourin, "Exploring diversity in data complexity and classifier decision spaces for pool generation," *Information Fusion*, vol. 89, pp. 567–587, 2023.
- [9] J. Elmi and M. Eftekhari, "Multi-layer selector (mls): Dynamic selection based on filtering some competence measures," *Applied Soft Computing*, vol. 104, p. 107257, 2021.
- [10] V. S. Costa, A. D. S. Farias, B. Bedregal, R. H. Santiago, and A. M. d. P. Canuto, "Combining multiple algorithms in classifier ensembles using generalized mixture functions," *Neurocomputing*, vol. 313, pp. 402–414, 2018.
- [11] B. Swaminathan, S. Palani, and S. Vairavasundaram, "Meta learning-based dynamic ensemble model for crop selection," *Applied Artificial Intelligence*, vol. 36, no. 1, p. 2145646, 2022.
- [12] L. I. Kuncheva, *Combining Pattern Classifiers: Methods and Algorithms: Second Edition*. Wiley Publishing, 2014.
- [13] L. Wang, T. Mo, X. Wang, W. Chen, Q. He, X. Li, S. Zhang, R. Yang, J. Wu, X. Gu *et al.*, "A hierarchical fusion framework to integrate homogeneous and heterogeneous classifiers for medical decision-making," *Knowledge-Based Systems*, vol. 212, p. 106517, 2021.
- [14] P. R. Cavalin, R. Sabourin, and C. Y. Suen, "Dynamic selection approaches for multiple classifier systems," *Neural Computing and Applications*, vol. 22, pp. 673–688, 2013.
- [15] D. Dua and C. Graff, "UCI machine learning repository," 2020. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [16] L. B. de Amorim, G. D. Cavalcanti, and R. M. Cruz, "The choice of scaling technique matters for classification performance," *Applied Soft Computing*, vol. 133, p. 109924, 2023.
- [17] T. T. Nguyen, A. V. Luong, M. T. Dang, A. W. C. Liew, and J. McCall, "Ensemble Selection based on Classifier Prediction Confidence," *Pattern Recognition*, vol. 100, p. 107104, 2020.
- [18] R. M. O. Cruz, L. G. Hafemann, R. Sabourin, and G. D. C. Cavalcanti,

TABLE IV

AVERAGE MCC AND RANKINGS OF THE EVALUATED METHODS FOR EACH DATASET. THE BEST RESULT PER DATASET IS PRESENTED IN BOLD.

Datasets	KNORA-U	KNOP	META-DES	DES-P	Random	PS-DES-MCC	PS-DES-F	PS-DES-acc
appendicitis	0.503	0.517	0.489	0.563	0.500	0.529	0.574	0.566
australian	0.694	0.691	0.694	0.691	0.690	0.695	0.697	0.695
balance	0.794	0.801	0.812	0.795	0.798	0.813	0.809	0.806
cmc	0.227	0.232	0.196	0.238	0.224	0.187	0.223	0.232
column_3C	0.755	0.753	0.748	0.759	0.756	0.744	0.755	0.757
diabetes	0.466	0.468	0.454	0.460	0.464	0.453	0.466	0.472
glass1	0.019	0.008	0.226	0.242	0.134	0.247	0.189	0.196
glass6	0.729	0.723	0.770	0.755	0.749	0.775	0.774	0.764
haberman	0.131	0.130	0.144	0.193	0.157	0.138	0.201	0.185
hayes	0.416	0.427	0.505	0.467	0.461	0.499	0.529	0.524
heart	0.673	0.677	0.667	0.682	0.673	0.674	0.677	0.674
led7digit	0.696	0.697	0.678	0.693	0.692	0.666	0.674	0.694
mammographic	0.660	0.662	0.644	0.655	0.660	0.652	0.647	0.650
musk	0.547	0.551	0.566	0.575	0.562	0.572	0.574	0.579
pima	0.471	0.466	0.451	0.462	0.461	0.452	0.455	0.464
sonar	0.547	0.544	0.594	0.570	0.576	0.575	0.595	0.600
vehicle	0.672	0.673	0.684	0.674	0.671	0.676	0.681	0.684
vehicle2	0.867	0.863	0.886	0.873	0.872	0.884	0.880	0.883
vowel	0.784	0.782	0.885	0.798	0.808	0.877	0.864	0.864
wdbc	0.933	0.933	0.931	0.934	0.933	0.931	0.932	0.929
ranking	5.45	5.30	4.80	3.80	5.35	4.65	3.55	3.10

TABLE V

STATISTICAL ANALYSES FOR PS-DES-ACC AGAINST STATE-OF-THE-ART DS METHODS. THE LINE (W/T/L) PRESENTS THE NUMBER OF WINS, TIES, AND LOSSES IT OBTAINED COMPARED TO THE COLUMN-WISE TECHNIQUE. THE P-VALUE LINE SHOWS THE RESULT OF APPLYING THE PAIRED WILCOXON STATISTICAL TEST. STATISTICALLY DIFFERENT RESULTS ($\alpha = 0.5$) ARE HIGHLIGHTED IN BOLD.

Metric		KNORA-U	KNOP	META-DES	DES-P	Random	PS-DES-MCC	PS-DES-F
Accuracy	w/t/l	16/0/4	14/0/6	13/0/7	13/0/7	17/1/2	12/1/7	12/0/8
	p-value	0.002	0.003	0.005	0.062	0.001	0.016	0.071
F-score	w/t/l	17/0/3	17/0/3	13/0/7	14/0/6	18/0/2	12/0/8	12/0/8
	p-value	0.000	0.000	0.049	0.020	0.000	0.066	0.139
MCC	w/t/l	16/0/4	14/0/6	13/0/7	13/0/7	18/0/2	13/0/7	11/0/9
	p-value	0.001	0.001	0.029	0.077	0.000	0.049	0.261

TABLE VI

STATISTICAL ANALYSES FOR PS-DES-F AGAINST STATE-OF-THE-ART DS METHODS. THE LINE (W/T/L) PRESENTS THE NUMBER OF WINS, TIES, AND LOSSES THAT IT OBTAINED COMPARED TO THE COLUMN-WISE TECHNIQUE. THE P-VALUE LINE SHOWS THE RESULT OF APPLYING THE PAIRED WILCOXON STATISTICAL TEST. STATISTICALLY DIFFERENT RESULTS ($\alpha = 0.5$) ARE HIGHLIGHTED IN BOLD.

Metric		KNORA-U	KNOP	META-DES	DES-P	Random	PS-DES-MCC	PS-DES-acc
Accuracy	w/t/l	13/1/6	13/1/6	13/0/7	10/0/10	13/0/7	15/0/5	8/0/12
	p-value	0.015	0.032	0.049	0.237	0.018	0.024	0.934
F-score	w/t/l	15/0/5	15/0/5	13/0/7	11/0/9	15/0/5	14/0/6	8/0/12
	p-value	0.003	0.005	0.174	0.147	0.003	0.066	0.869
MCC	w/t/l	13/0/7	14/0/6	14/0/6	11/0/9	14/0/6	14/0/6	9/0/11
	p-value	0.010	0.013	0.045	0.139	0.005	0.041	0.751

TABLE VII

STATISTICAL ANALYSES FOR PS-DES-MCC AGAINST STATE-OF-THE-ART DS METHODS. THE LINE (W/T/L) PRESENTS THE NUMBER OF WINS, TIES, AND LOSSES THAT IT OBTAINED COMPARED TO THE COLUMN-WISE TECHNIQUE. THE P-VALUE LINE SHOWS THE RESULT OF APPLYING THE PAIRED WILCOXON STATISTICAL TEST. STATISTICALLY DIFFERENT RESULTS ($\alpha = 0.5$) ARE HIGHLIGHTED IN BOLD.

Metric		KNORA-U	KNOP	META-DES	DES-P	Random	PS-DES-F	PS-DES-acc
Accuracy	w/t/l	12/0/8	10/0/10	10/0/10	9/0/11	12/0/8	5/0/15	7/1/12
	p-value	0.174	0.205	0.580	0.649	0.273	0.978	0.984
F-score	w/t/l	13/0/7	12/0/8	7/0/13	9/0/11	13/0/7	6/0/14	8/0/12
	p-value	0.032	0.045	0.861	0.522	0.101	0.938	0.938
MCC	w/t/l	13/0/7	12/0/8	9/0/11	9/0/11	11/0/9	6/0/14	7/0/13
	p-value	0.071	0.088	0.676	0.663	0.194	0.962	0.955