**IEEE** *Access*
Multidisciplinary : Rapid Review : Open Access Journal

# NE-nu-SVC: A new nested ensemble clinical decision support system for effective diagnosis of coronary artery disease

**MOLOUD ABDAR[1], U. RAJENDRA ACHARYA[2,3,4], NIZAL SARRAFZADEGAN[5,6], AND VLADIMIR MAKARENKOV[1]**

[1]Department of Computer Science, University of Quebec in Montreal, Montreal (QC), H2X 3Y7, Canada
[2]Department of Electronics and Computer Engineering, Ngee Ann Polytechnic, Singapore 599489
[3]Department of Biomedical Engineering, School of Science and Technology, SUSS University, Singapore 599491
[4]School of Medicine, Faculty of Health and Medical Sciences, Taylor's University, Subang Jaya 47500, Malaysia
[5]Isfahan Cardiovascular Research Center,Cardiovascular Research Institute,Isfahan University of Medical Sciences,Isfahan,Iran 8174673461
[6]School of Population and Public Health, Faculty of Medicine, University of British Columbia, 2206 East Mall, Vancouver, V6T, BC, Canada

Corresponding author: Moloud Abdar (e-mails: m.abdar1987@gmail.com and abdar.moloud@courrier.uqam.ca).

**ABSTRACT** Coronary artery disease (CAD) is one of the main causes of cardiac death around the world. Due to its significant impact on the society, early and accurate detection of CAD is essential. This study proposes a novel nested ensemble nu-Support Vector Classification (NE-nu-SVC) model which combines several traditional machine learning methods and ensemble learning techniques for effective diagnosis of CAD. We validated our model using two well-known CAD datasets (Z-Alizadeh Sani and Cleveland). To improve the performance of the model, we selected clinically significant features from the datasets using a genetic search algorithm. To further improve our results, we applied a multi-level filtering technique to balance the data using the ClassBlancer and Resample methods. Our base algorithm, nu-SVC, is performed using four well-known kernel functions (linear, polynomial, radial basis (RBF) and sigmoid). The proposed NE-nu-SVC model provided the highest accuracy of 94.66% and 98.60% to predict CAD entities in the Z-Alizadeh Sani and Cleveland CAD datasets, respectively. Our system can aid the clinicians to diagnose CAD accurately and may probably replace other invasive diagnostic techniques.

**INDEX TERMS** Coronary artery disease (CAD), machine learning, ensemble learning, nested ensemble (NE) model, genetic algorithm.

## I. INTRODUCTION

RECENT advances in artificial intelligence (AI) have led to the emergence of new intelligent automatic systems. AI, nowadays, acts as a bridge between different fields such as economics, biology, physics, mathematics, chemistry, etc [1]–[8]. AI methods can be applied to solve a variety of problems existing in those fields, providing more accurate solutions than standard classification methods.

Data mining and machine learning algorithms have been also widely used in the fields of bioinformatics and healthcare science. The volume of data in these fields has been growing very quickly [9]. Moreover, biological and medical data are usually heterogeneous and complex. All these factors motivate a prompt development of new data mining approaches, including the development of specialized machine learning techniques [10]. Machine learning algorithms, such as DTs (decision trees) [11], SVMs (support vector machines) [12]–[14] and ANNs (artificial neural networks) [15], combined to the powerful deep learning approach [16], can be used to tackle various challenging issues arising in bioinformatics and healthcare science, including protein structure prediction and disease identification [17], [18]. Nowadays, one of the most relevant challenges in healthcare science is the improvement of the performance of Clinical Decision

Support Systems (CDSS) [19] meant to predict and cure many important diseases such as Coronary Artery Disease (CAD) and other cardiovascular diseases [20], obesity [21], Chronic Obstructive Pulmonary Disease (COPD) [22], Alzheimer disease [23], prostate cancer [24], etc. CDSSs can be very beneficial to diagnose many of these diseases, including CAD, which is one of the major types of heart diseases [20] According to recent reports, CAD is the most common cardiovascular disease in the United States of America, being the leading cause of heart attacks among both male and female population [25]. As indicated in Nahar et al. [26], CAD is also one of the main causes of death in Australia and the United Kingdom. Therefore, it is very important to provide an efficient approach for CAD prediction, and we propose to do it here using different machine learning techniques.

The primary goal of this work is to introduce and test a new CDSS intended for accurate CAD prediction. Our new model, called Nested Ensemble nu-SVC (NE-nu-SVC), allows one to use different ensemble learning techniques with traditional machine learning algorithms. The proposed nu-SVC model including four different kernel functions was used to analyze two well-known CAD datasets (Z-Alizadeh Sani CAD [27] and Cleveland CAD datasets [25]). In order to eliminate redundant features and thus improve the model's performance, we carried out feature selection using a genetic search algorithm. In addition, the entities in both datasets were reweighted using a multi-level data balancing by the way of the supervised ClassBalancer (CB) method [28] and resampled using both supervised and unsupervised resample approaches [29]. The applied multi-level data balancing led to the prediction accuracy improvement for both minority and majority classes. Then, the proposed NE-nu-SVC model was applied. To achieve a better accuracy, four ensemble learning techniques at three different levels of the model were combined in the framework of NE-nu-SVC. As a result, we could greatly improve the accuracy of the traditional nu-SVC model with all kernel functions for both the Z-Alizadeh Sani and Cleveland CAD datasets [27], [25].

The rest of the paper is structured as follows. In Section II, we discuss the related work in the field. Section III describes materials, methods and the proposed approach. Section IV discusses the experimental results. Finally, Section V presents the conclusions of this study and the ideas for future research.

## II. RELATED WORK

Many recent studies address the problem of an early diagnosis of heart disease [27], [30]–[39]. Here, we briefly discuss those which are directly related to our work. Several traditional machine learning methods have been applied by Alizadehsani et al. [27] to predict CAD and some of its instances. The accuracy scores of 86.14%, 83.17% and 83.50% were obtained by these authors as to

an early recognition of the LAD (left anterior descending) artery, LCX (left circumflex) artery and RCA artery (right coronary artery) cases of CAD, respectively. In another study, Tayefi et al. [30] put forward a new model to predict CHD (Coronary Heart Disease). To this end, Tayefi et al. described a CHD prediction model using decision trees. Their model provided the average CHD prediction accuracy of 95.3%. Arabasadi et al. [31] introduced a new hybrid machine learning model to detect CAD. To do so, a genetic algorithm and an artificial neural network were combined. The method proposed by Arabasadi et al. provided a relatively good performance on the Z-Alizadeh Sani data with the accuracy, sensitivity and specificity scores equal to 93.85%, 97% and 92%, respectively. On the other hand, Alkeshuosh et al. [32] generated different rules for CAD detection using machine learning. For this purpose, the authors applied the well-known PSO (particle swarm optimization) evolutionary algorithm. Furthermore, the performance of the PSO algorithm was compared to that of the C4.5 algorithm. Alkeshuosh et al. showed that the average accuracy of their PSO method, which outperformed the C4.5 algorithm, was around 87%.

Abdar [33] applied four well-known DT algorithms, including CHAID (Chi-Square Automatic Interaction Detection), C5.0, QUEST (Quick, Unbiased and Efficient Statistical Tree) and CART (Classification And Regression Trees) to analyze the Cleveland CAD data. According to the results of that study, the C5.0 algorithm provided the greatest accuracy of 85.33% among the competing methods. Several simple rules were also generated by C5.0. Babič et al. [34] addressed the problem of CAD detection by applying different machine learning methods to three real CAD datasets. These authors focused on the two following directions: (1) a predictive CAD analysis with SVM, naïve Bayes classifier, neural networks and decision trees, and (2) a descriptive CAD analysis using association and decision rules. Babič et al. indicated that SVM was the best performer among the methods compared in this work. Polat et al. [35] applied the k-NN (k-nearest neighbour) algorithm as a preprocessing step for CAD detection. An AIRS (Artificial Immune Recognition System) based on a fuzzy resource allocation mechanism was then proposed to recognize CAD patients. The best CAD prediction accuracy reported by Polat et al. was 87%.

Acharya et al. [36] tackled the problem of CAD prediction by using the electrocardiogram (ECG) signals as input data for various machine learning techniques. Thus, a new automated diagnostic system for CAD and Myocardial Infarction (MI) detection was proposed. The statistical model described by Acharya et al. included three main methods: DCT (Discrete Cosine Transform), EMD (Empirical Mode Decomposition) and DWT (Discrete Wavelet Transform). The proposed system showed a good performance on real CAD data with an average accuracy of 98.5%. Patidar et al. [37] presented a new approach for CAD prediction using the tunable-Q wavelet

transform (TQWT) method. TQWT divides the heart rate signals into different sub-bands in order to improve the diagnostic feature selection. By using LS-SVM (Least-Squares Support Vector Machine) and PCA (Principal Component Analysis), Patidar et al. obtained the average CAD recognition accuracy of 99.72%. Kausar et al. [38] combined an unsupervised clustering method and a supervised classification technique for timely detection of CAD using an ensemble technique. At the first stage, PCA was carried out. Then, the SVM and K-means algorithms were applied. Mahmoodabadi and Tabrizi [39] introduced a new intelligent system, called Imperialist Competitive Algorithm (ICA), to predict CAD. This system included both the decision tree and evolutionary algorithms and provided the average CAD detection accuracy of 94.92%.

## III. MATERIALS AND METHODS

In this section, we first present the two real CAD datasets used in our study. Then, we describe our novel method based on the ensemble learning approach. Moreover, we shortly discuss some important features of the traditional machine learning algorithms considered in our work.

### A. DESCRIPTION OF THE Z-ALIZADEH SANI AND CLEVELAND CAD DATASETS

In order to test our new model, we decided to use two well-known CAD datasets available at the University of California, Irvine, machine learning repository (UCI). Specifically, the Z-Alizadeh Sani [27], [40] and Cleveland [25], [41] datasets were considered. The Z-Alizadeh Sani dataset includes 303 patients' records described by 56 features; 55 of them were selected as input and one as output in our prediction model. Namely, this dataset contains the data for 216 CAD (sick) patients and 87 non-CAD (healthy) patients. Four main types of features are available for these data: echo, symptom and examination, ECG and laboratory, and demographic features (for more information regarding the Z-Alizadeh Sani dataset, see Supplementary Material).

The angiography procedure was used in the study of Alizadeh Sani et al. [27] to measure the stenosis of each artery. When a patient had the diameter that was greater than or equal to 50%, he/she was categorized as a CAD patient, otherwise as a Normal patient. The Z-Alizadeh Sani data include 71% of positive records (CAD-affected patients) and 29% of negative records (Normal patients). This means that these data are not well-balanced. To classify such unbalanced data more effectively, we applied a multi-level balancing approach [28], [29] (the discretization ranges of the Z-Alizadeh Sani data features are reported by Alizadeh Sani et al. [27]). More information regarding the discretization ranges of heart disease features can be found in Braunwald's Heart Book [42].

The Cleveland CAD dataset is another well-known heart disease dataset, which has been widely studied in

the literature [25], [41]. This dataset contains 303 records as well, which are described 14 features; 13 of them were chosen as the input of our model and the remaining one was selected as our target attribute. The CAD data patients are categorized into 5 major classes. The first of them represents healthy patients (164 records; Class 1 in our work), whereas the four other classes correspond to different types of CAD patients (139 records in total). Here, we combined the data of the four latter classes into a general class of CAD patients (Class 2 in our work). The detailed information about the main features of the Cleveland CAD dataset is available in [25].

### B. DESCRIPTION OF THE NEW MODEL

In this study, we first applied several machine learning methods, including J48, Random Forest, NaiveBayes, BayesNet, Multilayer Perceptron, C-SVC and nu-SVC, to analyze the Z-Alizadeh Sani CAD dataset [27]. We found that the traditional nu-SVC model, which usually works well in practice, provides very average results for the Z-Alizadeh Sani data. This traditional nu-SVC model was used as a core of our new model (the nested ensemble (NE)) with a goal of improving its performance when predicting CAD.

The nested ensemble (NE) approach allows one to combine several ensemble learning techniques within one model [43]. The main idea of this approach consists of using an ensemble learning technique inside of another ensemble learning technique. The general view of an NE model is presented in Fig. 1 and Algorithm 1. Moreover, we can use multiple classifiers (algorithms) with each ensemble learning technique. In general, an NE model can include different numbers of ensemble learning techniques at different levels of the model. In this study, we considered a three-level NE model using four ensemble learning techniques (see Fig. 2).

At the first level, we used the stacking technique as our first ensemble learning technique. The stacking technique has two main components: "classifier" and "metaClassifier". Within the "classifier" component, we used the three following methods: the nu-SVC, SGD (Stochastic Gradient Descent) and Random Forest. Here, the Random Forest classifier was embedded into our stacking ensemble learning technique. The loss function we used for training in SGD was the Hinge Loss (SVM) function. As to "metaClassifier", we used the bagging ensemble learning technique. At the second level, we added one more ensemble learning technique to our model. Finally, we applied a voting ensemble technique (also called the Vote technique) as a classifier at the previous level (bagging technique). This voting technique included the SMO and Naïve Bayes classifiers. The NE model we consider in this study was combined with genetic-based feature selection and multi-level data balancing. The complete schematic view of the Nested Ensemble nu-SVC model proposed in our work is presented in Fig. 2.
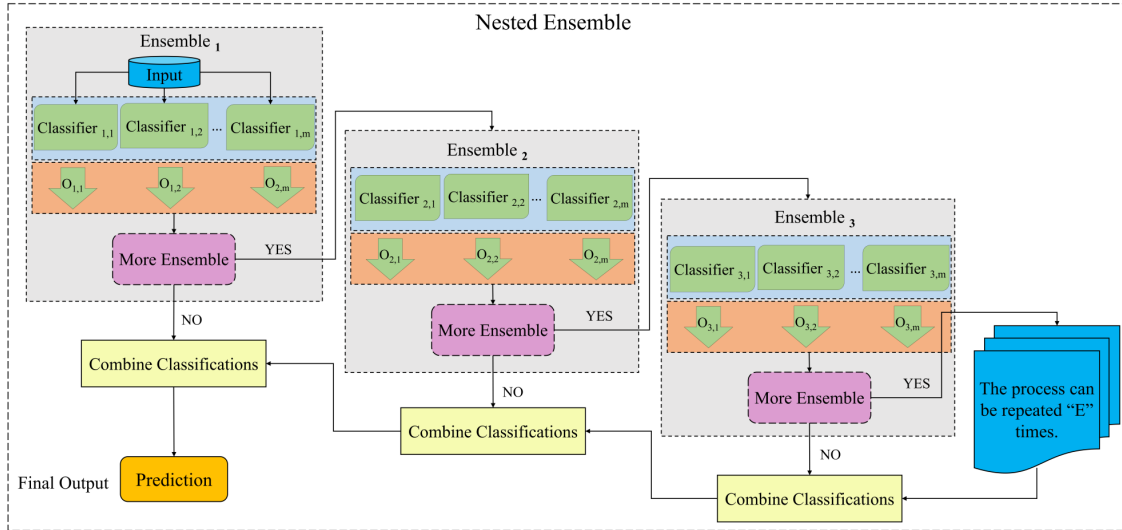
FIGURE 1: General scheme of a nested ensemble (NE) model.

---

**Algorithm 1:** A general Nested Ensemble (NE) model

**Input:** CAD dataset: $D = D_{train} \cup D_{test}$
**Output:** $O_f$: Best classification output

1 **Begin** Calculate the portance rate of each feature using a feature selection algorithm;
2 **if** $D$ *is not balanced* **then**
3     Balance $D$ using a balancing technique;
4 **end**
5 Select the number of $L$ levels in the NE model;
6 **for** $l = 1, \cdots, L$ **do**
7     Train the base machine learning algorithms, ensemble learning techniques and *metaClassifier* at different levels of the NE model using $D_{train}$;
8 **end**
9 Classify unseen records from $D_{test}$ using the NE model;
10 **return** Best classification result found;
11 **End**

---

As shown in Fig. 2, within our model, we also carry out the $K$-fold cross-validation technique, in which the value of $K$ is set to 10. Because both real datasets we consider here are not very large (i.e., 303 records in each of them), we can use the $K$-fold cross-validation technique which is generally very effective with this kind of data. By using the $K$-fold cross-validation technique, the problem of data bias can be minimized [44]. As mentioned previously, our base nu-SVC algorithm was used with four different kernel functions (linear, polynomial, RBF and sigmoid). We selected the most important data features using a genetic search algorithm. In order to balance the data (this was especially necessary for the Z-Alizadeh Sani

dataset), we used the multi-level balancing approach. Our tests suggested that a three-level NE model can generate accurate results while predicting CAD. Since both datasets include categorical features (e.g., gender), we applied one-hot encoding to deal with categorical features since they cannot be directly processed by machine learning algorithms.

### 1) Problem definition and base algorithms used

In this section, we discuss our nested ensemble model [45], [46]. First, we give the formal definition of a nested ensemble (NE) model and present the NE model parameters used in our work. Let $\mathscr{D} = \{\mathscr{D}_1, \mathscr{D}_2, \mathscr{D}_3,..., \mathscr{D}_n\}$ be a set of $n$ datasets (each dataset includes different numbers of $r$ records), $A = \{A_1, A_2, A_3,..., A_m\}$ be a set of $m$ machine learning algorithms (base algorithms) and $E = \{E_1, E_2, E_3,..., E_k\}$ be a set of $k$ ensemble learning techniques. Let $L$ be the number of levels in the NE model. It should be noted that different machine learning algorithms and ensemble learning techniques can be used at different levels of the model.

We can add ensemble learning techniques and machine learning algorithms to different levels of an NE model till the prediction results improve. By using the union of the sets $E_i$ ($i = 1, 2, \ldots, k$), we get a general NE model: NE $= \{E_1 \cup E_2 \cup E_3... \cup E_k\}$.

In our work, we used a three-level NE model including four ensemble learning techniques, four machine learning algorithms and two well-known CAD datasets (i.e., $L = 3$, $n = 2$, $m = 4$ and $k = 4$).

To the best of our knowledge, no existing work applied a multi-level NE model to analyze CAD data. The detailed information about the traditional machine learning algorithms and the ensemble learning techniques used in our work can be found in the Supplemental material. These
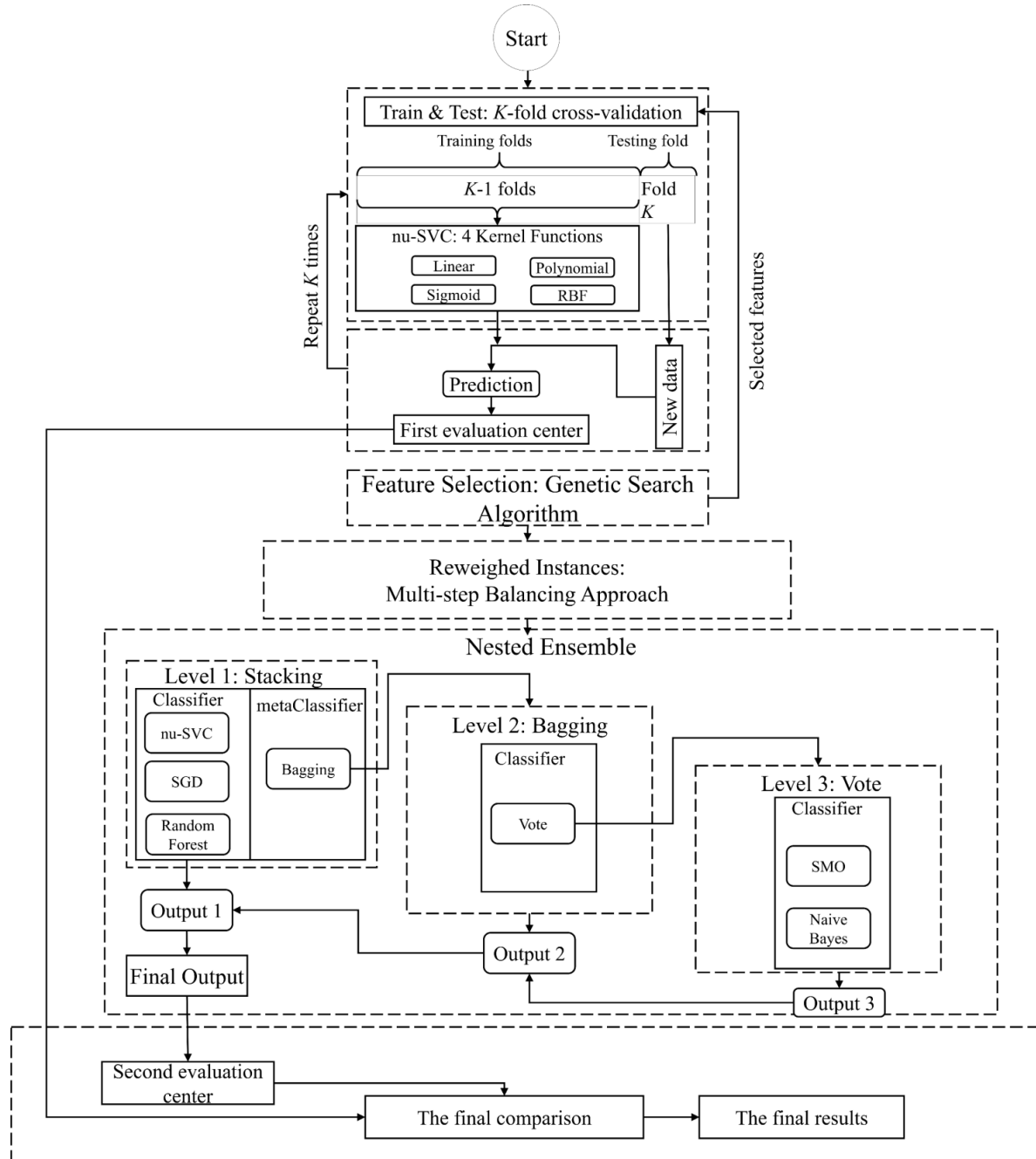
FIGURE 2: A block diagram of the proposed NE-nu-SVC model.

algorithms and techniques are also briefly described there.

out the nu-SVC algorithm.

## IV. EXPERIMENTAL RESULTS

In this section, we discuss the results obtained when analyzing the Z-Alizadeh Sani CAD data with our NE-nu-SVC model. An IBM PC computer, equipped with a 2.30 GHz Intel Core i7 CPU and 8 GB of RAM, was used in our simulations. Our model was implemented using the 3.9.1 version of the WEKA package [29]. In addition, the LIBSVM (a Library for Support Vector Machines) open source machine learning library [47] was used to carry

### A. PERFORMANCE METRICS

Various statistical measures can be considered to evaluate and compare the performance of machine learning methods [48]. Here, we used the following popular metrics: Recall (True Positive Rate (TPR)), False Positive Rate (FPR), Precision, F-Measure, Accuracy, ROC (receiver operating characteristic) area, Kappa statistic, Pc, MAE (Mean Absolute Error) and RMSE (Root Mean Squared Error). These metrics are listed in Eq. (1-9) below:

$$FPR = \frac{FP}{FP + TN}, \qquad (1)$$

$$Precision = \frac{TP}{TP + FP}, \qquad (2)$$

$$Recall = TPR = \frac{TP}{TP + FN}, \qquad (3)$$

$$F - Measure = \frac{2 \times (Recall \times Precision)}{Recall + Precision}, \qquad (4)$$

$$Accuracy = \frac{TP + TN}{P + N}, \qquad (5)$$

$$Kappa - statistic = \frac{Accuracy - P_c}{1 - P_c}, \qquad (6)$$

$$P_c = \frac{(TP + FP) \times (TP + FN) + (FN + TN) \times (FP + TN)}{(P + N)^2}, \qquad (7)$$

$$MAE = \frac{1}{r} \sum_{i=1}^{r} |(y_i - \widehat{y}_i)|, \qquad (8)$$

$$RMSE = \sqrt{\frac{1}{r} \sum_{i=1}^{r} |(y_i - \widehat{y}_i)|^2}, \qquad (9)$$

where $N$ stands for the number of negative records in the data (original), $P$ stands for the number of positive records in the data, $TP$ (true positives) stands for the number of positive records that are classified correctly, $FN$ (false negatives) stands for the number of positive records that are misclassified as negative, $FP$ (false positives) stands for the number of negative records that are misclassified as positive, $TN$ (true negatives) stands for the number of negative records that are classified correctly, $r$ denotes the number of records in the dataset, $i$ denotes the sample index, $y_i$ denotes the actual value of record $i$ and $\widehat{y}_i$ denotes the predicted value of record $i$.

### B. RESULTS OBTAINED PRIOR TO FEATURE SELECTION

This section presents the experimental results obtained using the traditional nu-SVC model when all original features of the Z-Alizadeh Sani CAD dataset were considered. The nu-SVC model was applied with four different kernel functions: linear, polynomial, RBF and sigmoid. The obtained results are presented in Table 1 and Fig. 3.

According to Table 1, the results provided by different kernel functions used within the traditional nu-SVC model vary a lot. On one hand, in terms of TP Rate, FP Rate, Precision, Recall and F-measure, the RBF function yielded good results for the CAD class, whereas it did not show good performance for the Normal class. On the other hand, the sigmoid function provided good results in terms of TP Rate, FP Rate and Recall for the Normal class, but not for the CAD class. Here, our main finding was that traditional nu-SVC had the greatest accuracy when the kernel function was polynomial. Moreover, Figure 3 shows that nu-SVC with the linear kernel function yielded

the best results according to the Kappa statistic, and the MAE and RMSE measures.

### C. FEATURE SELECTION PROCEDURE

In this step, a genetic search algorithm was carried out to select the most significant features of the Z-Alizadeh Sani CAD data. The applied genetic search algorithm by Goldberg [49], implemented in WEKA, uses a correlation-based variable selection approach to eliminate redundant features [50]. The values of population size, number of generations and report frequency during our feature selection were set to 20. The probability of crossover and the probability of mutation were set to 0.6 and 0.033, respectively. More details on the applied feature selection procedure are presented in Tables A.1 and A.2 (see Appendix A). Initial population features and generated features provided by the genetic search algorithm for the Z-Alizadeh Sani CAD dataset are presented in Tables A.1 and A.2. The numbers in the Subset column in both Tables A1 and A.2 show the feature order in the original data file. Applying this procedure, we selected 16 most significant features which were used in further investigation. The selected features were as follows: Age, DM (Diabetes Mellitus), HTN (Hyper Tension), CRF (Chronic Renal Failure), BP (Blood Pressure), Typical Chest pain, Dyspnea, Atypical, Nonanginal, Q Wave, T-inversion, ESR (Erythrocyte Sedimentation Rate), K (Potassium), EF (Ejection Fraction), RWMA (Regional Wall Motion Abnormality) and VHD (Valvular Heart Disease). It should be noted that we applied feature selection prior to one-hot encoding as suggested by Hasanin et al. [51].

### D. RESULTS OBTAINED AFTER FEATURE SELECTION

This section describes the results provided by the traditional nu-SVC model after feature selection. The nu-SVC model was applied with four kernel functions as discussed earlier. The results obtained when the 16 selected features of the Z-Alizadeh Sani CAD dataset were used are presented in Table 2 and Fig. 4.

As shown in Table 2 and Fig. 4, the general trends observed when using only the 16 selected features of the Z-Alizadeh Sani dataset are similar to those found for the original Z-Alizadeh Sani data. In other words, we can notice that traditional nu-SVC had a good performance either for the CAD class or for the Normal class. However, the accuracy results, especially in the case of the linear and polynomial kernel functions, were generally much better after the feature selection. According to Fig. 4, the traditional nu-SVC model with the linear kernel provides the best CAD detection having with the highest value of the Kappa statistic (0.6496), and the lowest values of MAE (0.1386) and RMSE (0.3723), followed by the polynomial kernel. The traditional nu-SVC model with the sigmoid kernel does not show a good performance compared to the three other kernels. We can see that the sigmoid

TABLE 1: The results provided by the traditional nu-SVC model when all original features of the Z-Alizadeh Sani CAD dataset were considered.

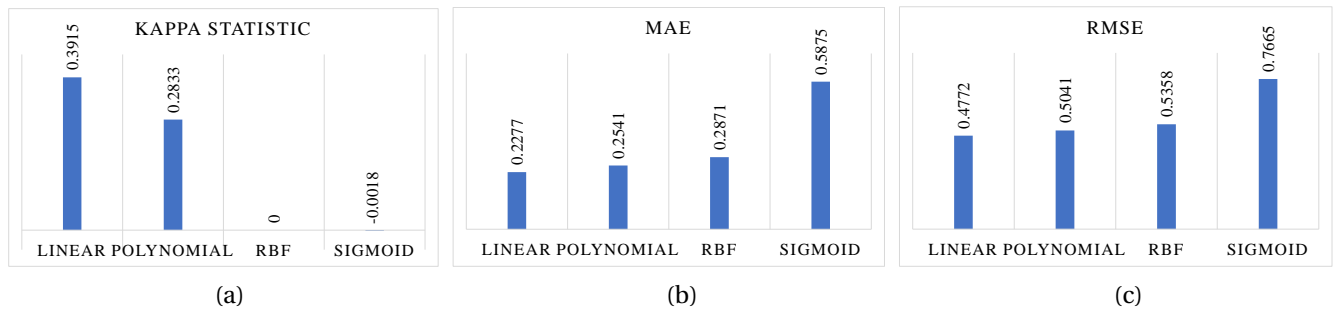| Measures | | nu-SVC | | | |
|---|---|---|---|---|---|
| | | **Linear** | **Polynomial** | **RBF** | **Sigmoid** |
| FPR | CAD (%) | 54.00 | 66.70 | 100 | 29.90 |
| | Normal (%) | 10.20 | 8.80 | 0.0 | 70.40 |
| | Average (%) | 41.40 | 50.10 | 71.30 | 41.50 |
| Precision | CAD (%) | 80.50 | 77.30 | 71.30 | 71.10 |
| | Normal (%) | 64.50 | 60.40 | 0.0 | 28.60 |
| | Average (%) | 75.90 | 72.40 | 50.80 | 58.90 |
| Recall | CAD (%) | 89.80 | 91.20 | 100 | 29.60 |
| | Normal (%) | 46.00 | 33.30 | 0.0 | 70.10 |
| | Average (%) | 77.20 | 74.60 | 71.30 | 41.30 |
| F-measure | CAD (%) | 84.90 | 83.70 | 83.20 | 41.80 |
| | Normal (%) | 53.70 | 43.00 | 0.0 | 40.70 |
| | Average (%) | 75.90 | 72.00 | 59.30 | 41.50 |
| ROC Area | CAD (%) | 67.90 | 62.30 | 50.00 | 49.90 |
| | Normal (%) | 67.90 | 62.30 | 50.00 | 49.90 |
| | Average (%) | 67.90 | 62.30 | 50.00 | 49.90 |
| Accuracy(%) | | **77.22** | 74.58 | 71.28 | 41.25 |



FIGURE 3: Comparison of the results provided by the traditional nu-SVC model with different kernel functions when all original features of the Z-Alizadeh Sani CAD dataset were considered: (a) Kappa statistic, (b) MAE and (c) RMSE.

TABLE 2: The results provided by the traditional nu-SVC model using the 16 selected features of the Z-Alizadeh Sani CAD dataset.

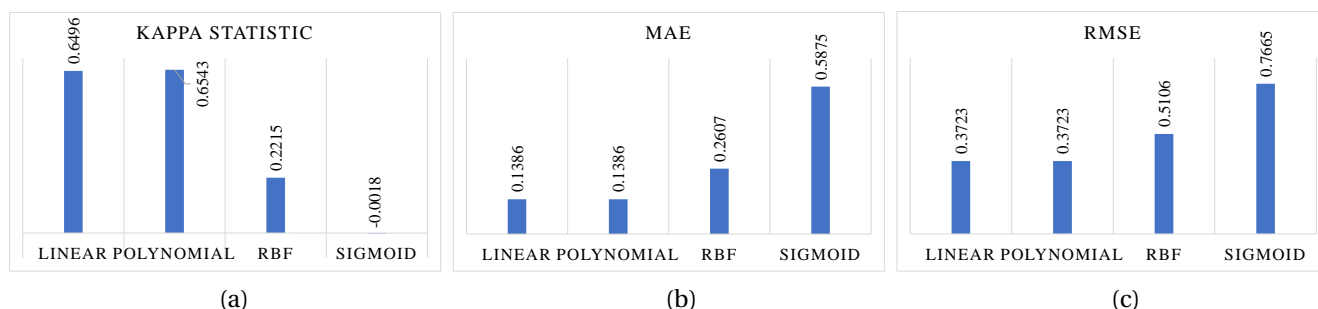| Measures | | nu-SVC | | | |
|---|---|---|---|---|---|
| | | **Linear** | **Polynomial** | **RBF** | **Sigmoid** |
| FPR | CAD (%) | 31.00 | 27.60 | 75.90 | 29.90 |
| | Normal (%) | 6.90 | 8.30 | 6.00 | 70.40 |
| | Average (%) | 24.10 | 22.10 | 55.80 | 41.50 |
| Precision | CAD (%) | 88.20 | 89.20 | 75.50 | 71.70 |
| | Normal (%) | 80.00 | 77.80 | 61.80 | 28.60 |
| | Average (%) | 85.80 | 85.90 | 71.50 | 58.90 |
| Recall | CAD (%) | 93.10 | 91.70 | 94.00 | 29.60 |
| | Normal (%) | 69.00 | 72.40 | 34.70 | 70.10 |
| | Average (%) | 86.10 | 86.10 | 73.90 | 41.30 |
| F-measure | CAD (%) | 90.50 | 90.40 | 83.70 | 41.80 |
| | Normal (%) | 74.10 | 75.00 | 34.70 | 40.70 |
| | Average (%) | 85.80 | 86.00 | 69.60 | 41.50 |
| ROC Area | CAD (%) | 81.00 | 82.00 | 59.10 | 49.90 |
| | Normal (%) | 81.00 | 82.00 | 59.10 | 49.90 |
| | Average (%) | 81.00 | 82.00 | 59.10 | 49.90 |
| Accuracy(%) | | **86.13** | **86.13** | 73.92 | 41.25 |

FIGURE 4: Comparison of the results provided by the traditional nu-SVC model with different kernel functions when the 16 selected features of the Z-Alizadeh Sani CAD dataset were considered: (a) Kappa statistic, (b) MAE and (c) RMSE.

kernel has generated a negative Kappa statistic value (-0.0018), and very high MAE (0.5878) and RMSE (0.7665) values.

### E. RESULTS OBTAINED AFTER DATA BALANCING

As mentioned earlier (see section III, subsection A), most of the entities (216 out of 303) of the Z-Alizadeh Sani CAD dataset belong to the CAD class and only 87 entities to the Normal class. In the previous section, we observed that the nu-SVC model provided different performances for these two classes (see the results in Tables 1 and 2). We can argue that the Z-Alizadeh Sani dataset is unbalanced (or weakly balanced). Hence, we used an approach to deal with such weakly balanced data.

Precisely, the multi-step balancing technique was carried out as follows. The Z-Alizadeh Sani dataset was first reweighted so that each class could get the same total weight. The supervised ClassBalancer (CB) method was used here. Since the CB method did not improve significantly the performance of the proposed model, the unsupervised resample method was also applied. Such a balancing approach is called a multi-step balancing. This two-step balancing allowed us to improve the prediction results for the Z-Alizadeh Sani dataset. It is worth mentioning that one can use different levels of balancing (two steps or more) to get a better performance. Since we used 10-fold cross validation, the balanced data were randomized using an unsupervised instance filter to avoid the overfitting problem. The results obtained after the supervised data balancing (using CB) and the two-step balancing (using both CB and resampling) are reported in Tables 3 and 4, respectively.

TABLE 3: Reweighted records obtained for the Z-Alizadeh Sani CAD data using the supervised ClassBalancer (CB) technique.

| No. | Label | Weight |
|-----|--------|--------|
| 1 | CAD | 151.5 |
| 2 | NORMAL | 151.5 |

After the data balancing, we applied the nu-SVC model once again using the same four kernel functions. We carried out our method 10 times with each kernel func-

TABLE 4: Final reweighted records obtained for the Z-Alizadeh Sani CAD data using two-step balancing, including both the supervised and unsupervised resample techniques.

| No. | Label | Weight |
|-----|--------|---------|
| 1 | CAD | 152.903 |
| 2 | NORMAL | 148.017 |

tion to find out whether the obtained results were stable or not. The average results generated after our two-step balancing are shown in Table 5 and Fig. 5.

It can be noted from Table 5 that the accuracy of the nu-SVC model with the RBF and sigmoid kernel functions increased (compared to the results presented in Table 2), whereas it decreased for the linear and polynomial kernels. In general, we can argue that the prediction of both CAD and Normal patient classes improved when the two-step balancing procedure was applied.

### F. RESULTS OBTAINED USING THE NESTED ENSEMBLE NU-SVC MODEL

A good machine learning model should correctly classify the entities of all classes present in a given dataset, whereas the classification results we have obtained so far have been good for one class only, either for the CAD class or for the Normal patient class. Thus, our new Nested Ensemble nu-SVC (NE-nu-SVC) model (see Fig. 2) was applied at this stage to improve the performance of the traditional nu-SVC model. This NE-nu-SVC model incorporates four well-known ensemble learning techniques: stacking, random forest, bagging and voting, which are used together. NE-nu-SVC was applied to the Z-Alizadeh Sani data after the feature selection and two-step balancing of the entities were carried out (as explained in the previous sections). The detailed results provided by NE-nu-SVC are presented in Table 6 and Fig. 6.

The results presented in Table 6 and Fig. 6 suggest that the application of the NE-nu-SVC model allowed us to improve drastically the prediction results for all kernel functions, compared to the traditional nu-SVC model (see Table 5 and Fig. 6). In order to verify the stability of the

**IEEE** *Access*

TABLE 5: The results provided by the traditional nu-SVC model using feature selection and two-step balancing for the Z-Alizadeh Sani CAD dataset.

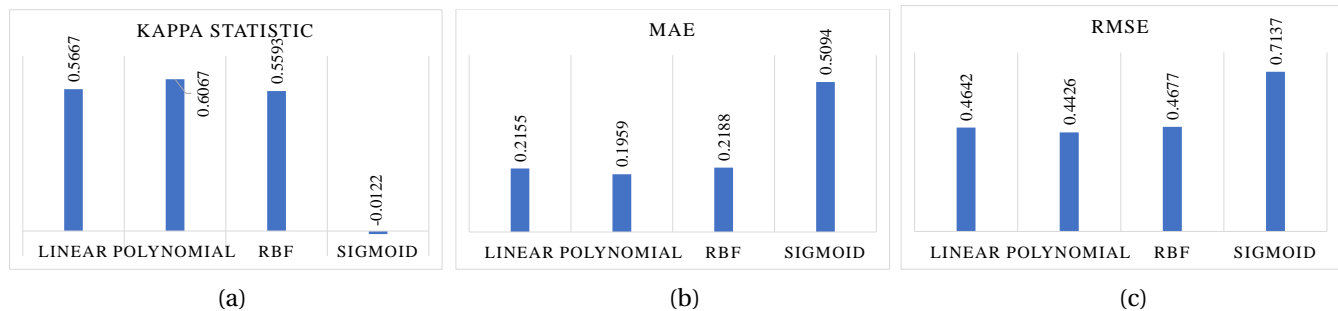| Measures | | nu-SVC | | | |
|---|---|---|---|---|---|
| | | Linear | Polynomial | RBF | Sigmoid |
| FPR | CAD (%) | 37.60 | 31.80 | 43.50 | 30.60 |
| | Normal (%) | 6.00 | 7.80 | 0.90 | 70.60 |
| | Average (%) | 22.10 | 20.00 | 22.60 | 50.30 |
| Precision | CAD (%) | 72.10 | 75.00 | 70.20 | 49.80 |
| | Normal (%) | 91.00 | 89.40 | 98.30 | 48.70 |
| | Average (%) | 81.40 | 82.10 | 84.00 | 49.30 |
| Recall | CAD (%) | 94.00 | 92.20 | 99.10 | 29.40 |
| | Normal (%) | 62.40 | 68.20 | 56.50 | 69.40 |
| | Average (%) | 78.50 | 80.40 | 78.10 | 49.10 |
| F-measure | CAD (%) | 81.60 | 82.70 | 82.20 | 36.90 |
| | Normal (%) | 74.00 | 77.40 | 71.70 | 57.30 |
| | Average (%) | 77.90 | 80.10 | 77.00 | 46.90 |
| ROC Area | CAD (%) | 78.20 | 80.20 | 77.80 | 49.40 |
| | Normal (%) | 78.20 | 80.20 | 77.80 | 49.40 |
| | Average (%) | 78.20 | 80.20 | 77.80 | 49.40 |
| Accuracy(%) | | 78.45 | **80.41** | 78.12 | 49.05 |



(a)



(b)



(c)

FIGURE 5: Comparison of the results provided by the traditional nu-SVC model with different kernel functions when the feature selection and two-step balancing of the entities of the Z-Alizadeh Sani CAD dataset were carried out: (a) Kappa statistic, (b) MAE and (c) RMSE.

TABLE 6: The results provided by the proposed NE-nu-SVC model for the Z-Alizadeh Sani data (after applying the feature selection and two-step balancing of the entities).

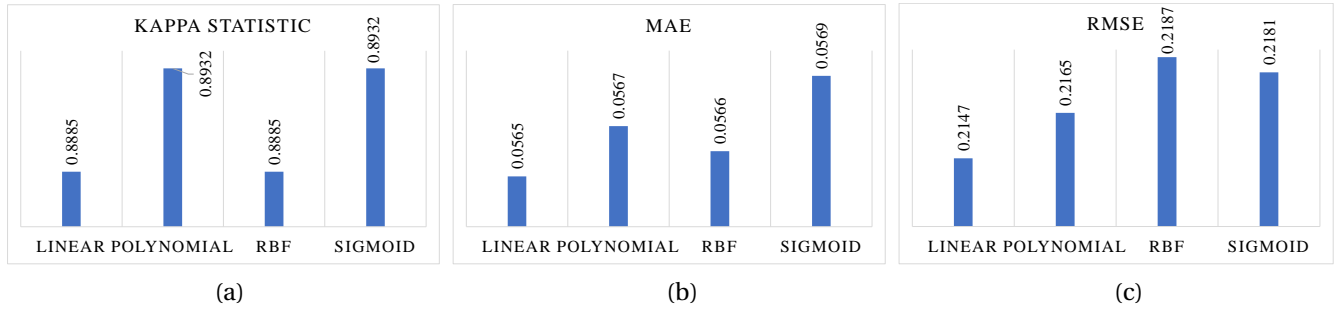| Measures | | NE-nu-SVC | | | |
|---|---|---|---|---|---|
| | | Linear | Polynomial | RBF | Sigmoid |
| FPR | CAD (%) | 7.10 | 7.00 | 7.10 | 7.10 |
| | Normal (%) | 4.10 | 3.70 | 4.10 | 3.70 |
| | Average (%) | 5.60 | 5.40 | 5.60 | 5.40 |
| Precision | CAD (%) | 93.30 | 93.40 | 93.30 | 93.40 |
| | Normal (%) | 95.60 | 96.10 | 95.60 | 96.10 |
| | Average (%) | 94.50 | 94.70 | 94.40 | 94.70 |
| Recall | CAD (%) | 95.90 | 96.30 | 95.90 | 96.30 |
| | Normal (%) | 92.90 | 92.90 | 92.90 | 92.90 |
| | Average (%) | 94.40 | 94.70 | 94.40 | 94.70 |
| F-measure | CAD (%) | 94.60 | 94.80 | 94.60 | 94.80 |
| | Normal (%) | 94.30 | 94.50 | 94.30 | 94.50 |
| | Average (%) | 94.40 | 94.70 | 94.40 | 94.70 |
| ROC Area | CAD (%) | 96.60 | 96.60 | 96.50 | 96.50 |
| | Normal (%) | 96.60 | 96.60 | 96.50 | 96.50 |
| | Average (%) | 96.60 | 96.60 | 96.50 | 96.50 |
| Accuracy(%) | | 94.43 | **94.66** | 94.43 | **94.66** |

FIGURE 6: Comparison of the results provided by the proposed NE-nu-SVC model with different kernel functions for the Z-Alizadeh Sani data (after applying the feature selection and two-step balancing of the entities): (a) Kappa statistic, (b) MAE and (c) RMSE.

proposed NE-nu-SVC model, it was applied 10 times for each kernel function. Very stable results were generated for all the evaluation metrics under consideration (see Eq.(1-9)). The running time is another important factor to be considered in CDSSs. Table 7 shows the average running time on the transformed records (obtained after applying the feature selection and two-step balancing of the entities) of the Z-Alizadeh Sani dataset before and after using the NE approach within the nu-SVC model.

Observing the results presented in Table 7, we can conclude that even though the proposed NE-nu-SVC model including several nested machine learning techniques requires more running time for the RBF and sigmoid kernel functions, compared to nu-SVC, it needs less running time for the linear and polynomial kernels. Thus, we can argue that, in general, our NE-nu-SVC model is not very time-consuming. The running time of our new model is very reasonable compared to the gain in accuracy it provides. Both good performance and low runtime are the key factors for a CDSS.

Moreover, we compared the accuracy provided by the proposed NE-nu-SVC model for the Z-Alizadeh Sani data with the results yielded by some recent studies dedicated to the analysis of this well-known CAD dataset. Table 9 presents a comprehensive comparison of the results generated by different classification models which were used to analyze the Z-Alizadeh Sani data. As reported in this table, our model has the highest accuracy (94.66%) among the competing approaches. It is worth mentioning that the proposed NE-nu-SVC model also provided very competitive results in terms of other metrics considered, including Recall, F-measure, ROC area, Kappa statistic, MAE and RMSE.

The time complexity of ensemble learning methods should be considered with pruning procedure and without pruning procedure [52]. Suppose that $m$ base algorithms are trained using an ensemble learning model. The total time complexity of training for the ensemble learning model without pruning procedure is $O(m \times t_{train})$, whereas the time complexity of prediction for unseen data (unknown instances) is $O(m \times t_{test})$ [52], where $t_{train}$

is average time required to train the model with one algorithm and $t_{test}$ is average time required to test the model with one algorithm. It should be noted that $t_{train}$ depends on two factors: 1) the size of training set, and 2) the specific base training algorithms being used, while $t_{test}$ depends on the specific machine learning algorithms being used. Table 8 provides the time complexity of the ensemble learning techniques and base machine learning algorithms used in our study.

In Table 8, $O(E_i)$ $(1 \leq i \leq k)$ is the time complexity of ensemble learning technique $E_i$, $k$ is the number of ensemble learning techniques used in stacking, $r$ is the number of records in the dataset and $v$ is the size of the adopted feature vector. In the SGD algorithm, $mtry$ is the number of used features, $ntree$ is the number of trees in Random Forest, $d$ is the depth of the Random Forest tree. Also, $T$ is the number of iterations and $t$ is the average running time of an individual machine learning algorithm used in Ne-nu-SVC, $T'$ is the number of trials, $h$ is the number of hypotheses for voting and $dl$ is the dictionary length. Thus, based on Table 8, we get the following time complexity for the proposed NE-nu-SVC model:

$$
\begin{aligned}
Time - Complexity&(NE - nu - SVC) = \\
& O(r \times v) + O(r) \\
& + O(ntree \times mtry \times d \times r) \\
& + O(T \times T' \times h \times t) \\
& + O(r^{2.3}) + O(dl \times r).
\end{aligned} \tag{10}
$$

### G. APPLICATION OF THE NE-NU-SVC MODEL TO THE CLEVELAND DATA

In order to confirm the effectiveness of the proposed methodology, the NE-nu-SVC model was also used to analyze the Cleveland CAD dataset [25], [41], [73]–[75].

As previously, the NE-nu-SVC model was used with the linear, polynomial, RBF and sigmoid kernel functions. We then applied the genetic search algorithm for feature selection, as explained in the previous section. As a result, seven original features of the Cleveland CAD dataset were selected for further analysis. These features are as follows: CP (Chest pain), Restecg (Results of resting electrocardiographic), Thalach (Maximum heart rate), Exang (Exercise

TABLE 7: Running times obtained for the Z-Alizadeh Sani dataset before and after using the Nested Ensemble method with different kernel functions within the nu-SVC model.

| Runtime (seconds) | Model | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | nu-SVC | | | | NE-nu-SVC | | | |
| | Linear | Polynomial | RBF | Sigmoid | Linear | Polynomial | RBF | Sigmoid |
| Time 1 | 1.65 | 1.99 | 0.04 | 0.01 | 0.61 | 0.64 | 0.64 | 0.59 |
| Time 2 | 1.38 | 2.07 | 0.04 | 0.01 | 0.66 | 0.67 | 0.65 | 0.59 |
| Time 3 | 1.42 | 2.09 | 0.04 | 0.01 | 0.62 | 0.66 | 0.65 | 0.63 |
| Time 4 | 1.39 | 2.08 | 0.04 | 0.01 | 0.64 | 0.65 | 0.65 | 0.60 |
| Time 5 | 1.38 | 2.12 | 0.04 | 0.01 | 0.63 | 0.65 | 0.65 | 0.59 |
| Time 6 | 1.38 | 2.08 | 0.04 | 0.01 | 0.64 | 0.64 | 0.65 | 0.60 |
| Time 7 | 1.39 | 2.05 | 0.04 | 0.01 | 0.63 | 0.65 | 0.65 | 0.61 |
| Time 8 | 1.38 | 2.06 | 0.04 | 0.01 | 0.62 | 0.65 | 0.66 | 0.59 |
| Time 9 | 1.61 | 2.05 | 0.04 | 0.01 | 0.62 | 0.65 | 0.66 | 0.60 |
| Time 10 | 1.42 | 2.03 | 0.04 | 0.01 | 0.74 | 0.65 | 0.66 | 0.59 |
| **Average** | 1.44 | 2.062 | 0.04 | **0.01** | 0.641 | 0.651 | 0.652 | **0.599** |

TABLE 8: Time complexity of all ensemble learning and base machine learning algorithms used in this study. See Section F for detailed information about the methods and variables being used.

| Algorithm | Study | Time complexity |
|---|---|---|
| Stacking | Zhao et al. (2018) [53] | $O(E_1 + E_2 + E_3 + ... + E_k)$ |
| nu-SVC (SVM) | Hsieh et al. (2016) [54] | $O(r \times v)$ |
| SGD | Chen et al. (2018) [55] | $O(r)$ |
| Random Forest | Gupta and Rana (2019) [56] | $O(ntree \times mtry \times d \times r)$ |
| Bagging | Zhao et al. (2018) [53] | $O(T \times t)$ |
| Vote | Mesterharm (2007) [57] | $O(T' \times h \times t)$ |
| SMO | Ouyang and Gray (2010) [58] | $O(r^{2.3})$ |
| Naïve Bayes | Jia et al. (2012) [59] | $O(dl \times r)$ |

TABLE 9: Comparison of the accuracy of the NE-nu-SVC model with state-of-art techniques for the Z-Alizadeh Sani CAD data.

| Study | Model | Accuracy in (%) |
|---|---|---|
| Alizadehsani et al. (2012) [60] | SMO | 82.16 |
| Alizadehsani et al. (2012) [61] | SMO 1-1 | 92.74 |
| Alizadehsani et al. (2012) [62] | SMO | 92.09 |
| Alizadehsani et al. (2012) [63] | Ensemble (Naïve Bayes-SMO) | 88.52 |
| Alizadehsani et al. (2013) [64] | Bagging-C4.5 | 79.54 (LAD) 61.46 (LCX) 68.96 (RCA) |
| Alizadehsani et al. (2013) [65] | Information gain-SMO | 94.08 |
| Yadav et al. (2014) [66] | Improved ARM | 93.75 |
| Alizadehsani et al. (2016) [27] | arteries-SVM Combined information gain for all | 86.14 (LAD) 83.17 (LCX) 83.50 (RCA) |
| Arabasadi et al. (2017) [31] | Neural network-genetic algorithm | 93.85 |
| Qin et al. (2017) [67] | EA-MFS | 93.70 |
| Babič et al. (2017) [34] | SVM | 86.67 |
| Hu et al. (2018) [68] | Variational finite inverted Beta-Liouville (IBL) Mixture Model (Var- IBLMM) | 81.84 |
| Kiliç and Kayakeles et al. (2018) [69] | Artificial Bee Colony+Sequential Minimal Optimization (ABCSMO) | 89.43 |
| Zhang et al. (2018) [70] | Extend correlation Restricted Boltzmann machine (Exp-CRBM) | 88.95±3.84 |
| Abdar et al. (2019) [71] | N2Genetic-nuSVM | 93.08 |
| Khan et al. (2019) [72] | Neural Network + Gini Index + Backward Weight Optimization | 88.49 |
| ***Proposed method*** | ***NE-nu-SVC + feature selection + multi-step balancing*** | ***94.66*** |

induced angina), Oldpeak (ST depression induced by exercises relevant to rest), Ca (Number of major vessels) and Thal. It should be noted that 6 original records included missing values. Therefore, they were eliminated from the dataset (i.e., 297 out of 303 original records were analyzed in our work). The obtained results are reported in Table 10 and Fig. 7.

The reduced dataset was balanced using the above-discussed multi-level balancing procedure. For the Cleveland data, we applied a five-step balancing in order to maximize the performance of NE-nu-SVC. The balancing procedure was carried out five times: one time using the ClassBalancer technique (supervised), three times using the resample technique (supervised) and once using the resample technique (unsupervised). The main reason for applying a five-step balancing was that a two-step balancing, which worked well with the Z-Alizadeh Sani data, led to a lower accuracy for the Cleveland data. The detailed results provided by NE-nu-SVC on the modified Cleveland data, obtained after applying the feature selection and five-step balancing of entities, are presented in Table 11 and Fig. 8. Table 12 reports the average running times of our program for the modified records of the Cleveland CAD dataset, obtained after applying the feature selection and five-step balancing of the entities when running the nu-SVC and NE-nu-SVC models.

A comparison of the accuracy provided by the proposed NE-nu-SVC model with state-of-art CAD detection methods for the modified Cleveland CAD data is presented in Table 13. It can be noted that our NE-nu-SVC model provided the highest accuracy (98.60%) among the competing methods. Even though our results are very promising, we cannot ignore the role of specialists in the diagnostic process. However, we can argue that this new model might be appropriate as an assistant during an implementation of clinical guidelines.

There are several advantages and disadvantages of the proposed model which we should report. For example, our model can be also used as a deep ensemble learning model which allows us to include in it different numbers of ensemble learning techniques and classical machine learning algorithms. Importantly, even though our new model comprises several algorithms at different levels, its running time remains reasonable. However, there are also some disadvantages of the proposed methodology which should be addressed in future work. First, the weights of classifiers have not been considered in this study. To do so, one has to apply different evolutionary algorithms to find the proper weights for classifiers at different levels. Moreover, the proposed model should be tested on different datasets, including big data. Finally, our study does not consider the impact of ECG signals [94]–[97] and ultrasound images [98], which should be investigated in the future.

## V. CONCLUSION

Nowadays, the impact of Clinical Decision Support Systems (CDSSs) on the individual's health increases gradually. Hence, the improvement of accuracy of statistical models included in CDSSs is a key challenge for clinical researchers, patients and physicians. Coronary artery disease (CAD), being one of the main causes of death worldwide, attracts valuable attention from many researchers worldwide. This study introduces a new hybrid ensemble learning model which can be used in the framework of a CDSS. The proposed model was tested on two well-known CAD datasets: the Z-Alizadeh Sani and Cleveland data from the UCI repository. Our model is a part of the Nested Ensemble (NE) approach. It relies on different traditional machine learning algorithms. In this study, the nu-SVC algorithm, including linear, polynomial, RBF and sigmoid kernels, was selected as the base algorithm of our NE-nu-SVC model. Within an NE model, nu-SVC was combined with other effective machine learning techniques such as SGD (Stochastic Gradient Descent), SMO (Sequential Minimal Optimization), random forest, Naïve Bayes, staking, bagging and voting. Furthermore, both the features selection and data balancing procedures were carried out to enhance the performance of the new model. We applied a genetic search algorithm for feature selection with both CAD datasets we analyzed (Z-Alizadeh Sani and Cleveland data). Since, these datasets were not well-balanced, we also carried out a multi-level balancing, using both ClassBlancer and Resample methods. The NE approach allows one to combine several ensemble learning techniques at different levels of the model. Here, we applied four ensemble learning techniques at three different levels. At the first level, the nu-SCV, SGD and random forest algorithms were combined using the stacking and bagging techniques. At the second level, the voting technique, and at the third level, the SMO and Naïve Bayes algorithms, were used. Our new model provided the accuracy of 94.66% for the Z-Alizadeh Sani data and of 98.60% for the modified Cleveland data. It allowed us to outperform the results provided by the existing machine learning algorithms for these well-known CAD datasets (see Tables 9 and 13). Moreover, we need to point out that the proposed NE-nu-SVC model is efficient in terms of running time. On average, over all four kernels, the execution of our program took 0.635 (s) for the Z-Alizadeh Sani dataset and 0.483 (s) for the Cleveland dataset. In the future, it would be interesting to apply the proposed model in the framework of other CDSSs aimed, for example, at the prediction of such important diseases as breast cancer or stroke. In this study, the number of levels in the NE model was selected manually. It would be essential to adapt an evolutionary algorithm (EA) to select the optimal number of levels automatically.

## APPENDIX A

TABLE 10: Experimental results provided by the proposed NE-nu-SVC model for the modified Cleveland CAD dataset after the feature selection and prior to data balancing.

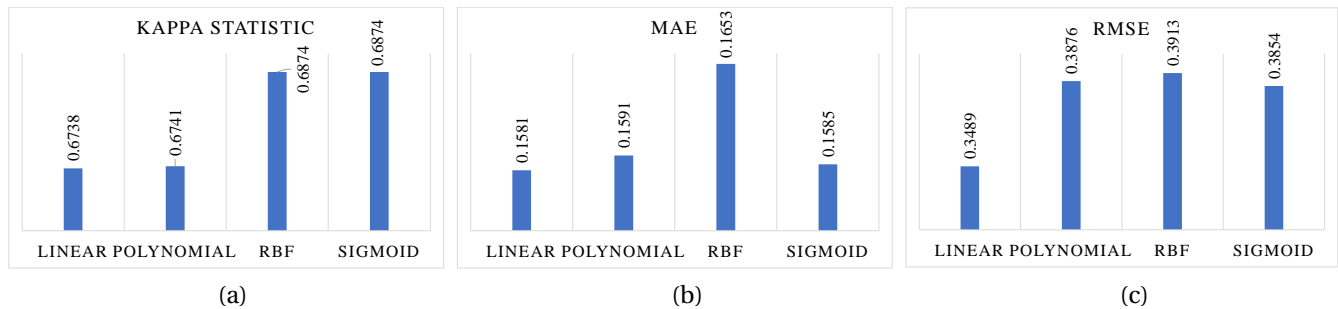| Measures | | NE-nu-SVC | | | |
|---|---|---|---|---|---|
| | | **Linear** | **Polynomial** | **RBF** | **Sigmoid** |
| FPR | Normal (%) | 19.70 | 19.00 | 19.00 | 19.00 |
| | CAD (%) | 13.10 | 13.80 | 12.50 | 12.50 |
| | Average (%) | 16.70 | 16.60 | 16.00 | 16.00 |
| Precision | Normal (%) | 83.70 | 84.10 | 84.30 | 87.50 |
| | CAD (%) | 84.00 | 83.50 | 84.70 | 81.00 |
| | Average (%) | 83.80 | 83.80 | 84.50 | 84.50 |
| Recall | Normal (%) | 86.90 | 86.30 | 87.50 | 85.90 |
| | CAD (%) | 80.30 | 81.00 | 81.00 | 82.80 |
| | Average (%) | 83.80 | 83.80 | 84.50 | 84.50 |
| F-measure | Normal (%) | 85.30 | 85.20 | 85.90 | 85.90 |
| | CAD (%) | 82.10 | 82.20 | 82.80 | 82.80 |
| | Average (%) | 83.80 | 83.80 | 84.50 | 84.50 |
| ROC Area | Normal (%) | 88.60 | 89.20 | 88.60 | 89.00 |
| | CAD (%) | 88.60 | 89.20 | 88.60 | 89.00 |
| | Average (%) | 88.60 | 89.20 | 88.60 | 89.00 |
| Accuracy(%) | | 83.84 | 83.83 | **84.51** | **84.51** |



FIGURE 7: Comparison of the results obtained using the proposed NE-nu-SVC model after the feature selection and prior to data balancing with different kernel functions when 7 selected features of the modified Cleveland CAD dataset were used: (a) Kappa statistic, (b) MAE and (c) RMSE.

TABLE 11: The results provided by the proposed NE-nu-SVC model for the modified Cleveland CAD data, obtained after applying feature selection and five-step balancing of the entities.

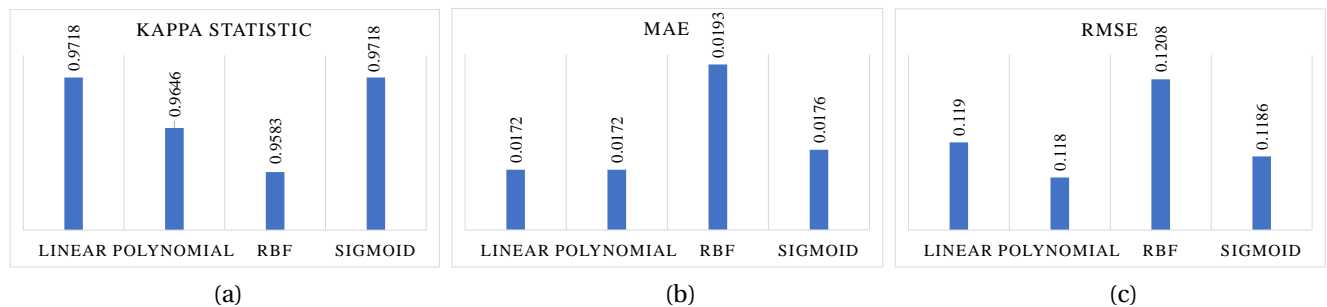| Measures | | NE-nu-SVC | | | |
|---|---|---|---|---|---|
| | | **Linear** | **Polynomial** | **RBF** | **Sigmoid** |
| FPR | Normal (%) | 2.00 | 2.60 | 2.60 | 2.00 |
| | CAD (%) | 0.70 | 0.70 | 1.40 | 0.70 |
| | Average (%) | 1.30 | 1.60 | 1.90 | 1.30 |
| Precision | Normal (%) | 97.60 | 96.90 | 96.80 | 97.60 |
| | CAD (%) | 99.40 | 99.40 | 98.90 | 99.40 |
| | Average (%) | 98.60 | 98.30 | 97.90 | 98.60 |
| Recall | Normal (%) | 99.30 | 99.30 | 98.60 | 99.30 |
| | CAD (%) | 98.00 | 97.40 | 97.40 | 98.00 |
| | Average (%) | 98.60 | 98.20 | 97.90 | 98.60 |
| F-measure | Normal (%) | 98.50 | 97.70 | 97.40 | 98.50 |
| | CAD (%) | 98.70 | 98.40 | 98.10 | 98.70 |
| | Average (%) | 98.60 | 98.20 | 97.90 | 98.60 |
| ROC Area | Normal (%) | 99.10 | 99.10 | 99.40 | 99.20 |
| | CAD (%) | 99.00 | 99.00 | 99.20 | 99.10 |
| | Average (%) | 99.00 | 99.00 | 99.30 | 99.10 |
| Accuracy(%) | | **98.60** | 98.24 | 97.93 | **98.60** |

FIGURE 8: Comparison of the results obtained using the proposed NE-nu-SVC model with different kernel functions when selected features of the Cleveland CAD dataset were considered (after applying feature selection and five-step balancing of the entities): (a) Kappa statistic, (b) MAE and (c) RMSE.

TABLE 12: Running times obtained using modified records of the Cleveland CAD dataset before and after using the Nested Ensemble (NE-nu-SVC) model with different kernel functions within nu-SVC.

| Runtime (seconds) | Models | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | nu-SVC | | | | NE-nu-SVC | | | |
| | Linear | Polynomial | RBF | Sigmoid | Linear | Polynomial | RBF | Sigmoid |
| Time 1 | 0.11 | 0.09 | 0.02 | 0.01 | 0.44 | 0.49 | 0.45 | 0.42 |
| Time 2 | 0.13 | 0.09 | 0.02 | 0.01 | 0.55 | 0.83 | 0.48 | 0.47 |
| Time 3 | 0.03 | 0.09 | 0.01 | 0.01 | 0.42 | 0.46 | 0.45 | 0.42 |
| Time 4 | 0.03 | 0.08 | 0.01 | 0.01 | 0.85 | 0.64 | 0.46 | 0.43 |
| Time 5 | 0.03 | 0.08 | 0.01 | 0.01 | 0.46 | 0.49 | 0.45 | 0.41 |
| Time 6 | 0.03 | 0.08 | 0.01 | 0.01 | 0.80 | 0.53 | 0.44 | 0.41 |
| Time 7 | 0.03 | 0.09 | 0.01 | 0.01 | 0.45 | 0.46 | 0.44 | 0.41 |
| Time 8 | 0.03 | 0.09 | 0.01 | 0.01 | 0.43 | 0.46 | 0.44 | 0.41 |
| Time 9 | 0.14 | 0.09 | 0.02 | 0.01 | 0.44 | 0.48 | 0.44 | 0.42 |
| Time 10 | 0.03 | 0.09 | 0.01 | 0.01 | 0.46 | 0.47 | 0.45 | 0.42 |
| Average | 0.059 | 0.087 | 0.013 | **0.01** | 0.530 | 0.531 | 0.450 | **0.422** |

TABLE 13: Comparison of the accuracy of the NE-nu-SVC model with state-of-art techniques for the Cleveland CAD dataset.

| Study | Model | Number of classes | Accuracy in (%) |
|---|---|---|---|
| Cheung (2001) [73] | Naive Bayes | N/A | 81.48 |
| Polat et al. (2005) [74] | AIRS | 2 | 84.50 |
| Polat et al. (2006) [75] | Fuzzy-AIRS-KNN-based system | 2 | 87.00 |
| Kahramanli and Allahverdi (2008) [76] | ANN and FNN (Fuzzy neural network) | 2 | 86.80 |
| Das et al. (2009) [25] | Neural networks ensembles | 2 | 89.01 |
| Polat et al. (2009) [77] | F-score approach for feature selection and LS-SVM with RBF kernel | 2 | 83.70 |
| Anooj et al. (2011) [78] | Weighted fuzzy rules | 2 | 86.35 |
| Srinivas et al. (2014) [79] | Rough-Fuzzy Classifier | 5 | 46.48 |
| Abdar (2015) [33] | C5.0 | 5 | 85.33 |
| El-Bialy et al. (2015) [80] | C4.5 | 2 | 78.54 |
| Elsayad and Fakhr, (2015) [81] | Markov Blanket Estimation (MBE) | 5 | 97.92 |
| Alizadehsani et al. (2017) [31] | Neural network and genetic algorithm | 2 | 89.40 |
| Paul et al. (2017) [82] | Weighted fuzzy system ensemble | 5 | 92.31 |
| Uyar and İlhan (2017) [83] | recurrent fuzzy neural networks (RFNN) and GA | 2 | 97.78 |
| Karayılan and Kılıç (2017) [84] | Multilayer Perceptron Neural Network | 2 | 95.55 |
| Alizadehsani et al. (2018) [85] | Feature engineering and SVM | 2 | 93.06 |
| Amin et al. (2018) [86] | Vote with Naive Bayes and logistic regression | 2 | 87.41 |
| Haq et al. (2018) [87] | Logistic regression after features selection | 2 | 89.00 |
| Gokulnath and Shantharajah (2018) [88] | Genetic algorithm with SVM | 2 | 88.34 |
| Khan et al. (2019) [72] | Neural Network and Gini Index and Backward Weight Optimization | 5 | 95.01 |
| Burse et al. (2019) [89] | Multi-Layer Pi-Sigma Neuron Model (MLPSNM) and Principal Component Analysis (PCA) | N/A | 94.53 |
| Rajab et al. (2019) [90] | Kernel-based FCM (KFCM)-based ANFIS | N/A | 86.00 |
| Terrada et al. (2019) [91] | ANN, KNN, K-means, and K-medoids | 2 | 96.01 |
| Ali et al. (2019) [92] | $L_1$ Linear SVM, $L_2$ Linear, and RBF SVM | 2 | 92.22 |
| Akgül et al. (2019) [93] | ANN-GA | 2 | 95.82 |
| ***Proposed method*** | ***NE-nu-SVC + feature selection + multi-step balancing*** | ***2*** | ***98.60*** |

TABLE A.1: Initial population features generated by the applied genetic algorithm for the Z-Alizadeh Sani CAD dataset.

| Merit | Scaled | Subset |
|---|---|---|
| 0.01958 | 0.02731 | 4, 5, 7, 9, 11, 12, 13, 20, 21, 23, 24, 27, 29, 30, 32, 34, 36, 37, 38, 39, 41, 43, 44, 46, 47, 49 |
| 0.04188 | 0.03687 | 1, 4, 6, 7, 9, 11, 13, 14, 16, 18, 21, 22, 23, 25, 26, 28, 32, 34, 36, 40, 41, 43, 44, 45, 46, 48, 49, 52, 53 |
| 0.04686 | 0.03901 | 5, 6, 7, 8, 11, 12, 14, 21, 23, 25, 26, 27, 28, 30, 31, 32, 34, 37, 38, 39, 45, 46, 47, 48, 49, 50, 52, 53, 54 |
| 0.03643 | 0.03453 | 1, 2, 5, 7, 8, 12, 15, 17, 18, 22, 23, 24, 25, 27, 29, 31, 33, 35, 38, 41, 42, 44, 46, 48, 49, 50, 51, 52, 53, 54, 55 |
| 0.11029 | 0.06621 | 7, 12, 13, 16, 18, 25, 26, 28, 34, 35, 37, 47, 50 |
| 0.03707 | 0.03481 | 1, 2, 4, 5, 7, 9, 10, 11, 12, 15, 16, 17, 19, 20, 21, 23, 24, 28, 29, 31, 32, 33, 35, 36, 37, 38, 39, 40, 45, 46, 49, 50, 53, 54 |
| 0.01659 | 0.02602 | 9, 10, 29, 44, 46 |
| 0.01996 | 0.02747 | 1, 2, 5, 7, 8, 9, 10, 12, 13, 14, 15, 17, 20, 21, 24, 27, 29, 30, 31, 33, 34, 35, 36, 37, 38, 43, 44, 46, 47, 49, 50, 51, 55 |
| 0.02493 | 0.02960 | 5, 12, 13, 16, 24, 26, 49, 52, 54 |
| 0.01767 | 0.02649 | 5, 11, 24, 55 |
| 0.04273 | 0.03724 | 4, 13, 18, 21, 23, 28, 29, 34, 35, 38, 42, 43, 46 |
| 0.02747 | 0.03069 | 1, 2, 3, 4, 5, 7, 9, 11, 12, 14, 15, 17, 18, 20, 21, 24, 28, 29, 30, 32, 33, 36, 39, 41, 42, 44, 45, 49, 50, 51, 52, 55 |
| 0.02911 | 0.03139 | 1, 6, 14, 20, 23, 27, 30, 36, 41, 48, 54 |
| 0.03895 | 0.03561 | 1, 4, 6 ,8 ,9, 10, 18, 19, 20, 21, 23, 25, 28, 29, 31, 32, 35, 40, 41, 42, 43, 44, 45, 46, 47, 48, 51, 52, 53, 55 |
| 0.02372 | 0.02908 | 4, 16, 17, 19, 21, 25, 38, 41, 42, 43, 44, 46, 49 |
| 0.00021 | 0.01900 | 11, 20 |
| 0.03611 | 0.03439 | 1, 4, 8, 9, 10, 11, 15, 17, 18, 19, 22, 24, 25, 29, 31, 32, 33, 35, 38, 40, 43, 45, 48, 49, 51 |
| 0.02587 | 0.03000 | 2, 4, 5, 8, 10, 11, 15, 19, 22, 23, 24, 27, 28, 32, 33, 35, 36, 37, 38, 39, 41, 42, 47, 48, 51, 52, 54, 55 |
| 0.02565 | 0.02991 | 2, 5, 6, 13, 16, 18, 19, 20, 21, 23, 24, 29, 31, 33, 36, 39, 40, 41, 42, 43, 47, 49, 50, 54 |
| 0.04104 | 0.03651 | 1, 3, 6, 10, 11, 13, 14, 15, 16, 18, 21, 27, 28, 34, 37, 39, 41, 44, 45, 46, 47, 50, 52, 53, 54, 55 |

TABLE A.2: Generated features provided by the applied genetic algorithm for the Z-Alizadeh Sani CAD dataset.

| | Merit | Scaled | Subset |
|---|---|---|---|
| | 0.24654 | 0.3294 | 6, 7, 12, 18, 25, 26, 28, 29, 32, 35, 37, 53, 54, 55 |
| | 0.24654 | 0.3294 | 6, 7, 12, 18, 25, 26, 28, 29, 32, 35, 37, 53, 54, 55 |
| | 0.22603 | 0.28407 | 7, 12, 25, 26, 28, 29, 32, 34, 37, 38, 53, 54 |
| | 0.09750 | 0.0000 | 7, 12, 18, 24, 25, 26, 27, 28, 34, 36, 37, 38, 48, 53, 54, 55 |
| | 0.24626 | 0.32878 | 6, 7, 12, 18, 25, 26, 28, 32, 35, 47, 53, 54 |
| | 0.10218 | 0.01034 | 7, 12, 16, 25, 26, 28, 29, 37, 38, 41, 49, 53, 54, 55 |
| | 0.10619 | 0.01920 | 7, 12, 18, 24, 25, 26, 28, 34, 36, 37, 38, 40, 45, 55 |
| | 0.14178 | 0.09785 | 6, 7, 12, 18, 20, 21, 25, 26, 28, 29, 31, 37, 53, 54, 55 |
| | 0.09832 | 0.00180 | 4, 6, 7, 10, 12, 18, 19, 25, 26, 28, 29, 32, 45, 55 |
| [H] | 0.22071 | 0.27232 | 7, 12, 24, 25, 26, 28, 29, 37, 38, 53, 54, 55 |
| | 0.24654 | 0.32940 | 6, 7, 12, 18, 25, 26, 28, 29, 32, 35, 37, 53, 54, 55 |
| | 0.10440 | 0.01525 | 5, 7, 12, 18, 22, 24, 25, 26, 28, 29, 37, 38, 42, 53, 54, 55 |
| | 0.23369 | 0.30101 | 7, 12, 25, 26, 28, 29, 32, 35, 37, 38, 53, 54 |
| | 0.11147 | 0.03087 | 6, 7, 12, 18, 25, 28, 29, 32, 35, 37, 40, 52, 53, 55 |
| | 0.20152 | 0.22990 | 6, 7, 12, 18, 24, 25, 26, 28, 29, 34, 37, 55 |
| | 0.13983 | 0.09355 | 4, 6, 7, 8, 12, 15, 16, 18, 25, 26, 28, 29, 32, 35, 37, 53, 54, 55 |
| | 0.24940 | 0.33572 | 6, 7, 12, 18, 25, 26, 28, 29, 32, 35, 53, 54, 55 |
| | 0.24654 | 0.32940 | 6, 7, 12, 18, 25, 26, 28, 29, 32, 35, 37, 53, 54, 55 |
| | 0.15224 | 0.12099 | 6, 7, 12, 18, 25, 26, 28, 29, 32, 35, 37, 44, 53, 54, 55 |
| | 0.14370 | 0.10211 | 6, 7, 12, 17, 19, 25, 26, 28, 29, 32, 35, 37, 53, 54, 55 |

## REFERENCES

[1] Y. Wang, L. Kung, W. Y. C. Wang, and C. G. Cegielski, "An integrated big data analytics-enabled transformation model: Application to health care," Information & Management, vol. 55, no. 1, pp. 64–79, 2018.

[2] M. Abdar, M. Zomorodi-Moghadam, R. Das, and I.-H. Ting, "Performance analysis of classification algorithms on early detection of liver disease," Expert Systems with Applications, vol. 67, pp. 239–251, 2017.

[3] M. Abdar and N. Y. Yen, "Design of a universal user model for dynamic crowd preference sensing and decision-making behavior analysis," IEEE Access, vol. 5, pp. 24 842–24 852, 2017.

[4] H.-H. Yang, M.-L. Huang, C.-M. Lai, and J.-R. Jin, "An approach combining data mining and control charts-based model for fault detection in wind turbines," Renewable energy, vol. 115, pp. 808–816, 2018.

[5] G. Sun, C. Jiang, P. Cheng, Y. Liu, X. Wang, Y. Fu, and Y. He, "Short-term wind power forecasts by a synthetical similar time series data mining method," Renewable energy, vol. 115, pp. 575–584, 2018.

[6] Y. Cong, G. Sun, J. Liu, H. Yu, and J. Luo, "User attribute discovery with missing labels," Pattern Recognition, vol. 73, pp. 33–46, 2018.

[7] G. Sun, Y. Cong, and X. Xu, "Active lifelong learning with" watchdog"," in Thirty-Second AAAI Conference on Artificial Intelligence, 2018.

[8] G. Sun, Y. Cong, D. Hou, H. Fan, X. Xu, and H. Yu, "Joint household characteristic prediction via smart meter data," IEEE Transactions on Smart Grid, vol. 10, no. 2, pp. 1834–1844, 2017.

[9] H. Kashyap, H. A. Ahmed, N. Hoque, S. Roy, and D. K. Bhattacharyya, "Big data analytics in bioinformatics: A machine learning perspective," arXiv preprint arXiv:1506.05101, 2015.

[10] L. Chen, T. Huang, C. Lu, L. Lu, and D. Li, "Machine learning and network methods for biology and medicine," Computational and mathematical methods in medicine, vol. 2015, 2015.

[11] M. A. Karaolis, J. A. Moutiris, D. Hadjipanayi, and C. S. Pattichis, "Assessment of the risk factors of coronary heart events based on data mining with decision trees," IEEE Transactions on information technology in biomedicine, vol. 14, no. 3, pp. 559–566, 2010.

[12] W. Książek, M. Abdar, U. R. Acharya, and P. Pławiak, "A novel machine learning approach for early detection of hepatocellular carcinoma patients," Cognitive Systems Research, vol. 54, pp. 116–127, 2019.

[13] P. Pławiak, "Novel methodology of cardiac health recognition based on ecg signals and evolutionary-neural system," Expert Systems with Applications, vol. 92, pp. 334–349, 2018.

[14] P. Plawiak, "Novel genetic ensembles of classifiers applied to my-

ocardium dysfunction recognition based on ecg signals," Swarm and evolutionary computation, vol. 39, pp. 192–208, 2018.

[15] D. Roffman, G. Hart, M. Girardi, C. J. Ko, and J. Deng, "Predicting non-melanoma skin cancer via a multi-parameterized artificial neural network," Scientific reports, vol. 8, no. 1, p. 1701, 2018.

[16] D. Medved, M. Ohlsson, P. Höglund, B. Andersson, P. Nugues, and J. Nilsson, "Improving prediction of heart transplantation outcome using deep learning techniques," Scientific reports, vol. 8, no. 1, p. 3613, 2018.

[17] H. Li, J. Peng, Y. Leung, K.-S. Leung, M.-H. Wong, G. Lu, and P. Ballester, "The impact of protein structure and sequence similarity on the accuracy of machine-learning scoring functions for binding affinity prediction," Biomolecules, vol. 8, no. 1, p. 12, 2018.

[18] C. N. Magnan and P. Baldi, "Sspro/accpro 5: almost perfect prediction of protein secondary structure and relative solvent accessibility using profiles, machine learning and structural similarity," Bioinformatics, vol. 30, no. 18, pp. 2592–2597, 2014.

[19] E. S. Berner, Clinical decision support systems. Springer, 2007, vol. 233.

[20] S. Harjai and S. K. Khatri, "An intelligent clinical decision support system based on artificial neural network for early diagnosis of cardiovascular diseases in rural areas," in 2019 Amity International Conference on Artificial Intelligence (AICAI). IEEE, 2019, pp. 729–736.

[21] A. Sheikhtaheri, A. Orooji, A. Pazouki, and M. Beitollahi, "A clinical decision support system for predicting the early complications of one-anastomosis gastric bypass surgery," Obesity surgery, pp. 1–11, 2019.

[22] E. Iadanza, V. Mudura, P. Melillo, and M. Gherardelli, "An automatic system supporting clinical decision for chronic obstructive pulmonary disease," Health and Technology, pp. 1–12, 2019.

[23] E. Sanchez, C. Toro, E. Carrasco, P. Bonachela, C. Parra, G. Bueno, and F. Guijarro, "A knowledge-based clinical decision support system for the diagnosis of alzheimer disease," in 2011 IEEE 13th International Conference on e-Health Networking, Applications and Services. IEEE, 2011, pp. 351–357.

[24] H. A. L. Rocha, I. Dankwa-Mullan, S. F. Juacaba, V. Willis, Y. E. Arriaga, G. P. Jackson, and P. Meneleu, "Shared-decision making in prostate cancer with clinical decision-support." 2019.

[25] R. Das, I. Turkoglu, and A. Sengur, "Effective diagnosis of heart disease through neural networks ensembles," Expert systems with applications, vol. 36, no. 4, pp. 7675–7680, 2009.

[26] J. Nahar, T. Imam, K. S. Tickle, and Y.-P. P. Chen, "Association rule mining to detect factors which contribute to heart disease in males and females," Expert Systems with Applications, vol. 40, no. 4, pp. 1086–1093, 2013.

[27] R. Alizadehsani, M. H. Zangooei, M. J. Hosseini, J. Habibi, A. Khosravi, M. Roshanzamir, F. Khozeimeh, N. Sarrafzadegan, and S. Nahavandi, "Coronary artery disease detection using computational intelligence methods," Knowledge-Based Systems, vol. 109, pp. 187–197, 2016.

[28] J. Li, S. Fong, S. Mohammed, and J. Fiaidhi, "Improving the classification performance of biological imbalanced datasets by swarm optimization algorithms," The Journal of Supercomputing, vol. 72, no. 10, pp. 3708–3728, 2016.

[29] E. Frank, M. Hall, G. Holmes, R. Kirkby, B. Pfahringer, I. H. Witten, and L. Trigg, "Weka-a machine learning workbench for data mining," in Data mining and knowledge discovery handbook. Springer, 2009, pp. 1269–1277.

[30] M. Tayefi, M. Tajfard, S. Saffar, P. Hanachi, A. R. Amirabadizadeh, H. Esmaeily, A. Taghipour, G. A. Ferns, M. Moohebati, and M. Ghayour-Mobarhan, "hs-crp is strongly associated with coronary heart disease (chd): A data mining approach using decision tree algorithm," Computer methods and programs in biomedicine, vol. 141, pp. 105–109, 2017.

[31] Z. Arabasadi, R. Alizadehsani, M. Roshanzamir, H. Moosaei, and A. A. Yarifard, "Computer aided decision making for heart disease detection using hybrid neural network-genetic algorithm," Computer methods and programs in biomedicine, vol. 141, pp. 19–26, 2017.

[32] A. H. Alkeshuosh, M. Z. Moghadam, I. Al Mansoori, and M. Abdar, "Using pso algorithm for producing best rules in diagnosis of heart disease," in 2017 international conference on computer and applications (ICCA). IEEE, 2017, pp. 306–311.

[33] M. Abdar, "Using decision trees in data mining for predicting factors influencing of heart disease." Carpathian Journal of Electronic & Computer Engineering, vol. 8, no. 2, 2015.

[34] F. Babič, J. Olejár, Z. Vantová, and J. Paralič, "Predictive and descriptive analysis for heart disease diagnosis," in 2017 Federated Conference on Computer Science and Information Systems (FedCSIS). IEEE, 2017, pp. 155–163.

[35] K. Polat, S. Şahan, and S. Güneş, "Automatic detection of heart disease using an artificial immune recognition system (airs) with fuzzy resource allocation mechanism and k-nn (nearest neighbour) based weighting preprocessing," Expert Systems with Applications, vol. 32, no. 2, pp. 625–631, 2007.

[36] U. R. Acharya, H. Fujita, M. Adam, O. S. Lih, V. K. Sudarshan, T. J. Hong, J. E. Koh, Y. Hagiwara, C. K. Chua, C. K. Poo et al., "Automated characterization and classification of coronary artery disease and myocardial infarction by decomposition of ecg signals: A comparative study," Information Sciences, vol. 377, pp. 17–29, 2017.

[37] S. Patidar, R. B. Pachori, and U. R. Acharya, "Automated diagnosis of coronary artery disease using tunable-q wavelet transform applied on heart rate signals," Knowledge-Based Systems, vol. 82, pp. 1–10, 2015.

[38] N. Kausar, A. Abdullah, B. B. Samir, S. Palaniappan, B. S. AlGhamdi, and N. Dey, "Ensemble clustering algorithm with supervised classification of clinical data for early diagnosis of coronary artery disease," Journal of Medical Imaging and Health Informatics, vol. 6, no. 1, pp. 78–87, 2016.

[39] Z. Mahmoodabadi and S. S. Tabrizi, "A new ica-based algorithm for diagnosis of coronary artery disease," in Intelligent Computing, Communication and Devices. Springer, 2015, pp. 415–427.

[40] Z. A. Sani, R. Alizadehsani, and M. Roshanzamir, "Z-alizadeh sani data set," https://archive.ics.uci.edu/ml/datasets/Z-Alizadeh+Sani, January 2018.

[41] D. W. Aha, "Heart disease data set, 4 databases: Cleveland, hungary, switzerland, and the va long beach," https://archive.ics.uci.edu/ml/datasets/heart+disease, January 2018.

[42] D. L. Mann, D. P. Zipes, P. Libby, and R. O. Bonow, Braunwald's heart disease e-book: a textbook of cardiovascular medicine. Elsevier Health Sciences, 2014.

[43] M. Abdar, M. Zomorodi-Moghadam, X. Zhou, R. Gururajan, X. Tao, P. D. Barua, and R. Gururajan, "A new nested ensemble technique for automated diagnosis of breast cancer," Pattern Recognition Letters, 2018.

[44] H. Byliński, A. Sobecki, and J. Gębicki, "The use of artificial neural networks and decision trees to predict the degree of odor nuisance of post-digestion sludge in the sewage treatment plant process," Sustainability, vol. 11, no. 16, p. 4407, 2019.

[45] T. Leathart, E. Frank, B. Pfahringer, and G. Holmes, "Ensembles of nested dichotomies with multiple subset evaluation," in Pacific-Asia Conference on Knowledge Discovery and Data Mining. Springer, 2019, pp. 81–93.

[46] H. R. Roth, L. Lu, N. Lay, A. P. Harrison, A. Farag, A. Sohn, and R. M. Summers, "Spatial aggregation of holistically-nested convolutional neural networks for automated pancreas localization and segmentation," Medical image analysis, vol. 45, pp. 94–107, 2018.

[47] C.-C. Chang and C.-J. Lin, "Libsvm: A library for support vector machines," ACM transactions on intelligent systems and technology (TIST), vol. 2, no. 3, p. 27, 2011.

[48] M. Abdar and M. Zomorodi-Moghadam, "Impact of patients' gender on parkinson's disease using classification algorithms," Journal of AI and Data Mining, vol. 6, no. 2, pp. 277–285, 2018.

[49] D. E. Goldberg, Genetic algorithms in search, optimization and machine learning. Addison-Wesley, 1989.

[50] M. A. Hall, "Correlation-based feature selection for machine learning," 1999.

[51] T. Hasanin, T. M. Khoshgoftaar, J. L. Leevy, and N. Seliya, "Examining characteristics of predictive models with imbalanced big data," Journal of Big Data, vol. 6, no. 1, p. 69, 2019.

[52] Z. Ma and Q. Dai, "Selected an stacking elms for time series prediction," Neural Processing Letters, vol. 44, no. 3, pp. 831–856, 2016.

[53] Z. Zhao, M. Karimzadeh, F. Gerber, and T. Braun, "Mobile crowd location prediction with hybrid features using ensemble learning," Future Generation Computer Systems, 2018.

[54] C.-C. Hsieh, M.-H. Hsih, M.-K. Jiang, Y.-M. Cheng, and E.-H. Liang, "Effective semantic features for facial expressions recognition using svm," Multimedia Tools and Applications, vol. 75, no. 11, pp. 6663–6682, 2016.

[55] X. Chen, J. D. Lee, X. T. Tong, and Y. Zhang, "Statistical inference for model parameters in stochastic gradient descent," arXiv preprint arXiv:1610.08637, 2016.

[56] V. K. Gupta and P. S. Rana, "Activity assessment of small drug molecules in estrogen receptor using multilevel prediction model," IET systems biology, vol. 13, no. 3, pp. 147–158, 2019.

[57] C. Mesterharm, "Improving on-line learning," Ph.D. dissertation, Rutgers University-Graduate School-New Brunswick, 2007.
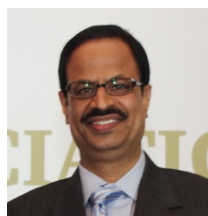
[58] H. Ouyang and A. Gray, "Fast stochastic frank-wolfe algorithms for non-linear svms," in Proceedings of the 2010 SIAM International Conference on Data Mining. SIAM, 2010, pp. 245–256.

[59] Z. Jia, R. Zhou, C. Zhu, L. Wang, W. Gao, Y. Shi, J. Zhan, and L. Zhang, "The implications of diverse applications and scalable data sets in benchmarking big data systems," in Specifying Big Data Benchmarks. Springer, 2012, pp. 44–59.

[60] R. Alizadehsani, J. Habibi, Z. A. Sani, H. Mashayekhi, R. Boghrati, A. Ghandeharioun, and B. Bahadorian, "Diagnosis of coronary artery disease using data mining based on lab data and echo features," Journal of Medical and Bioengineering, vol. 1, no. 1, 2012.

[61] R. Alizadehsani, M. J. Hosseini, R. Boghrati, A. Ghandeharioun, F. Khozeimeh, and Z. A. Sani, "Exerting cost-sensitive and feature creation algorithms for coronary artery disease diagnosis," International Journal of Knowledge Discovery in Bioinformatics (IJKDB), vol. 3, no. 1, pp. 59–79, 2012.

[62] R. Alizadehsani, M. J. Hosseini, Z. A. Sani, A. Ghandeharioun, and R. Boghrati, "Diagnosis of coronary artery disease using cost-sensitive algorithms," in 2012 IEEE 12th International Conference on Data Mining Workshops. IEEE, 2012, pp. 9–16.

[63] R. Alizadehsani, J. Habibi, M. J. Hosseini, R. Boghrati, A. Ghandeharioun, B. Bahadorian, and Z. A. Sani, "Diagnosis of coronary artery disease using data mining techniques based on symptoms and ecg features," European Journal of Scientific Research, vol. 82, no. 4, pp. 542–553, 2012.

[64] R. Alizadehsani, J. Habibi, Z. A. Sani, H. Mashayekhi, R. Boghrati, A. Ghandeharioun, F. Khozeimeh, and F. Alizadeh-Sani, "Diagnosing coronary artery disease via data mining algorithms by considering laboratory and echocardiography features," Research in cardiovascular medicine, vol. 2, no. 3, p. 133, 2013.

[65] R. Alizadehsani, J. Habibi, M. J. Hosseini, H. Mashayekhi, R. Boghrati, A. Ghandeharioun, B. Bahadorian, and Z. A. Sani, "A data mining approach for diagnosis of coronary artery disease," Computer methods and programs in biomedicine, vol. 111, no. 1, pp. 52–61, 2013.

[66] C. Yadav, S. Lade, and M. K. Suman, "Predictive analysis for the diagnosis of coronary artery disease using association rule mining," International Journal of Computer Applications, vol. 87, no. 4, 2014.

[67] C.-J. Qin, Q. Guan, and X.-P. Wang, "Application of ensemble algorithm integrating multiple criteria feature selection in coronary heart disease detection," Biomedical Engineering: Applications, Basis and Communications, vol. 29, no. 06, p. 1750043, 2017.

[68] C. Hu, W. Fan, J.-X. Du, and N. Bouguila, "A novel statistical approach for clustering positive data based on finite inverted beta-liouville mixture models," Neurocomputing, vol. 333, pp. 110–123, 2019.

[69] Ü. Kılıç and M. K. Keleş, "Feature selection with artificial bee colony algorithm on z-alizadeh sani dataset," in 2018 Innovations in Intelligent Systems and Applications Conference (ASYU). IEEE, 2018, pp. 1–3.

[70] N. Zhang, S. Ding, H. Liao, and W. Jia, "Multimodal correlation deep belief networks for multi-view classification," Applied Intelligence, vol. 49, no. 5, pp. 1925–1936, 2019.

[71] M. Abdar, W. Książek, U. R. Acharya, R.-S. Tan, V. Makarenkov, and P. Pławiak, "A new machine learning technique for an accurate diagnosis of coronary artery disease," Computer Methods and Programs in Biomedicine, p. 104992, 2019.

[72] Y. Khan, U. Qamar, M. Asad, and B. Zeb, "Applying feature selection and weight optimization techniques to enhance artificial neural network for heart disease diagnosis," in Proceedings of SAI Intelligent Systems Conference. Springer, 2019, pp. 340–351.

[73] N. Cheung, "Machine learning techniques for medical analysis," School of Information Technology and Electrical Engineering, 2001.

[74] K. Polat, S. Sahan, H. Kodaz, and S. Güneş, "A new classification method to diagnosis heart disease: supervised artificial immune system (airs)," in proceedings of the turkish symposium on artificial intelligence and neural networks (TAINN), 2005.

[75] K. Polat, S. Güneş, and S. Tosun, "Diagnosis of heart disease using artificial immune recognition system and fuzzy weighted pre-processing," Pattern Recognition, vol. 39, no. 11, pp. 2186–2193, 2006.

[76] H. Kahramanli and N. Allahverdi, "Design of a hybrid system for the diabetes and heart diseases," Expert systems with applications, vol. 35, no. 1-2, pp. 82–89, 2008.

[77] K. Polat and S. Güneş, "A new feature selection method on classification of medical datasets: Kernel f-score feature selection," Expert Systems with Applications, vol. 36, no. 7, pp. 10 367–10 373, 2009.

[78] P. Anooj, "Clinical decision support system: risk level prediction of heart disease using weighted fuzzy rules and decision tree rules," Open Computer Science, vol. 1, no. 4, pp. 482–498, 2011.

[79] K. Srinivas, G. R. Rao, and A. Govardhan, "Rough-fuzzy classifier: a system to predict the heart disease by blending two different set theories," Arabian Journal for Science and Engineering, vol. 39, no. 4, pp. 2857–2868, 2014.

[80] R. El-Bialy, M. A. Salamay, O. H. Karam, and M. E. Khalifa, "Feature analysis of coronary artery heart disease data sets," Procedia Computer Science, vol. 65, pp. 459–468, 2015.

[81] A. M. Elsayad and M. Fakhr, "Diagnosis of cardiovascular diseases with bayesian classifiers." JCS, vol. 11, no. 2, pp. 274–282, 2015.

[82] A. K. Paul, P. C. Shill, M. R. I. Rabin, and K. Murase, "Adaptive weighted fuzzy rule-based system for the risk level assessment of heart disease," Applied Intelligence, vol. 48, no. 7, pp. 1739–1756, 2018.

[83] K. Uyar and A. İlhan, "Diagnosis of heart disease using genetic algorithm based trained recurrent fuzzy neural networks," Procedia computer science, vol. 120, pp. 588–593, 2017.

[84] T. Karaylan and Ö. Kılıç, "Prediction of heart disease using neural network," in 2017 International Conference on Computer Science and Engineering (UBMK). IEEE, 2017, pp. 719–723.

[85] R. Alizadehsani, M. J. Hosseini, A. Khosravi, F. Khozeimeh, M. Roshanzamir, N. Sarrafzadegan, and S. Nahavandi, "Non-invasive detection of coronary artery disease in high-risk patients based on the stenosis prediction of separate coronary arteries," Computer methods and programs in biomedicine, vol. 162, pp. 119–127, 2018.

[86] M. S. Amin, Y. K. Chiam, and K. D. Varathan, "Identification of significant features and data mining techniques in predicting heart disease," Telematics and Informatics, vol. 36, pp. 82–93, 2019.

[87] A. U. Haq, J. P. Li, M. H. Memon, S. Nazir, and R. Sun, "A hybrid intelligent system framework for the prediction of heart disease using machine learning algorithms," Mobile Information Systems, vol. 2018, 2018.

[88] C. B. Gokulnath and S. Shantharajah, "An optimized feature selection based on genetic approach and support vector machine for heart disease," Cluster Computing, pp. 1–11, 2018.

[89] K. Burse, V. P. S. Kirar, A. Burse, and R. Burse, "Various preprocessing methods for neural network based heart disease prediction," in Smart Innovations in Communication and Computational Sciences. Springer, 2019, pp. 55–65.

[90] W. Rajab, S. Rajab, and V. Sharma, "Kernel fcm-based anfis approach to heart disease prediction," in Emerging Trends in Expert Applications and Security. Springer, 2019, pp. 643–650.

[91] O. Terrada, B. Cherradi, A. Raihani, and O. Bouattane, "Classification and prediction of atherosclerosis diseases using machine learning algorithms," in 2019 5th International Conference on Optimization and Applications (ICOA). IEEE, 2019, pp. 1–5.

[92] L. Ali, A. Niamat, J. A. Khan, N. A. Golilarz, X. Xingzhong, A. Noor, R. Nour, and S. A. C. Bukhari, "An optimized stacked support vector machines based expert system for the effective prediction of heart failure," IEEE Access, vol. 7, pp. 54 007–54 014, 2019.

[93] M. Akgül, Ö. E. Sönmez, and T. Özcan, "Diagnosis of heart disease using an intelligent method: A hybrid ann–ga approach," in International Conference on Intelligent and Fuzzy Systems. Springer, 2019, pp. 1250–1257.

[94] M. Sharma and U. R. Acharya, "A new method to identify coronary artery disease with ecg signals and time-frequency concentrated antisymmetric biorthogonal wavelet filter bank," Pattern Recognition Letters, vol. 125, pp. 235–240, 2019.

[95] U. R. Acharya, Y. Hagiwara, J. E. W. Koh, S. L. Oh, J. H. Tan, M. Adam, and R. San Tan, "Entropies for automated detection of coronary artery disease using ecg signals: A review," Biocybernetics and Biomedical Engineering, vol. 38, no. 2, pp. 373–384, 2018.

[96] U. R. Acharya, H. Fujita, O. S. Lih, M. Adam, J. H. Tan, and C. K. Chua, "Automated detection of coronary artery disease using different durations of ecg segments with convolutional neural network," Knowledge-Based Systems, vol. 132, pp. 62–71, 2017.

[97] M. Kumar, R. B. Pachori, and U. R. Acharya, "Characterization of coronary artery disease using flexible analytic wavelet transform applied on ecg signals," Biomedical signal processing and control, vol. 31, pp. 301–308, 2017.

[98] U. Raghavendra, H. Fujita, A. Gudigar, R. Shetty, K. Nayak, U. Pai, J. Samanth, and U. R. Acharya, "Automated technique for coronary artery disease characterization and classification using dd-dtdwt in ul-

trasound images," Biomedical Signal Processing and Control, vol. 40, pp. 324–334, 2018.

**MOLOUD ABDAR** received the bachelor's degree in computer engineering from Damghan University, Iran, in 2015. He also received the master's degree in computer science and engineering from the University of Aizu, Aizu, Japan, in 2018. Currently, he is working toward his Ph.D. degree at the University of Quebec in Montreal, Montreal (QC), Canada. He has written several papers in the fields of data mining, machine learning, and user modeling in some refereed international journals and conferences.

He is also very active in several international conferences (e. g., TPC in IEEE AINA 2018, IEEE AINA 2019 and IEEE AINA 2020) and some referred international journals, such as the IEEE Access, the Applied Soft Computing, the Future Generation Computer Systems (outstanding reviewer in October 2017), the Neurocomputing (outstanding reviewer in January 2017), the Neural Computing and Applications, etc., as a Reviewer. His research interests include data mining, machine learning, ensemble learning, evolutionary algorithms and user modeling. He is a recipient of the Fonds de Recherche du Quebec—Nature et Technologies Award (ranked 5th among 20 candidates in the second round of selection process) in 2019.

**U. RAJENDRA ACHARYA** received the Ph.D. degree from the National Institute of Technology Karnataka, Surathkal, India, and the D.Eng. degree from Chiba University, Japan. He is currently a Senior Faculty Member with Ngee Ann Polytechnic, Singapore. He is also an Adjunct Professor with Taylor's University, Malaysia, an Adjunct Faculty with the Singapore Institute of Technology–University of Glasgow, Singapore, and an Associate Faculty with the Singapore University of Social Sciences, Singapore. He has published more than 400 papers, in refereed international SCI-IF journals (345), international conference proceedings (42), books (17) with more than 24,500 citations in Google Scholar (with h-index of 83), and ResearchGate RG Score of 47.66.

He is ranked in the top 1% of the Highly Cited Researchers for the last three consecutive years (2016, 2017, and 2018) in computer science according to the Essential Science Indicators of Thomson. He has worked on various funded projects, with grants worth more than 2 million SGD. He has three patents and an editorial board member of many journals. His research interests include biomedical signal processing, biomedical imaging, data mining, visualization and biophysics for better healthcare design, delivery, and therapy. He has served as a Guest Editor for many journals.
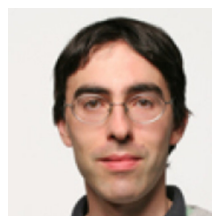
**NIZAL SARRAFZADEGAN** is a professor of Internal Medicine & Cardiology, Isfahan University of medical Sciences in Iran, and Affiliate Professor of the Faculty of Medicine, School of Population and Public Health (SPPH) in the University of British Columbian (UBC) in Vancouver, Canada. She is currently teaching a post graduate course on "NCD Epidemiology and Prevention" in the SPPH. She is the founder director of Isfahan Cardiovascular Research Institute (ICRI), a WHO collaborating center in the Eastern Meditreanian region (EMR).

She started her national and international studies on Cardiovascular Diseases (CVD) prevention, healthy lifestyle promotion and rehabilitation of cardiac patients since 1992 and published more than 450 articles and chapter books in peer-reviewed journals. She took part in more than 80 international meetings as international invited speaker or WHO advisor. She is the founder and President of the "Eastren Meditreanian Network on Heart Health", "Iranin Heart Foundation", founder and Co-Chair of the "National Network of CVD research" and the "Food, Industry and Healthy Community". She won the 2016 WHO/EMR award for her extensive research in CVD.

**VLADIMIR MAKARENKOV** is a Full professor and Director of the graduate Bioinformatics program at the Department of Computer Science of the Université du Québec à Montréal (Canada). His research interests are in the fields of bioinformatics, data mining and software engineering.

• • •