Research Article

Performance evaluation of optimal feature selectionbased machine learning for heart disease diagnosis

Doi: 10.30508/kdip.2023.400967.1072

Mohamad Hasanvand^{1*}, Arezu Selyari², Hamideh Jashn³, Zeinab Ghaseminejad⁴, Mahdi Nooshyar⁵

Abstract

Today, heart disease is one of the main causes of morbidity and mortality. Since the initial and early diagnosis of this disease is very important and vital and the usual methods used in the medical industry need to spend a lot of time and money to diagnose this disease, accurate prediction of this disease has become a challenge. According to the huge amount of hospital data, which is added to its volume every day, the importance of data mining, which is one of the important techniques for discovering knowledge and hidden patterns, is increasing. Many studies have been done based on data mining to predict heart disease; according to their solution, each one pursues goals such as increasing speed, increasing accuracy, reducing the volume of calculations, and error coefficient. This research aims to increase the reliability and accuracy of heart disease diagnosis using the feature selection technique by meta-heuristic algorithms to extract useful features and to reduce the computational burden, and we use machine-learning algorithms to evaluate the proposed method. Based on the obtained results, the proposed system diagnoses people suffering from cardiovascular disease with relatively high accuracy and precision.

Keywords: Machine Learning; Feature Selection; Heart Patients; Meta-heuristic; Data Mining.

¹⁻ Msc, Faculty of Computer Engineering, University of Mohagegh Ardebili, Ardebil, Iran.

²⁻ Msc, Faculty of Computer Engineering, Islamic Azad University, Gorgan Branch, Iran.

³⁻ MSc, Faculty of Computer Engineering, Science and Research Branch, Islamic Azad University, Khuzestan, Iran.

⁴⁻ MSc, Faculty of Computer Engineering, Vali-e-Asr University, Rafsanhan, Iran.

⁵⁻ Assistant Prof, Faculty of Computer Engineering, University of Mohagegh Ardebili, Ardebil, Iran.

1- Introduction

Healthcare is one of the most important issues that have been considered by data mining experts in recent years. One of the important challenges of medical organizations is the quality of their services against the affordable cost for the users of these services. The quality of services depends on the correct diagnosis and prescription of correct treatment and poor clinical decisions can result in undesirable results, which may have unexpected outcomes. The quality of services includes proper identification of the disease and its effective treatment. Poor clinical decisions can lead to serious and unacceptable consequences, as well as hospitals should minimize the cost of clinical trials. By using computer-based decisionmaking systems and information, this can be achieved (Koh, & Tan, 2011). All over the world, heart disease is the leading cause of mortality in men and women, and more than half of the deaths occur in men. One out of four people has heart disease and dies in the United States. In the United States, more than 610,000 Americans are affected by heart disease and lose their lives every year (Masethe, & Masethe, 2014). It is necessary to use computer technology to help physicians diagnose heart disease with greater speed and accuracy. Equipping medical science with smart tools in the diagnosis and treatment of diseases can reduce doctors' mistakes and loss of life and money. Because correctly predicting people's disease status is very important. In terms of information, medical environments are rich environments. Estimates show that a special hospital for acute diseases produces five terabytes of data per year. In front of this massive amount of information, we are facing the poverty of knowledge. There is a huge amount of data in medical systems; This is an important advantage to being able to increase the quality of services provided in the field of health and treatment to a very desirable level by extracting the knowledge hidden in the heart of this information using data mining techniques and applying this knowledge in carrying out processes. In this regard, the field of heart disease is of double importance due to the sensitivity of its health in the continuation of human life, and the improvement of diagnoses and treatments in this field can save many human lives (Ramakrishnan, Hanauer, & Keller, 2010). Data mining is the process of discovering knowledge such as associations, patterns, anomalies, significant changes and structures, from a large amount of information stored in a database or other information repositories (Dey, Singh, J., & Singh, N, 2016). In other words, data mining is a process that analyzes a set of data to find the relationships in them. Data mining is a way to summarize the data, which in addition to making it easier to understand the data, provides the possibility of using the knowledge in the data (Mukhopadhyay, Maulik, Bandyopadhyay, & Coello, 2013). Therefore, data mining is used to overcome this problem and obtain useful relationships between risk factors in diseases with regard to the prevalence and their contribution to human mortality (Masethe, & Masethe, 2014).

The innovation of our work can be summarized as follows: We implemented a comprehensive fourphase approach for diagnosing and evaluating heart patients. In the second phase, we employed three powerful met heuristic algorithms, namely Grey Wolf Optimizer (GWO), Particle Swarm Optimization (PSO) and Genetic Algorithm (GA), to effectively select relevant features. In the third phase, we utilized various machine learning algorithms to evaluate the performance of our proposed method.

The main goals of this paper are as follows:

• Examining existing studies in the diagnosis of heart disease

■ Introducing the proposed approach

 Modeling using machine learning and metaheuristics

■ The outcomes of the evaluation

The following is the paper's continuation: In the second section, we examined at related works. The proposed method was introduced in the third section. The results have been assessed in the fourth section and then, in the fifth section, we have a conclusion.

2- Literature Review

Machine learning is a developing subset of computing algorithms that aims to imitate human intelligence through environmental learning. In the brand-new era of "big data," they are regarded as the workhorse. Various fields, including pattern recognition, computer vision, spacecraft engineering, finance, entertainment, and computational biology, as well as biological medical applications, have effectively and used machine learning techniques (El Naga, & Murphy, 2015). As reported in (Moharamkhani, Yahyaei Feriz Hendi, Bandar, Izadkhasti, & Sirwan Raza, 2022) Moharamkhani et al. proposed a new approach for intrusion detection in cloud computing environments using a combination of the Firefly Algorithm and Random Forest. Hasanvand, Nooshyar, Moharamkhani, & Selvari (2023) proposed a different approach that focuses on identifying vehicles using machine learning and image processing techniques. Used in (Gavari Bami, Moharamkhani, Zadmehr, Najafpoor, & Shokouhifar, 2022) Gavari Bami et al. to accurately identify attacks.

Anbarasi, Anupriya, & Iyengar (2010)in advanced prediction of heart diseases by selecting a subset of genetic algorithm features using genetic algorithm, they predicted coronary artery disease. Using the genetic algorithm, they reduced the number of features from 13 to six and performed their processing with six features that had a greater impact on disease diagnosis. After that, they used three classifiers such as Naive Bayes, classification by clustering and decision tree to predict the diagnosis of patients with the same accuracy as before reducing the number of features. The observations showed that the decision tree data extraction technique is relatively faster and better than the other two data-mining methods after combining the feature subset selection with the time of building the model.

According to the Bhatla, & Jyoti, (2012), the proposed strategy can increase diagnosis accuracy and reduce errors in medical decisionmaking. They give an outline of related research on heart disease detection and explain how data mining and fuzzy logic approaches were applied to the problem. The report contains experimental results that demonstrate the efficacy of their proposed approach. Overall, the findings indicate that data mining and fuzzy logic can be valuable techniques for improving heart disease diagnosis.

Shouman, Turner, & Stocker, (2011), have presented a model to increase the accuracy of the decision tree classification method in the diagnosis of patients with coronary arteries, which by applying various types of decision trees on the data set and using the majority voting technique, they achieved an accuracy of 84.1 and showed that this model has better accuracy compared to Bagging methods and J48 decision tree.

Cardiovascular disorders are a primary cause of death worldwide, according to Subbalakshmi, G., Ramesh, K., & Rao, M. C. (2011), and early identification is crucial in reducing mortality rates. They propose a system that uses patient data and the Naive Bayes algorithm to estimate a patient's risk of acquiring heart disease. The study contains a full discussion of the Naive Bayes algorithm and its application to the prediction of heart disease. They also give experimental results that demonstrate the efficacy of their proposed method. Overall, the research indicates that the suggested decision support system has the potential to help medical practitioners make accurate diagnoses and improve patient outcomes.

Bashir, Qamar, Khan, & Javed, (2014), have proposed a clinical decision support system for heart disease diagnosis. First, five simple Bayesian classification methods, decision tree based on Gini index, decision tree based on information gain, learner based on memory and support vector machine are applied on the data set; the final result of diagnosis is obtained by applying the majority vote on the results of these five classes. The application of this system on several datasets of heart disease, in the best case, has an accuracy of 90.3, which has higher accuracy than the new ones.

Negahbani, Joulazadeh, Marateb, & Mansourian, (2015), have proposed a combined system of automatic CAD diagnosis using the fuzzy Means-K clustering method, which has an accuracy of 87% in correct diagnosis. This method uses a statistical selection method of features in order to reduce the number of features and select the most important ones. Considering that the types of features are of different numerical and nominal types, the generalized Minkowski distance criterion is used to calculate the distance between samples.

Koluksa, Haclar, Ku, Bakr-Güngör, Aral, & Güngör, (2019), tried to improve the method of

diagnosing coronary artery disease by classifying algorithms and new features of method selection. In this study, a set of different classification algorithms along with a new combined feature selection method was proposed for the diagnosis of coronary heart disease; KNN, Naive Bayes, Random Forest, Bagging and MLP algorithms were used based on the selection of different features that included one to several features in different datasets. Finally, the SVM algorithm with 81.84% accuracy in the UCI data set and the Z-Alizadehsani data set and the SVM algorithm with 87.12% accuracy were able to diagnose coronary heart disease patients and performed better than other algorithms.

Kalaiarasi, Maheswari, Selvi, Yogitha, & Devadas, (2022), found it a mistake for clinicians to prove the diagnosis of coronary artery disease in an individual, as it requires long periods of experience and extraordinary clinical tests to be conducted. At present, data mining classification calculations, such as decision tree, naive Bayes, and random forest are used to develop a framework of expectations for predicting the chance of coronary artery disease. The main objective of this paper is to identify the best suitable feature calculation to obtain the maximum accuracy when the normal and odd person order are completed. The results have been done using the coronary artery disease benchmark data from the artificial intelligence archive; it showed that random forest calculation performed the best with 81% accuracy compared to different calculation methods for predicting coronary artery disease.

Aliyar Vellameeran, & Brindha, (2022), proposed a novel type of deep belief network (DBN) for diagnosing heart disease using IoT wearable medical devices and optimal feature selection. The study's primary goal was to train the DBN model by examining and choosing the most crucial features from a sizable dataset. The proposed method showed promising results in accurately diagnosing cardiac illness when tested using real data collected from wearable medical devices connected to the Internet of Things.

Shorewala (2021), investigated the use of ensemble techniques in the early identification of coronary heart disease. The purpose of the study was to evaluate various models and algorithms, including decision tree, random forest, support vector machine, k-nearest neighbor, and neural network models, in order to determine the most effective technique for early disease diagnosis. The results showed that multi-algorithm ensemble techniques led to the highest degree of disease prediction accuracy.

Latha, & Jeeva, (2019), investigated the efficacy of several machine learning algorithms such as support vector machines, decision trees, and random forests. They compared the various approaches with an ensemble approach that combined these models. The outcomes demonstrated that in terms of prediction accuracy, sensitivity, and specificity, the ensemble technique performed better than the individual algorithms. The study also stressed how crucial feature choice is to improving model accuracy.

Waigi, Choudhary, Fulzele, & Mishra, (2020), used a variety of algorithms to assess a dataset comprising multiple medical indicators of patients with and without heart disease, including decision trees, random forests, and support vector machines. The results showed that the suggested machine learning algorithms could predict the risk of heart disease successfully and precisely.

Next, in table 1, we have discussed the advantages and disadvantages of the existing works.

Table (1): An overview of existing works.					
Ref Main objectives Advantages Disadvantages					
Anbarasi etal (2010)	-Improve heart disease prediction accuracy by selecting feature subsets based on genetic algorithms.	-Increased prediction accuracy -Decreased computing complexity -Identification of a smaller collection of important features	-The requirement for suitable parameter tuning and the potential of overfitting the data		

Table (1): An overview of existing works.					
Shouman etal (2011)	-The study's major objective is to provide a decision tree-based model that can help clinicians effectively diagnose cardiac problems in patients.	-Because decision trees are simple to comprehend and use, they are a useful tool for clinicians of all levels of experience. -Based on easily available clinical data, the model can deliver reliable diagnoses.	 -Decision trees have the potential to oversimplify complex medical issues and overlook essential details. -The quality of the input data, which may contain errors or inconsistencies, may have an impact on the model's accuracy. 		
Subbalakshmi etal (2011)	-The study's major goal was to create a heart disease prediction system based on the naïve Bayes algorithm and examine its usefulness in providing clinicians with decision support.	-The naive Bayes algorithm used in the heart disease prediction system is fast and accurate. -The system uses a simple and easy-to-understand user interface.	 The study was conducted on a relatively small dataset and may not be representative of larger populations. The system only considers a limited set of factors for heart disease prediction and may not be comprehensive enough for a full diagnosis. 		
Bhatla, & Jyoti, (2012)	 -To develop a novel approach for heart disease diagnosis using data mining and fuzzy logic. -To evaluate the effectiveness of the proposed approach in diagnosing heart disease. -To compare the proposed approach to other diagnostic methods and demonstrate its superiority. 	 The method analyzes a variety of parameters such as age, gender, blood pressure, and cholesterol levels, making it a thorough diagnostic tool. When compared to standard diagnostic procedures, the use of fuzzy logic allows for more flexible decision-making. 	 The study does not provide information on the number or demographics of participants, which could impact the generalizability of the results. The study does not compare the proposed approach to other diagnostic methods, so its efficacy relative to other approaches is unknown. 		
Bashir etal (2014)	-The study's primary objective is to develop the MV5 clinical decision support system for predicting heart disease using a majority vote based classifier ensemble.	 -MV5 achieved high accuracy levels in predicting heart disease, which can help healthcare professionals make more informed decisions about patient care. -The model is based on an ensemble approach, which combines the strengths of multiple classifiers to improve overall prediction performance. -The authors used a large and diverse dataset to train and test their model, which enhances its generalizability to different populations. 	 -The study did not compare the performance of MV5 with other existing heart disease prediction models, which limits the ability to determine whether it outperforms them. -The paper did not discuss the interpretability of the model or how the predictions were obtained, which may be important for gaining trust from healthcare professionals and patients. -The study was conducted using data from a single hospital, which may limit the generalizability of the findings to other healthcare settings or populations. 		

Table (1): An overview of existing works.					
Negahbani etal (2015)	-The study's major objective was to develop a new approach for identifying coronary artery disease using supervised fuzzy c-means and differential search algorithm-based generalized Minkowski metrics. The study intended to increase diagnosis accuracy while maintaining computational performance.	 The proposed method is able to accurately diagnose coronary artery disease. The use of supervised fuzzy c-means with differential search algorithm-based generalized Minkowski metrics improves the accuracy of the diagnosis. The proposed method is computationally efficient. 	 The study has a relatively small sample size, which may affect the generalizability of the results. The proposed method requires expert knowledge in both fuzzy clustering and differential search algorithms, which may limit its applicability to non- expert users. 		
Koluksa etal (2019)	The study's primary objective was to provide a novel way for identifying coronary heart disease using classification algorithms and a new feature selection methodology. The study's goal was to enhance diagnosis accuracy with a smaller collection of features and to determine the most essential criteria in the diagnosis of coronary heart disease.	 -Using classification algorithms and a novel feature selection methodology, the suggested method may accurately diagnose coronary heart disease. -The application of feature selection aids in the identification of the most relevant characteristics in the diagnosis of coronary heart disease, which can aid in treatment decisions. -The proposed method is computationally efficient and has clinical relevance. 	 -The study had a limited sample size, which may limit the generalizability of the findings. -Because the suggested feature selection methodology necessitates expert knowledge in both data mining and medicine, it may be inapplicable to non-expert users. 		
Latha, & Jeeva (2019)	-This study developed an improved model for identifying high-risk individuals for personalized health care.	-Ensemble classification techniques enhance heart disease risk prediction accuracy by combining multiple classifiers, resulting in more reliable and precise predictions. - Robustness	-Complex Implementation -Computational Resources		
Waigi, Choudhary, Fulzele, & Mishra, (2020)	 -Applying advanced machine learning algorithms to predict the risk of heart disease. -Evaluating the effectiveness and accuracy of the proposed model in comparison to traditional risk assessment methods. 	-Advanced machine learning algorithms predict heart disease risk, outperforming traditional risk assessment methods by considering complex interactions.	-Data limitations impact predictive model effectiveness; sample size and representativeness may affect training and testing.		
Shorewala (2021)	-This study proposesd ensemble technique for early detection of coronary heart disease, improving accuracy and efficiency using multiple classifiers and machine learning.	 -Ensemble technique improves coronary heart disease diagnosis accuracy by combining classifier predictions. -Paper focused on early detection of coronary heart disease for improved patient outcomes and reduced complications. 	-Ensemble techniques require significant computational resources, impacting resource- constrained environments and large datasets. -Complexity		

Table (1): An overview of existing works.				
Aliyar Vellameeran etal (2022)	-The paper proposed a new Deep Belief Network for optimal feature selection in IoT wearable medical devices, improving heart disease diagnosis accuracy and efficiency.	-Improved Accuracy -Incorporation of IoT Wearable Devices	-Limited Generalizability: -Model Complexity and Computational Cost	
Kalaiarasi etal (2022)	 -To present a data mining approach for detecting cardiac disease. -Using real-world data, assess the performance of the proposed approach. -To examine the efficacy of various data mining algorithms for detecting heart disease. 	 The proposed method employs a variety of data mining techniques, including Decision Tree (DT), Random Forest (RF), Nave Bayes (NB), and K-Nearest Neighbor (KNN), to improve the accuracy of the results. The study uses real-world evidence to evaluate the performance of the proposed approach, making the conclusions more credible and applicable in practice. 	-The research does not contrast the proposed approach to other known methods for detecting cardiac disease, making it difficult to determine its superiority over other ways. -The study does not give information on the demographic features of the participants, which may impair the findings' generalizability.	

3- Method

Figure 1 depicts the proposed approach, which is separated into four steps.



Figure (1): Flowchart of the proposed approach $% \left(f_{1}, f_{2}, f_{3}, f_{3$

Phase 1 (Data preprocessing)

Data preprocessing is the first step, and we try to identify noisy, missing and inconsistent data. At this stage, we will test and train data with 80 and 20 percent. The dataset was 1025 people with 13 characteristics including age, gender, type of pain in the chest, blood pressure, fasting blood glucose, etc (Oliullah, Barros, & Whaiduzzaman, 2023). Before choosing the features, we show all the significant characteristics of heart patients in figure number 2.





فصلنامه علمی مؤسسه آموزش عالی فردوس ۱ بامشارکت انجـمن علمـی مـدیـریت دانش ایـــران ۲۰۰ ۸۷ –۱

















فصلنامەعلمىمۇسسەآموزش عالى فردوس — ۸۸ — بامشاركتانجـمن علمىمديـريت دانش ايــــران ۲

l MonyAdditaalbiture h. Restecg



Figure (2): 13 Features before selection

phase 2 (Feature selection)

Using meta-heuristic algorithms such as PSO with three thresholds (0.3, 0.4, and 0.5) and GWO and GA algorithms, we have selected features on the desired data (Kennedy, & Eberhart, 1995., Mirjalili, Mirjalili, & Lewis, 2014., Mitchell, 1998). Phase 3 (Modeling with machine-learning

algorithms)

Our proposed modeling method is using machinelearning algorithms such as Logistic Regression, Decision Tree, Random Forest, Gaussian Naïve Bayes, SVC, and KNN.

Evaluation criteria

$$\operatorname{Re}call = \frac{TP}{FN + TP} \tag{1}$$

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

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$
(3)

$$F_Value = \frac{2*Precision*recall}{Precision+recall}$$
(4)

$$DR = \frac{TP}{TP + FN} \tag{5}$$

■ The most important aspect to take into account when assessing a classification algorithm's effectiveness is accuracy.

• The most important factor in determining how good a classification algorithm is precision, which displays what proportion of the entire set of test records is correctly categorized.

■ As you can recall, it shows the proportion of

valid data that is accurately tagged.

■ The f-value measure is the harmonic average of recall or precision.

■ Additionally, the DR criterion considers the proportion of anomalous samples that the IDS discovered out of all of the test set's samples that were abnormal.

Phase 4 (Performance evaluation)

In this stage, we assess our results and the output graphs' accuracy, precision, recall, and f-value.

4- Research Findings

for which we have separated the output from each of the three meta-heuristic methods. The 13 features that were utilized to create the dataset are first presented in table 2.

In this section, we assess our suggested approach,

Table 2. Display of the dataset for heart disease utilized in this research.				
Featurename	Featurename Description		Values	
Age	Age of the patients	Numeric	Years	
Ca	Number of major vessels	Numeric	04-	
Chol	Serum cholesterol	Numeric	mg/dl	
Ср	Chest pain type	Numeric	Male=1; Female=0	
Exang	Exercise induced angina	Numeric	Yes=1; N0=0	
Fbs	Fasting blood sugar	Numeric	mg/dl	
Oldpeak	ST depression induced by exercise relative to rest	Numeric	06.2-	
Restecg	Resting electrocardiographic	Numeric	0,1,2	
Sex	Gender of patients	Numeric	M,F	
Slope	the slope of the peak exercise ST segment	Numeric	0,1,2	
Thal	Normal; Fixed defect; Reversible defect	Numeric	0,1,2,3	
Thalach	Maximum heart rate achieved	Numeric	71202-	
Trestbps	Resting blood pressure	Numeric	94200-	

We have summarized a list of the simulation parameters for meta-heuristic algorithms in Table 3.

Table (3): Simulation parameters			
Parameter	value		
Number of nodes	100		
Number of iteration	100		

The outcomes of the GA

Each row of the table represents a separate algorithm, while the columns relate to several algorithm evaluation measures.

The name of the algorithm being evaluated is listed in the first column, which includes Logistic Regression, Decision Tree, Random Forest, Gaussian Nave Bayes, SVC (Support Vector Classifier), and KNN (K-Nearest Neighbors).

The accuracy score for each algorithm is shown in the second column. Accuracy is calculated as the ratio of correct predictions to total predictions and measures how often the algorithm successfully predicts the class label.

The F1-value, which is a measure of the balance between precision and recall, is given in the third column. F1-score combines precision and recall into a single statistic and allows you to compare the performance of classifiers with different trade-offs between accuracy and recall.

The precision and recall scores are shown in the fourth and fifth columns, respectively. Precision evaluates how many predicted positive occurrences were actually positive, whereas recall indicates how many actual positive cases were properly predicted by the model.

The output results of feature selection using a GA algorithm are displayed in Table 4 in this section. The Decision Tree and Random Forest algorithms had the greatest accuracy scores (98%), according to the results in the table, suggesting that they accurately predicted outcomes for a significant portion of the dataset's data points. It is crucial to keep in mind, though, that these algorithms can be prone to overfitting the data, which could prevent them from generalizing successfully to fresh, untested data.

accuracy score of 95.1%, the KNN algorithm also did well. When the data is organized into distinct clusters or groups, this technique performs quite well.

Although having varying precision and recall scores, the Logistic Regression, Gaussian Naive Bayes, and SVC algorithms all had relatively good accuracy scores (ranging from 79.5% to 87.3%) and F1-values. These algorithms could be beneficial for various datasets or applications.

In general, while choosing a machine learning algorithm for a specific task, it is crucial to take into account variables like interpretability, computational efficiency, and generalizability.

Table 4. Comparing the outputs of ML algorithms with GA feature selection.				
Algorithms	Accuracy	F1-value	Precision	Recall
Logistic Regression	82.4	83.5	82	85
Decision Tree	98	98.2	96.4	100
Random Forest	98	98.2	96.4	100
Gaussian Naïve Bayes	79.5	80.2	81	79.4
SVC	87.3	88.6	83.5	94.4
KNN	95.1	95.1	100	90.7

With a precision score of 100 percent and an

The outcomes of the GWO

The output results of feature selection using a GWO algorithm are displayed in Table 5 in this section. With an accuracy of 92.2%, F1-value of 92.4%, precision of 94.2%, and recall of 90.7%, the Decision Tree method did the best. With an accuracy of 91.7%, F1-value of 91.9%, precision of 93.3%, and recall of 90.7%, Random Forest performed well as well.

A decent performance was achieved by Logistic

Regression, SVC, and KNN, with accuracy ranging

from 82% to 83.8% and F1-values from 81.8% to 90%. Of all methods, Gaussian Naive Bayes performed the worst, with an accuracy of 80.5%, an F1-value of 81.8%, a precision of 79.6%, and a recall of 84.1%.

In conclusion, the study's Decision Tree and Random Forest algorithms outperformed the others, while Logistic Regression, SVC, and KNN did only moderately well. The performance of Gaussian Naive Bayes was the worst of all the algorithms.

Table (5): Comparing the outputs of ML algorithms with GWO feature selection.				
Algorithms	Accuracy	F1-value	Precision	Recall
Logistic Regression	82	81.8	86.5	77.6
Decision Tree	92.2	92.4	94.2	90.7
Random Forest	91.7	91.9	93.3	90.7
Gaussian Naïve Bayes	80.5	81.8	79.6	84.1
SVC	83.4	83.8	85.4	82.2
KNN	89.8	90	91.3	88.8

فصلنامه علمى مؤسسه آموزش عالى فردوس

I بامشارکت انجـ من علمـی مـدیـریت دانش ایـــــران 📕 🖣 🗕

The outcomes of the PSO 0.3

Table 6 in this section shows the output results of feature selection using a PSO 0.3 algorithm. The performance of the Decision Tree and Random Forest algorithms is flawless, with 100% F1-value, precision, and recall ratings. This can mean that the model has been overfit to the training data, so it's critical to validate the model on a different test set to make sure it generalizes properly.

SVC outperforms the other three algorithms in terms of accuracy, F1-value, precision, and recall, with a score of 86.8%, 88.1%, 83%, and 93.5% respectively. KNN also performs well, with 99% accuracy, 99.1% F1 value, 100% precision, and 98.1% recall.

When compared to the other algorithms, Logistic Regression and Gaussian Naive Bayes perform worse, with accuracy ranging from 81.5% to 82.7% and F1-values from 80% to 82.9%. Nonetheless, Logistic Regression has a lower recall score (85%) than Gaussian Naive Bayes (86%).

In conclusion, SVC and KNN perform better than Decision Tree and Random Forest, which appear to have overfitted to the training data. Both Gaussian Naive Bayes and Logistic Regression perform poorly, although Gaussian Naive Bayes outperforms Logistic Regression in terms of recall.

Table 6. Comparing the outputs of ML algorithms with PSO 0.3 feature selection.				
Algorithms	Accuracy	F1-value	Precision	Recall
Logistic Regression	81.5	82.7	80.5	85
Decision Tree	100	100	100	100
Random Forest	100	100	100	100
Gaussian Naïve Bayes	81.5	82.9	80	86
SVC	86.8	88.1	83	93.5
KNN	99	99.1	100	98.1

The outcomes of the PSO 0.4

Table 7 in this section shows the output results of feature selection using a PSO 0.4 algorithm. Both the Random Forest and Decision Tree algorithms received perfect scores for each of the four measures, demonstrating that they were both able to accurately categorize every instance in the dataset without making any mistakes. Perfect scores, however, may be a sign of overfitting, which could lead to subpar performance on fresh, untested data.

The Logistic Regression algorithm classified a sizable percentage of the instances in the dataset correctly, but there was still space for improvement. Its accuracy score was 80.5%, and its F1-value was 82.1%. It may have more false negatives than false positives since its precision score (78.6%) is slightly lower than its recall score (86%).

The accuracy score and F1-value for the Gaussian Naive Bayes algorithm were both 79%

and 81.1%, respectively. Also, it had a recall score of 86% and a precision score of 76.7%, indicating a potential problem with false positives.

The SVC algorithm had a good F1-value of 85.6% and the greatest recall score (91.6%). Although it has space for growth, especially in terms of reducing false positives, its accuracy score (83.9%) and precision score (80.3%) indicate that it has room for progress.

A high accuracy score of 97.1% and an F1value of 97.1% were attained by the KNN algorithm. Although it had a recall score of 94.4%and a precision score of 100%, this means that it may have more false negatives than false positives.

Overall, even though the Decision Tree and Random Forest algorithms got faultless marks, overfitting should be taken into account. Although the other algorithms performed largely satisfactorily, there is still potential for improvement with regard to lowering false positives or false negatives.

Table (7): Comparing the outputs of ML algorithms with PSO 0.4 feature selection.				
Algorithms	Accuracy	F1-value	Precision	Recall
Logistic Regression	80.5	82.1	78.6	86
Decision Tree	100	100	100	100
Random Forest	100	100	100	100
Gaussian Naïve Bayes	79	81.1	76.7	86
SVC	83.9	85.6	80.3	91.6
KNN	97.1	97.1	100	94.4

The outcomes of the PSO 0.5

Table 8 in this section shows the output results of feature selection using a PSO 0.5 algorithm. With perfect scores across all four measures, the Decision Tree and Random Forest algorithms were able to accurately categorize every instance in the dataset without making any mistakes. Perfect scores, however, can point to overfitting, as was already indicated.

With an accuracy score of 78.5% and an F1value of 82.2%, the Logistic Regression algorithm was able to accurately categorize a sizable part of the dataset's instances, but there was still space for improvement. It may have more false negatives than false positives because its precision score (77.4%) is slightly lower than its recall score (83.2%).

The accuracy score and F1-value for the Gaussian Naive Bayes algorithm were 79.5% and 79.8%, respectively. It may have more false positives than false negatives because its precision

score (82.2%) is greater than recall score (77.6%).

The accuracy score for the SVC method was 81%, the F1-value was 83%, the precision score was 77.9%, and the recall score was 88.8%. Despite having a high recall score, it suggests that there can be more false negatives than false positives. It may also have more false positives than false negatives, according to its precision score.

The accuracy score and F1-value for the KNN algorithm were 93.2% and 93%, respectively. Its recall score (86.9%) is lower than those of the other algorithms assessed in this table, even if its precision score was 100%, suggesting that it may have more false negatives than false positives.

Ultimately, the Random Forest and Decision Tree algorithms outperformed all others, scoring 100 percent on every criterion. Although the other algorithms performed largely satisfactorily, there is still potential for improvement with regard to lowering false positives or false negatives.

Table (8): Comparing the outputs of ML algorithms with PSO 0.5 feature selection.				
Algorithms	Accuracy	F1-value	Precision	Recall
Logistic Regression	78.5	82.2	77.4	83.2
Decision Tree	100	100	100	100
Random Forest	100	100	100	100
Gaussian Naïve Bayes	79.5	79.8	82.2	77.6
SVC	81	83	77.9	88.8
KNN	93.2	93	100	86.9

5- Conclusion

In recognition of the fact that heart disease is one of the most prevalent illnesses and a leading cause of mortality, researchers in this area have recently proposed a number of strategies and algorithms. In this paper, we sought to use machine learning algorithms to diagnose patients with great precision and accuracy. The detection rate of our suggested approach is then demonstrated in Figures 3-7.



Figure (5): Detection rate with PSO 0.3



Figure (7): Detection rate with PSO 0.5.

فصلنامه علمى مؤسسه آموزش عالى فردوس

🗕 ۹۴ 🛏 بامشارکت انجـ من علمــی مـدیــریت دانش ایـ

Diagnosing the disease is the most important and the first step in the treatment process. One of the most widely used methods in the diagnosis of a disease is to create rules based on which a person is classified into one of two classes, sick or healthy. In fact, by having a set of data that includes the signs and characteristics of the patient and healthy individuals, the diagnosis procedure is accelerated. In addition, the science of data mining in medicine is very much appreciated by those interested and attracted by scholars of this field due to its unique function. In this research, by using this science and using meta-heuristic algorithms in the feature selection section, and using machine-learning algorithms in the modeling section, we were able to present a method that can detect heart diseases with high accuracy.

References

1-Koh, H. C., & Tan, G. (2011). Data mining applications in healthcare. *Journal of healthcare information management*, *19*(2), 65.

2-Masethe, H. D., & Masethe, M. A. (2014, October). Prediction of heart disease using classification algorithms. In *Proceedings of the world Congress on Engineering and computer Science* (Vol. 2, No. 1, pp. 25-29).

3-Ramakrishnan, N., Hanauer, D., & Keller, B. (2010). Mining electronic health records. *Computer*, 43(10), 77-81.

4-Dey, A., Singh, J., & Singh, N. (2016). Analysis of supervised machine learning algorithms for heart disease prediction with reduced number of attributes using principal component analysis. *International Journal of Computer Applications*, *140*(2), 27-31.

5-Mukhopadhyay, A., Maulik, U., Bandyopadhyay, S., & Coello, C. A. C. (2013). A survey of multiobjective evolutionary algorithms for data mining: Part I. *IEEE Transactions on Evolutionary Computation*, *18*(1), 4-19.

6-El Naqa, I., & Murphy, M. J. (2015). *What is machine learning?* (pp. 3-11). Springer International Publishing.

7-Moharamkhani, E., Yahyaei Feriz Hendi, M., Bandar, E., Izadkhasti, A., & Sirwan Raza, R. (2022). Intrusion detection system based firefly algorithmrandom forest for cloud computing. *Concurrency and Computation: Practice and Experience*, *34*(24), e7220.

8-Hasanvand, M., Nooshyar, M., Moharamkhani, E., & Selyari, A. (2023, April). Machine Learning Methodology for Identifying Vehicles Using Image Processing. In *Artificial Intelligence and Applications* (Vol. 1, No. 3, pp. 170-178).

9-Gavari Bami, H., Moharamkhani, E., Zadmehr, B., Najafpoor, V., & Shokouhifar, M. (2022). Detection of zeroday attacks in computer networks using combined classification. *Concurrency and Computation: Practice and Experience, 34*(27), e7312.

10-Anbarasi, M., Anupriya, E., & Iyengar, N. C. S. N. (2010). Enhanced prediction of heart disease with feature subset selection using genetic algorithm. *International Journal of Engineering Science and*

فصلنامه علمى مؤسسه آموزش عـالى فردوس

Technology, 2(10), 5370-5376.

11-Bhatla, N., & Jyoti, K. (2012). A Novel Approach for heart disease diagnosis using Data Mining and Fuzzy logic. *International Journal of Computer Applications*, *54*(17).

12-Shouman, M., Turner, T., & Stocker, R. (2011). Using Decision Tree for Diagnosing Heart Disease Patients. *AusDM*, *11*, 23-30.

13-Subbalakshmi, G., Ramesh, K., & Rao, M. C. (2011). Decision support in heart disease prediction system using naive bayes. *Indian Journal of Computer Science and Engineering (IJCSE)*, 2(2), 170-176.

14-Bashir, S., Qamar, U., Khan, F. H., & Javed, M. Y. (2014). MV5: a clinical decision support framework for heart disease prediction using majority vote based classifier ensemble. *Arabian Journal for Science and Engineering*, *39*, 7771-7783.

15-Negahbani, M., Joulazadeh, S., Marateb, H. R., & Mansourian, M. (2015). Coronary artery disease diagnosis using supervised fuzzy c-means with differential search algorithm-based generalized Minkowski metrics. *Peertechz Journal of Biomedical Engineering*, 1(1), 006-014.

16-Koluksa, B., Haclar, H., Ku, M., Bakr-Güngör, B., Aral, A., & Güngör, V. Ç. (2019). Diagnosis of coronary heart disease via classification algorithms and a new feature selection methodology. *International Journal of Data Mining Science*, 1(1), 8-15.

17-Kalaiarasi, G., Maheswari, M., Selvi, M., Yogitha, R., & Devadas, P. (2022). Detection of Heart Disease Using Data Mining. In *Biologically Inspired Techniques in Many Criteria Decision Making: Proceedings of BITMDM 2021* (pp. 627-637). Singapore: Springer Nature Singapore.

18-Aliyar Vellameeran, F., & Brindha, T. (2022). A new variant of deep belief network assisted with optimal feature selection for heart disease diagnosis using IoT wearable medical devices. *Computer Methods in Biomechanics and Biomedical Engineering*, 25(4), 387-411.

19-Shorewala, V. (2021). Early detection of coronary heart disease using ensemble techniques. *Informatics inMedicineUnlocked, 26*, 100655.

20-Latha, C. B. C., & Jeeva, S. C. (2019). Improving the accuracy of prediction of heart disease risk based on ensemble classification techniques. *Informatics in Medicine Unlocked*, *16*, 100203.

21-Waigi, D., Choudhary, D. S., Fulzele, D. P., & Mishra, D. (2020). Predicting the risk of heart disease using advanced machine learning approach. *Eur. J. Mol. Clin. Med*, 7(7), 1638-1645.

22-Oliullah, K., Barros, A., & Whaiduzzaman, M. (2023, May). Analyzing the Effectiveness of Several Machine Learning Methods for Heart Attack Prediction. In *Proceedings of the Fourth International Conference on Trends in Computational and Cognitive Engineering: TCCE 2022* (pp. 225-236). Singapore: Springer Nature Singapore.

23-Kennedy, J., & Eberhart, R. (1995, November). Particle swarm optimization. In *Proceedings of ICNN'95-international conference on neural networks* (Vol. 4, pp. 1942-1948). IEEE.

24-Mirjalili, S., Mirjalili, S. M., & Lewis, A. (2014). Grey wolf optimizer. *Advances in engineering* software, 69, 46-61.

25-Mitchell, M. (1998). An introduction to genetic algorithms. MIT press.

©Authors, Published by Journal of Intelligent Knowledge Exploration and Processing. This is an openaccess paper distributed under the CC BY (license http://creativecommons.org/licenses/by/4.0/).

