

Research Article

Feature Subset Selection with Optimal Adaptive Neuro-Fuzzy Systems for Bioinformatics Gene Expression Classification

Anwer Mustafa Hilal ¹, **Areej A. Malibari**,² **Marwa Obayya**,³ **Jaber S. Alzahrani**,⁴ **Mohammad Alamgeer**,⁵ **Abdullah Mohamed**,⁶ **Abdelwahed Motwakel**,¹ **Ishfaq Yaseen**,¹ **Manar Ahmed Hamza** ¹ and **Abu Sarwar Zamani**¹

¹Department of Computer and Self Development, Preparatory Year Deanship, Prince Sattam Bin Abdulaziz University, AlKharj, Saudi Arabia

²Department of Industrial and Systems Engineering, College of Engineering, Princess Nourah Bint Abdulrahman University, P.O. Box 84428, Riyadh 11671, Saudi Arabia

³Department of Biomedical Engineering, College of Engineering, Princess Nourah Bint Abdulrahman University, P.O.Box 84428, Riyadh 11671, Saudi Arabia

⁴Department of Industrial Engineering, College of Engineering Alqunfudah, Umm Al-Qura University, Mecca, Saudi Arabia

⁵Department of Information Systems, College of Science & Art Mahayil, King Khalid University, Abha, Saudi Arabia

⁶Research Centre, Future University, Egypt, New Cairo 11845, Egypt

Correspondence should be addressed to Anwer Mustafa Hilal; a.hilal@psau.edu.sa

Received 9 March 2022; Revised 20 April 2022; Accepted 27 April 2022; Published 14 May 2022

Academic Editor: Laxmi Lydia

Copyright © 2022 Anwer Mustafa Hilal et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Recently, bioinformatics and computational biology-enabled applications such as gene expression analysis, cellular restoration, medical image processing, protein structure examination, and medical data classification utilize fuzzy systems in offering effective solutions and decisions. The latest developments of fuzzy systems with artificial intelligence techniques enable to design the effective microarray gene expression classification models. In this aspect, this study introduces a novel feature subset selection with optimal adaptive neuro-fuzzy inference system (FSS-OANFIS) for gene expression classification. The major aim of the FSS-OANFIS model is to detect and classify the gene expression data. To accomplish this, the FSS-OANFIS model designs an improved grey wolf optimizer-based feature selection (IGWO-FS) model to derive an optimal subset of features. Besides, the OANFIS model is employed for gene classification and the parameter tuning of the ANFIS model is adjusted by the use of coyote optimization algorithm (COA). The application of IGWO-FS and COA techniques helps in accomplishing enhanced microarray gene expression classification outcomes. The experimental validation of the FSS-OANFIS model has been performed using Leukemia, Prostate, DLBCL Stanford, and Colon Cancer datasets. The proposed FSS-OANFIS model has resulted in a maximum classification accuracy of 89.47%.

1. Introduction

Microarray is an advanced technology that helps to recognize the pattern of gene expression of various genes at a time at the genomic level. It supports the researcher to investigate and analyze millions of genes in a single experiment [1]. It identifies many present diseases connected to each individual gene such as anaemia and

cancer. Analysis of Gene Expression provides a method to recognize the gene that is differentially expressed [2], which is accountable to develop some diseases. Also, it shows the difference between normal and abnormal genes through a mathematical model [3, 4]. Many openly accessible datasets such as Array Express and Gene Expression Omnibus (GEO) make the task easier to identify gene patterns of rare diseases. Classification of gene

expression data splits cancer samples from healthy samples that are utilized in response to treatment prediction. Due to the smaller amount of samples with a larger amount of features in the gene expression information, the standard ML method disappoints to implement better for cancer classification [5].

Recently, there has been tremendous growth in the medical field around the world. There are several computational approaches utilized in the bioinformatics field in the last few decades, for example, data mining and pattern recognition, to deal with higher-dimensional problems but still unsuccessful [6]. Thus, recently, machine learning (ML), a branch of artificial intelligence, has received considerable attention from researchers in gene expression and genomics [7]. Also, ML is a branch of data science; the main goal is to allow a model for training and learning to make decisions by itself in the future. Machine learning is widely classified into semisupervised, semi-supervised, supervised, and unsupervised learning [8]. For microarray data classification, the ML-based feature selection (FS) techniques such as gene selection techniques assist in selecting the essential gene [9]. Feature selection assists to preserve useful attributes. It is mainly utilized for the higher-dimensional data; simply, FS is a dimensionality reduction method. Feature selection significantly assists in the field that has relatively scarce and samples too many features, e.g., DNA Microarray and RNA sequencing [10]. This approach assists in better understanding of the feature space, preventing the scare of model overfitting, maximizing the model training time, handling the dimension, and maximizing the prediction accuracy. The results of FS are the optimum amount of features that are related to the provided class label that contributed to the prediction process.

This study introduces a novel feature subset selection with optimal adaptive neuro-fuzzy inference system (FSS-OANFIS) for gene expression classification. The FSS-OANFIS model designs an improved grey wolf optimizer-based feature selection (IGWO-FS) model to derive an optimal subset of features. Besides, the OANFIS model is employed for gene classification and the parameter tuning of the ANFIS model is adjusted by the use of coyote optimization algorithm (COA). The application of IGWO-FS and COA techniques helps in accomplishing enhanced microarray gene expression classification outcomes. For examining the enhanced outcomes of the FSS-OANFIS model, a comprehensive simulation analysis was performed on distinct datasets.

2. Related Works

In reference [11], a two-phase approach named as ML-integrated ensemble of feature selection (FS) technique is used, and then a survival study was presented. In a primary stage, it can be chosen the optimum amongst 7 ML approaches dependent upon classifier accuracy, utilizing the whole group of features (under this case miRNAs). In the secondary stage, dependent upon classifier accuracy values, the top feature in all the FS approaches is assumed for making an ensemble to offer

more categorization of miRNAs. Ayyad et al. [12] presented a novel classifier approach to gene expression data. Both executions are assumed that improve the performance of KNN. An important idea is for utilizing robust neighbors in trained data with utilizing a novel weighting approach. The authors in reference [13] presented a recently developed classification named Forest DNN (fDNN) for integrating the DNN structure with a supervised forest feature detector. Utilizing this built-in feature detector, this technique is capable of learning sparse feature representation and feeding the representation to NN for mitigating the overfitting problem. Dwivedi [14] developed a structure of approaches dependent upon supervised ML with utilizing the ANN approach for gene classification.

Shukla [15] established a novel gene selection (GS) approach by integrating minimum redundancy maximum relevance (mRMR) and teaching learning-based optimization (TLBO) for accurate cancer prediction. Primarily, during the presented method, mRMR was executed for determining one of the discriminative genes in the original feature set. In SVM, mRMR was utilized as a fitness function (FF) under the presented technique for selecting relevant features that are used for estimating the prediction accuracy and classifying cancer correctly. In reference [16], a novel social network analysis-based GS method was presented. The presented approach contains 2 important objectives: relevance maximization and redundancy minimization of chosen genes. During this approach, on all iterations, a maximal community was chosen repetitively. Next amongst the present genes under this community, the suitable genes were chosen by utilizing the node centrality-based condition.

In reference [17], an ensemble DL approach was presented for reducing the dimensional features. Primarily, the reduction of dimensional with utilize of auto-encoder (AE) by utilizing several hidden layers have occurred and under the next step, a folded AE is also utilized for reducing the dimensional of identical original data. Eventually, both are combined and top feature is chosen on the fundamental of T-score value. Forestiero et al. [18] presented a multiagent technique to create a distributing approach to DNA microarray management. The group of agents, whereas all one signifying a Microarray (or chip), implement from the parallel a sequence of easy functions exploiting local data and organized virtual infrastructure was created at a global level. The word embedded method, capable of capturing the semantic context and signifying microarray with vector, was utilized for mapping the chip, thus permitting advanced agent functions.

3. The Proposed Model

In this study, a new FSS-OANFIS model has been developed for microarray gene expression data classification. The presented FSS-OANFIS model encompasses a series of processes, namely, data preprocessing, IGWO-FS-based election of features, ANFIS classification, and

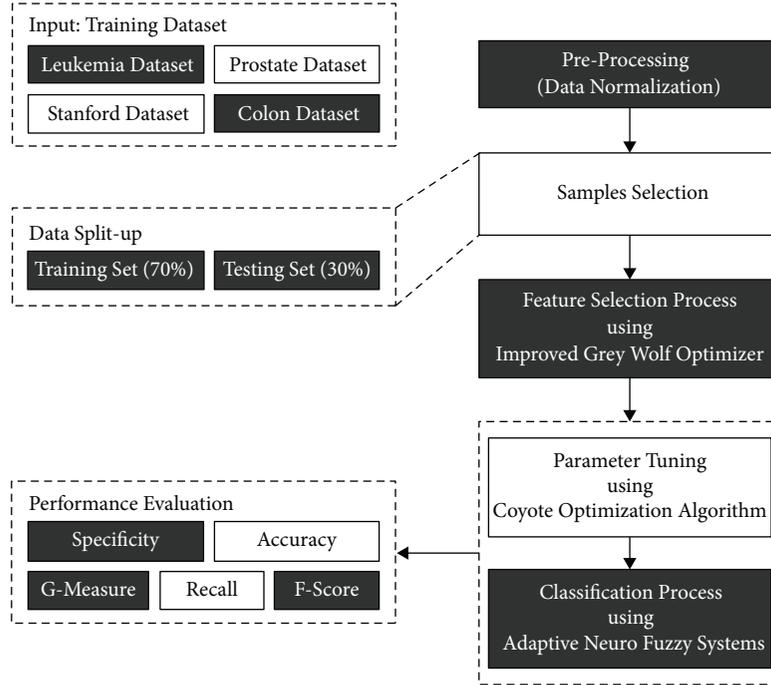


FIGURE 1: Overall process of FSS-OANFIS technique.

COA-based parameter optimization. The application of IGWO-FS and COA techniques helps in accomplishing enhanced microarray gene expression classification outcomes. Figure 1 shows the overall process of FSS-OANFIS technique.

3.1. Preprocessing. The z-score is a normalized and standardized system, which describes the count of standard deviation (SD), a raw data point, which is below or above the population mean [18]. It ideally lies in the range of -3 and $+3$. It standardizes the dataset to the aforementioned scale to change data with distinct scales to default scale. Thus, reflecting that several SD a point is below/above the mean as follows, but x refers to the value of particular instance, μ signifies the mean, and σ depicts the SD:

$$Z - \text{score} = \frac{(x - \mu)}{\sigma}. \quad (1)$$

3.2. Steps Involved in IGWO-FS Technique. Once the input data is preprocessed, the next stage is to choose an optimal subset of features. The GWO algorithm is naturally inspired by the behavior and social leadership of the grey wolves [19]. The population of wolves can be classified into alpha, beta, delta, and omega for establishing the social hierarchy of wolves. The fittest solution is called alpha (α), whereas beta (β) and delta (δ) represent the 2nd and 3rd most efficient options, respectively. Omega (ω) represents semblance of a hopeful solution. The arithmetical expression of readapting position 0 is shown as follows:

$$\begin{aligned} \vec{D}_\alpha &= |\vec{C}_1 \cdot \vec{X}_\alpha - \vec{X}|, \\ \vec{D}_\beta &= |\vec{C}_2 \cdot \vec{X}_\beta - \vec{X}|, \\ \vec{D}_\delta &= |\vec{C}_3 \cdot \vec{X}_\delta - \vec{X}|, \end{aligned} \quad (2)$$

$$\begin{aligned} \vec{X}_1 &= \vec{X}_\alpha - \vec{A}_1 \cdot (\vec{D}_\alpha), \\ \vec{X}_2 &= \vec{X}_\beta - \vec{A}_2 \cdot (\vec{D}_\beta), \\ \vec{X}_3 &= \vec{X}_\delta - \vec{A}_3 \cdot (\vec{D}_\delta), \end{aligned} \quad (3)$$

where \vec{X}_α denotes the location of the alpha, \vec{X}_β represent the location of the beta, \vec{X}_δ indicates the location of the delta, and \vec{C}_1, \vec{C}_2 , and \vec{C}_3 and \vec{A}_1, \vec{A}_2 , and \vec{A}_3 signifies random vector, that is, the location of the existing solution, and shows the amount of iterations. It can be expressed in the following equation:

$$\vec{T}(u+1) = \vec{T}_p(u) + \vec{B} \cdot \vec{E}, \quad (4)$$

where \vec{E} is represented in equation (3), u indicates the iteration number, \vec{B}, \vec{D} denote the vector coefficient, and \vec{T}_p and \vec{T} represent the praise and grey wolf locations. The \vec{B}, \vec{D} vectors are calculated in the following equation:

$$\begin{aligned} \vec{E} &= |\vec{D} \cdot \vec{T}_p(u) - \vec{T}(u)|, \\ \vec{B} &= 2\vec{b} \cdot \vec{s}_1 - \vec{b}, \\ \vec{D} &= 2\vec{s}_2. \end{aligned} \quad (5)$$

s_1 and s_2 denote vectors with arbitrary numbers within $[0, 1]$ and \vec{b} parameter is linearly reduced from 2 to 0 all

over the iteration. Usually, the alpha is responsible for the chase. To change the position with the optimal searching agent position, the first three optimal solutions attained up until now compel another searching agent. Then, the wolves position can be upgraded as follows:

$$\begin{aligned}\vec{E}_\alpha &= |\vec{D}_1 \cdot \vec{t}_\alpha - \vec{Y}|, \vec{E}_\beta = |\vec{D}_2 \cdot \vec{t}_\beta - \vec{Y}|, \vec{E}_\delta = |\vec{D}_3 \cdot \vec{t}_\delta - \vec{Y}|, \\ \vec{t}_1 &= |\vec{t}_\alpha - \vec{B}_1 \cdot \vec{E}_\alpha|, \vec{t}_2 = |\vec{t}_\beta - \vec{B}_2 \cdot \vec{E}_\beta|, \vec{t}_3 = |\vec{t}_\delta - \vec{B}_3 \cdot \vec{E}_\delta|, \\ \vec{t}(u+1) &= \frac{\vec{t}_1 + \vec{t}_2 + \vec{t}_3}{3}.\end{aligned}\quad (6)$$

The variable b governs the balance between exploration and exploitation. Here, the variable b can be updated linearly in all the iterations, which ranges from 2 to 0, with u being the iteration number and m_i be the overall iterations allowed for the optimization:

$$\vec{b} = 2 - u \cdot \frac{2}{m_i}. \quad (7)$$

The wolf's location reflects attribute set selection and the solution space can be made by each probable attribute selection. The fitness function of the IGWO-FS technique can be utilized for determining whether an attribute subset would be selected or not:

$$\text{Fitness} = \alpha^* \gamma_S(E) + \beta^* \frac{|D - S|}{|D|}. \quad (8)$$

$|S|$ represents the selected attribute length subset and $\gamma_S(E)$ denotes the classification quality of attribute set S in relation to decision E . The overall amount of quality indicates the letter $|D|$ $\alpha \in 0, 1$ and $\beta = 1 - \alpha$, are two respective values for the attribute subset length and classification quality, respectively. Both have dissimilar implications for the attribute reduction task. The set $\alpha = 0.9, \beta = 0.1$ and attribute subset length are less important when compared to the quality classification. The higher ensure that the ideal location is a rough set reduction as a minimum. The fitness function evaluates the quality of location. After defining the fitness level, important feature is taken as well as removing the unwanted feature.

The performance of the GWO algorithm can be improved by the design of IGWO algorithm with the inclusion of adaptive β -hill climbing ($A\beta HC$). It is a recently presented local search-based technique, that is, basically, a modified version of β -hill climbing (βHC) [20]. The study has established that $A\beta HC$ gives optimum performance than several other famous local search techniques. For boosting the techniques exploitation capability and the quality of last solutions, $A\beta HC$ has been combined with the fundamental GTO for support searching the neighborhoods of optimum solution under this study. The explanation of $A\beta HC$ has been demonstrated mathematically as follows:

In order to provide an existing solution $X_i = (x_{i,1}, x_{i,2}, \dots, x_{i,D})$, $A\beta HC$ is iteratively created an improved solution $X''_i = (x''_{i,1}, x''_{i,2}, \dots, x''_{i,D})$ on the

fundamental of 2 control operators: \mathcal{N} -operator and β -operator. The \mathcal{N} -operator primarily transfers X_i to a novel neighborhood solution.

$X'_i = (x'_{i,1}, x'_{i,2}, \dots, x'_{i,D})$ that is demonstrated in equations (9) and (10) as follows:

$$x'_{ij} = x_{ij} \pm U(0, 1) \times \mathcal{N}, j = 1, 2, \dots, D, \quad (9)$$

$$\mathcal{N}(t) = 1 - \frac{t^{1/K}}{\text{Maxiter}^{1/K}}, \quad (10)$$

where $U(0, 1)$ refers the arbitrary number between the interval of 0 and 1, x_{ij} represents the value of decision variable from the j^{th} dimensional, t stands for the existing iteration, Maxiter denotes the maximal count of iterations, N signifies the bandwidth distance amongst the existing solution and their neighbor, D refers to the spatial dimensionality, and the parameter K is a constant.

3.3. Optimal ANFIS-Based Classification. At the final stage, the OANFIS model has been employed for the detection and classification of gene expression data into multiple classes. A network with 2 inputs, x and y and one output, f is considered. The ANFIS is a fuzzy Sugeno method. For presenting the ANFIS structure, 2 fuzzy if-then rules dependent upon a first-order Sugeno method are assumed as follows [21]:

- (i) Rule 1: if x is A_1 and y is B_1 , then $f_1 = p_1x + q_1y + r_1$
- (ii) Rule 2: if x is A_2 and y is B_2 , then $f_2 = p_2x + q_2y + r_2$

Where x and y are inputs, A_1 and B_i imply fuzzy sets, f_i is the output of fuzzy system, and p_i, q_i , and r_i represent the design parameters that are defined in the training procedure. The ANFIS structure for implementing these 2 rules, whereas a circle represents the set node and a square refers the adaptive node. The ANFIS infrastructure has 5 layers. Figure 2 showss the framework of ANFIS.

Layer 1. All nodes from layer1 are adaptation nodes. The resultant of layer 1 is are fuzzy membership grade of the inputs that are provided as follows:

$$\begin{aligned}O_i^1 &= \mu_{A_i}(x) i = 1, 2, \\ O_i^1 &= \mu_{B_{i-2}}(x) i = 3, 4,\end{aligned}\quad (11)$$

where x and y refer the inputs to node i , A refers the linguistic label, and $\mu_{A_i}(x)$ and $\mu_{B_{i-2}}(x)$ are some fuzzy membership functions. Generally, $\mu_{A_i}(x)$ is chosen as

$$\mu_{A_i} = \frac{1}{1 + \left\{ \left[\frac{(x - c_i)}{a_i} \right]^2 \right\}^{b_i}}, \quad (12)$$

where a_i, b_i , and c_i are the parameters of membership bell-shaped function.

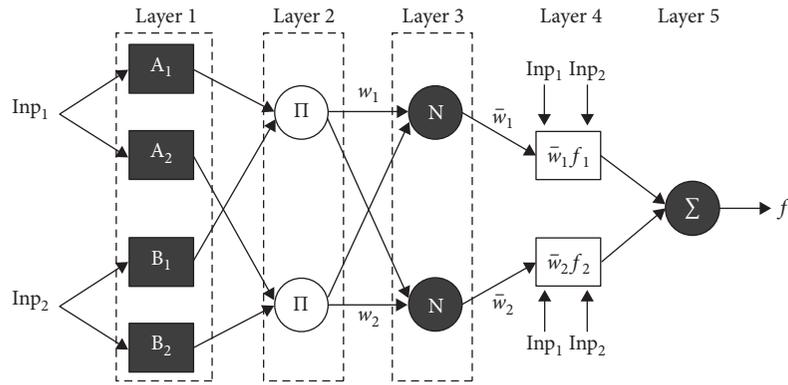


FIGURE 2: Structure of ANFIS.

TABLE 1: Dataset details.

Dataset	Leukemia	Prostate	DLBCL Stanford	Colon Cancer
No. of genes	7129	12600	4026	2000
Class 0	27	52	24	40
Class 1	11	50	23	22
Total no. of samples	38	102	47	62

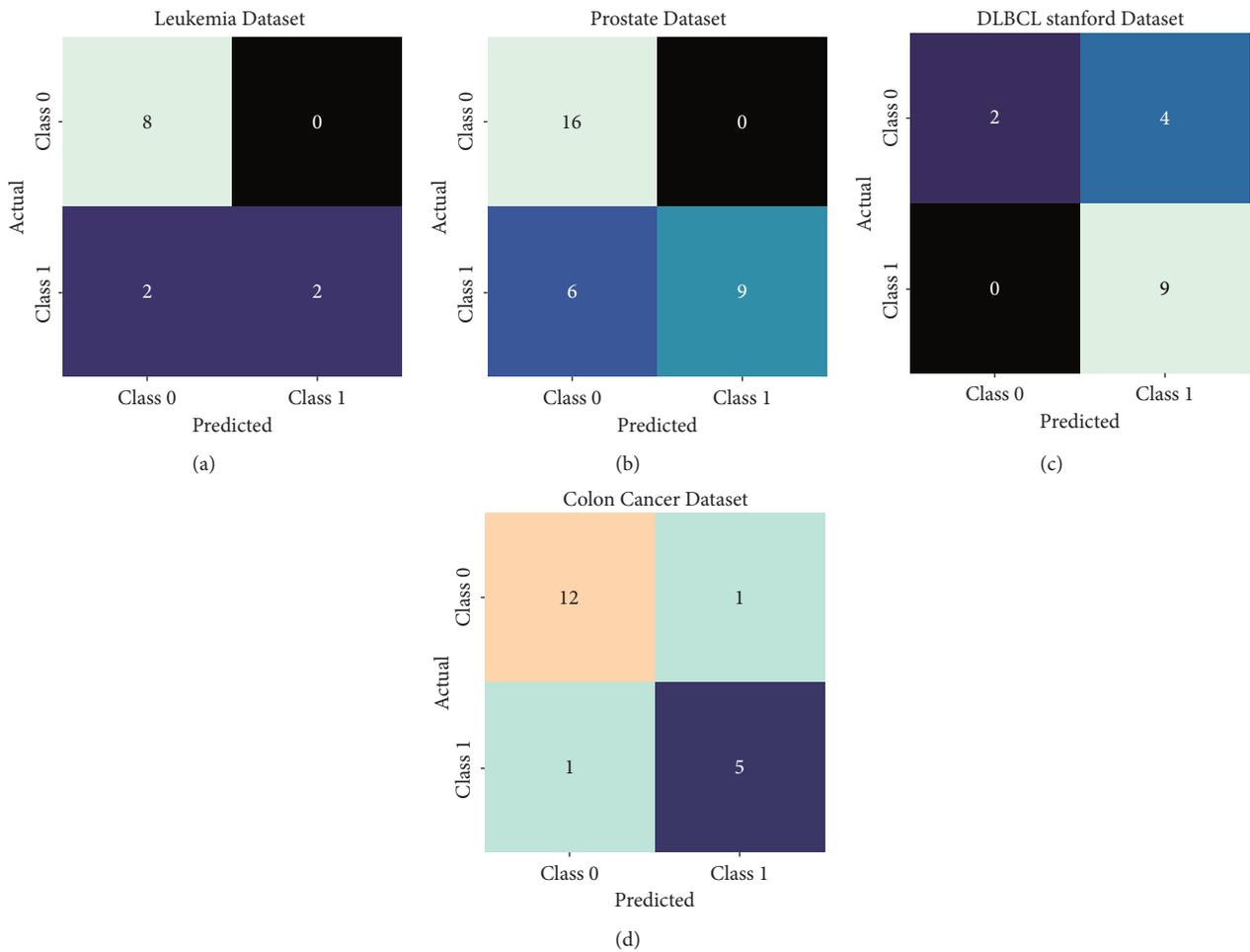


FIGURE 3: Confusion matrix of FSS-OANFIS technique under four datasets.

TABLE 2: Result analysis of FSS-OANFIS technique with different measures and datasets.

Class labels	Accuracy	Recall	Specificity	F-score	G-measure
Leukemia dataset					
Class 0	83.33	100.00	50.00	88.89	89.44
Class 1	83.33	50.00	100.00	66.67	70.71
Average	83.33	75.00	75.00	77.78	80.08
Prostate dataset					
Class 0	80.65	100.00	60.00	84.21	85.28
Class 1	80.65	60.00	100.00	75.00	77.46
Average	80.65	80.00	80.00	79.61	81.37
Stanford dataset					
Class 0	73.33	33.33	100.00	50.00	57.74
Class 1	73.33	100.00	33.33	81.82	83.21
Average	73.33	66.67	66.67	65.91	70.47
Colon dataset					
Class 0	89.47	92.31	83.33	92.31	92.31
Class 1	89.47	83.33	92.31	83.33	83.33
Average	89.47	87.82	87.82	87.82	87.82

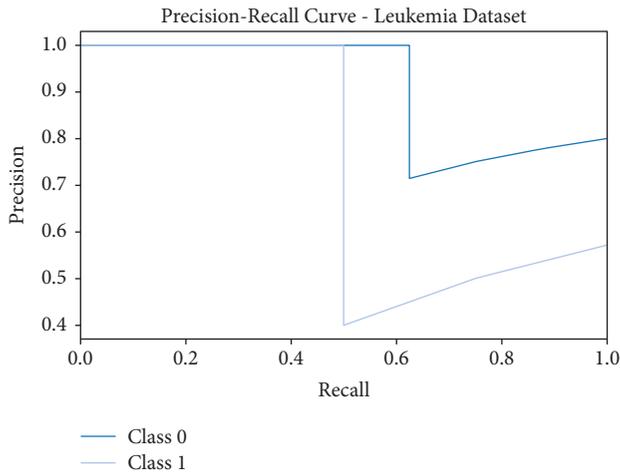


FIGURE 4: Precision-recall analysis of FSS-OANFIS technique under the Leukemia dataset.

Layer 2. The node of this layer is labeled M , signifying that it can be implemented as an easy multiplier. The resultant of this layer is demonstrated as follows:

$$O_i^2 = w_i = \mu_{A_i}(x)\mu_{B_i}(y) \quad i = 1, 2. \quad (13)$$

Layer 3. It comprises set nodes which compute the ratio of firing strength of the rules as follows:

$$O_i^3 = \bar{w}_i = \frac{w_i}{w_1 + w_2} \quad i = 1, 2. \quad (14)$$

Layer 4. During this layer, the adaptive node is used. The resultants of this layer are calculated by the following equation:

$$O_i^4 = \bar{w}_i f_i = \bar{w}_i (p_i x + q_i y + r_i) \quad i = 1, 2. \quad (15)$$

\bar{w}_i signifies the normalized firing strength in layer 3.

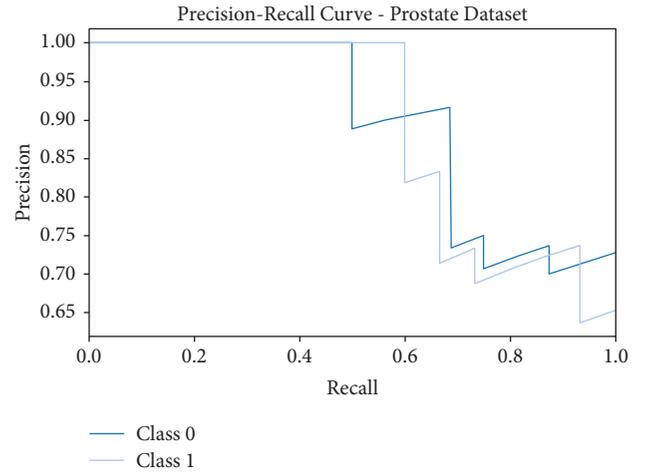


FIGURE 5: Precision-recall analysis of FSS-OANFIS technique under the Prostate dataset.

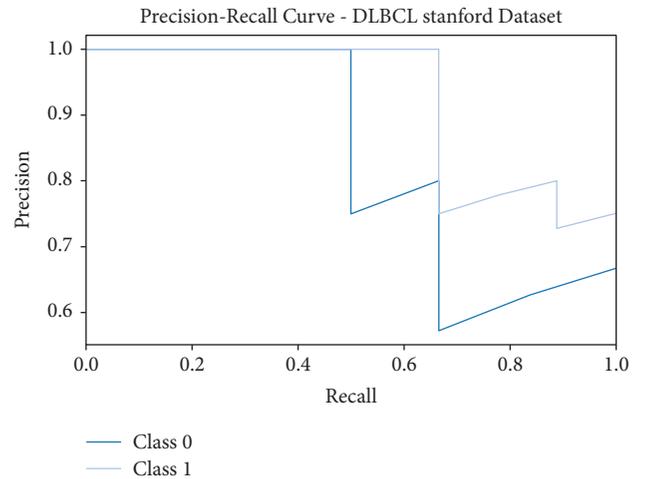


FIGURE 6: Precision-recall analysis of FSS-OANFIS technique under the DLBCL Stanford dataset.

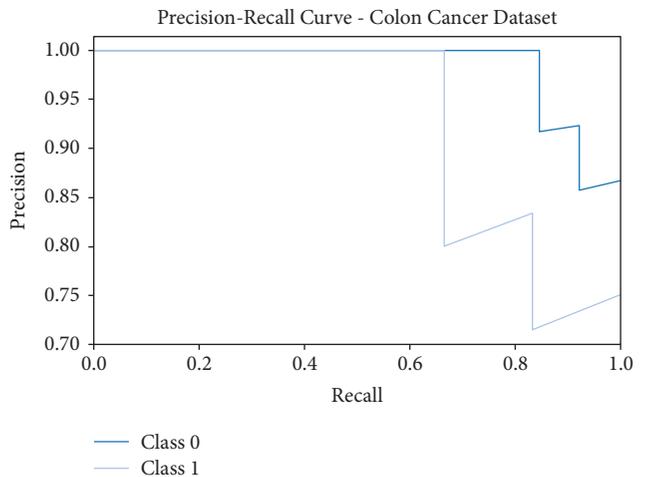


FIGURE 7: Precision-recall analysis of FSS-OANFIS technique under the Colon Cancer dataset.

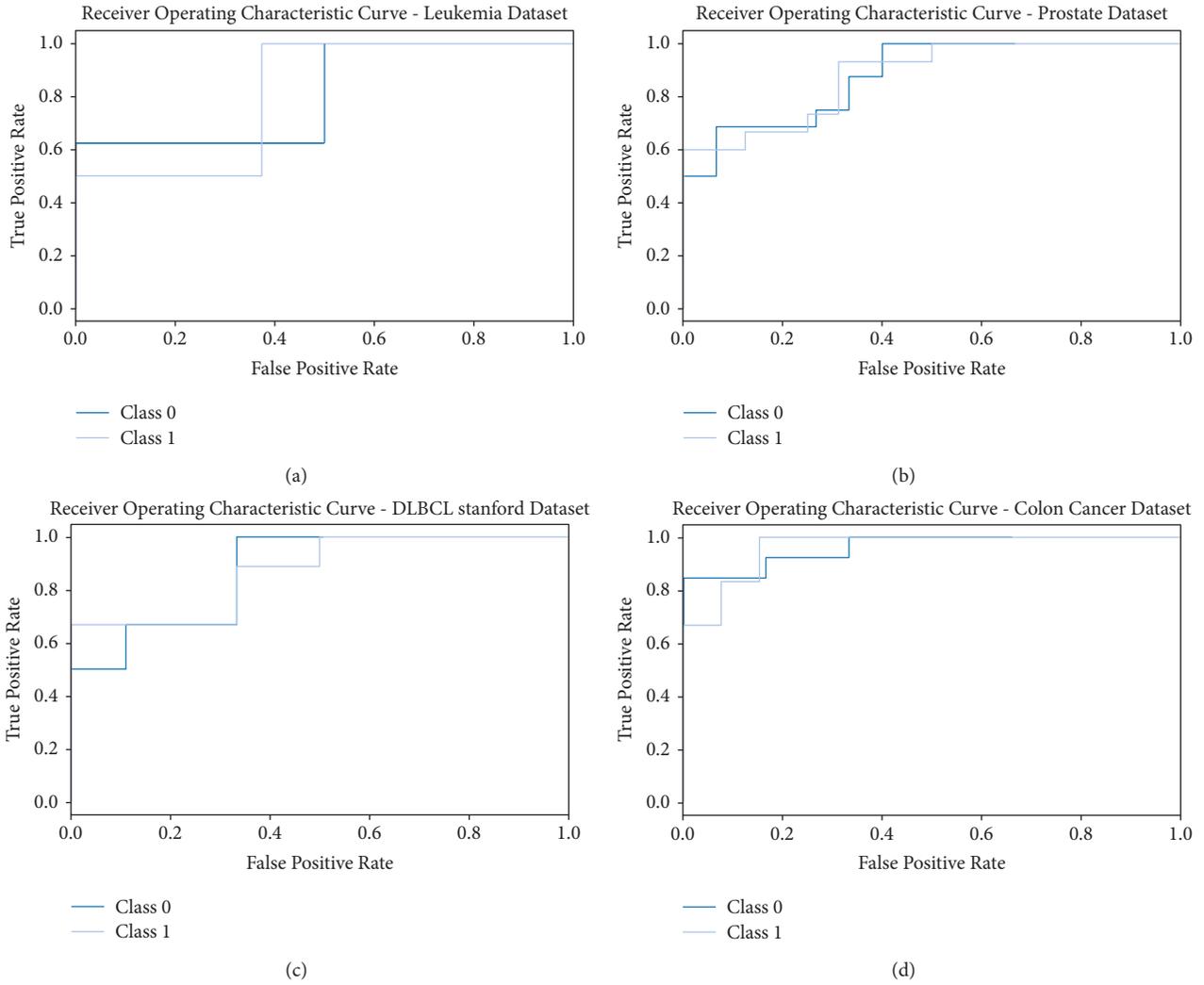


FIGURE 8: ROC analysis of FSS-OANFIS technique under different datasets.

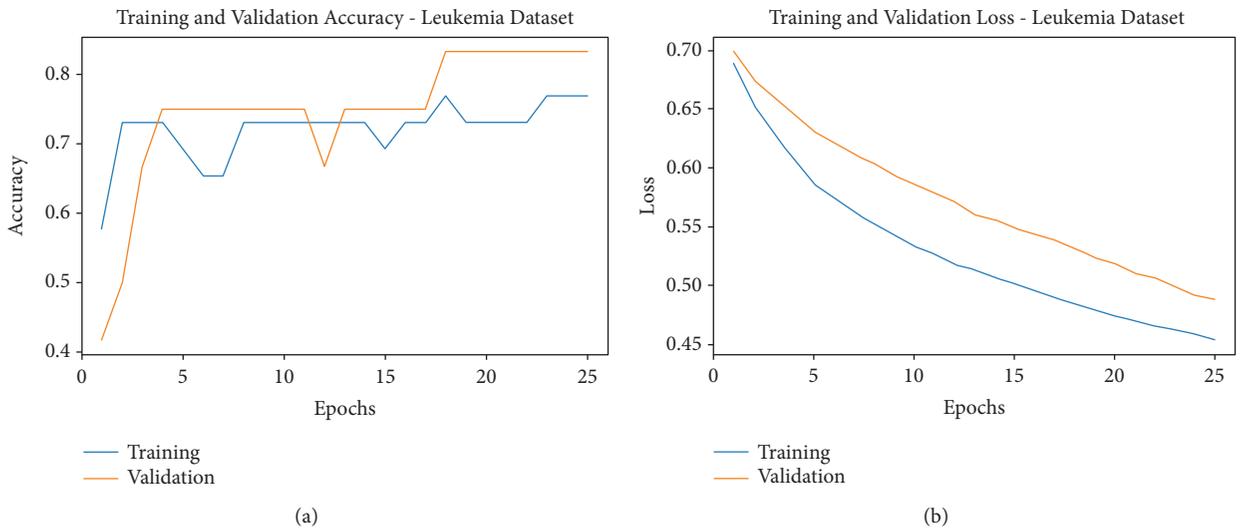
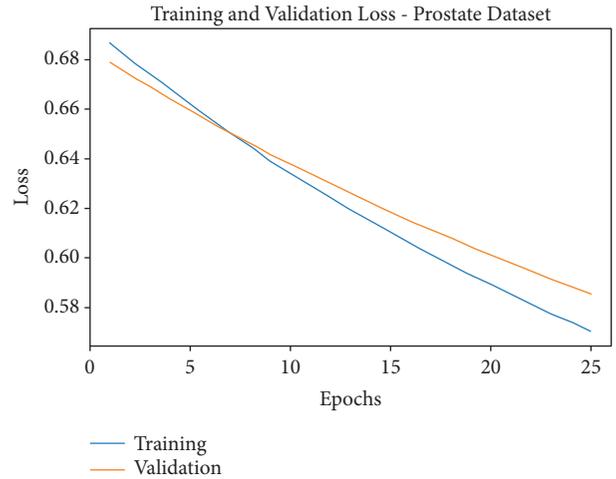
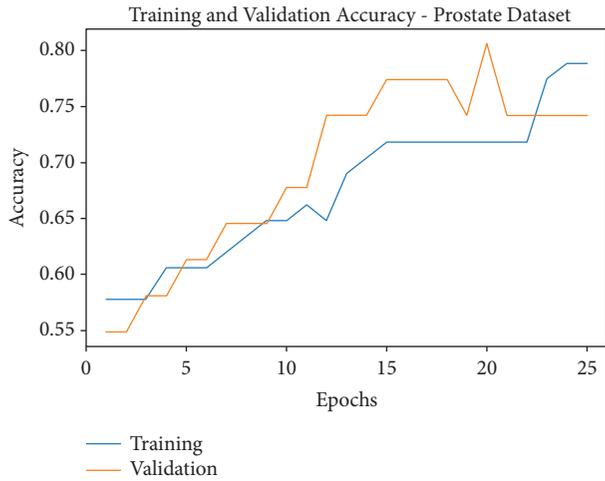
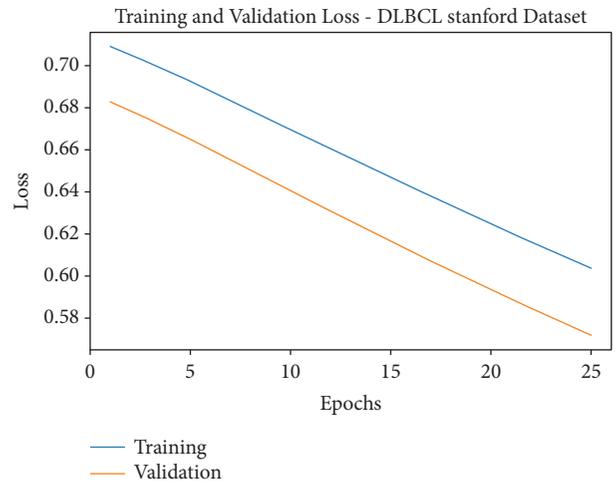
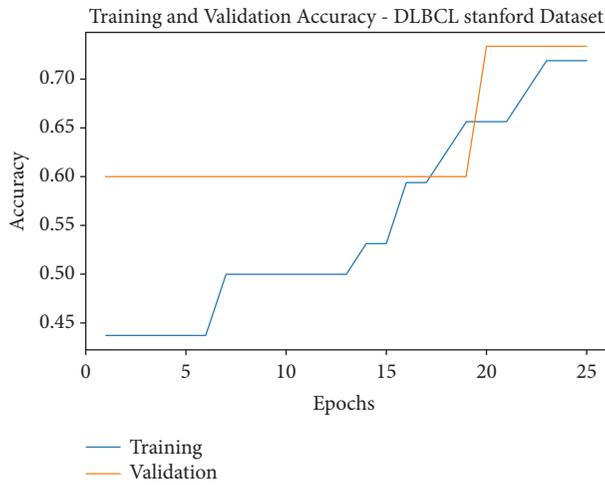


FIGURE 9: Continued.



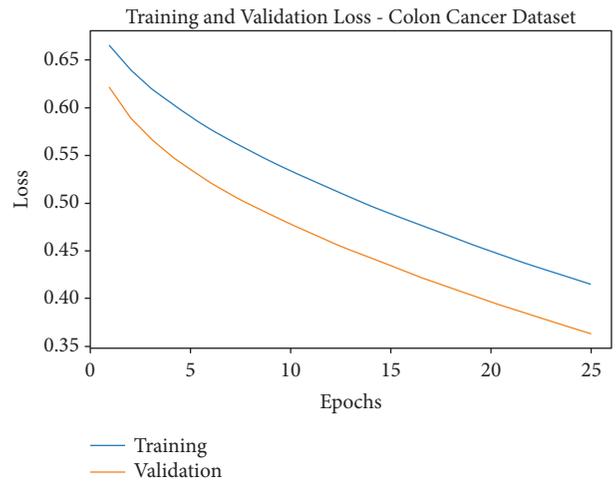
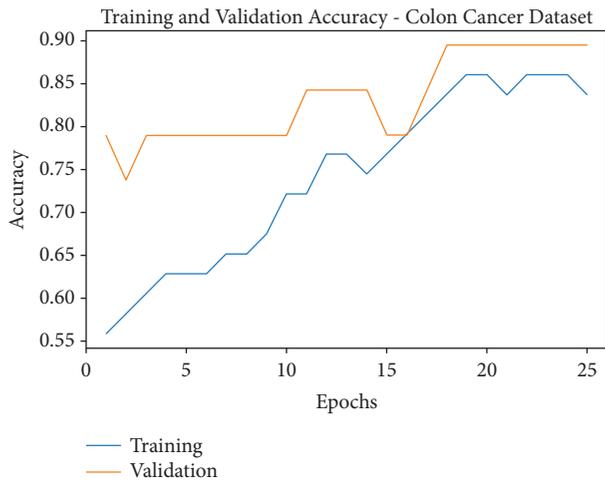
(c)

(d)



(e)

(f)



(g)

(h)

FIGURE 9: Accuracy and loss analysis of FSS-OANFIS technique under various datasets.

TABLE 3: Comparative analysis of FSS-OANFIS technique with existing approaches.

Methods	Accuracy	Sensitivity	Specificity	G-measure
Leukemia dataset				
FSS-OANFIS	83.33	75.00	75.00	80.08
AHSA-GS	75.49	69.66	74.81	45.94
PSO algorithm	80.59	74.95	73.96	68.07
DE algorithm	68.67	63.01	63.62	64.80
Prostate dataset				
FSS-OANFIS	80.65	80.00	80.00	81.37
AHSA-GS	71.19	53.82	79.79	79.84
PSO algorithm	68.78	63.63	70.15	66.01
DE algorithm	62.77	60.37	63.22	67.94
Stanford dataset				
FSS-OANFIS	73.33	66.67	66.67	70.47
AHSA-GS	71.27	62.64	65.15	68.82
PSO algorithm	72.80	60.82	60.17	61.74
DE algorithm	66.93	63.80	59.16	61.48
Colon dataset				
FSS-OANFIS	89.47	87.82	87.82	87.82
AHSA-GS	61.02	48.07	64.04	43.62
PSO algorithm	59.00	43.02	58.34	36.66
DE algorithm	50.38	33.63	38.76	58.07

Layer 5. The node executes the summation of every incoming signal. Therefore, an entire result of the model is provided as follows:

$$O_i^5 = \sum_i \bar{w}_i f_i = \frac{\sum_i w_i f_i}{\sum_i w_i}. \quad (16)$$

It could be realized that there are 2 adaptive layers under this ANFIS structure such as the 1st layer and 4th layer. During the 1st layer, there are 3 modifiable parameters $\{a, b, c_i\}$ that are connected to the input membership function. These parameters are usually named as premise parameters. During the 4th layer, there are also 3 modifiable parameters $\{p, q, r_i\}$ relating to the first-order polynomial. This parameter is supposed the consequent parameter.

For tuning the ANFIS parameters, the COA is applied to it. The COA is a mathematical model that depends on smart diversity [22]. Chasing, driving, attacking, and blocking are archived by four distinct kinds of chimps that are attained by chasers, drivers, attackers, and obstacles. These hunting steps are accomplished in two phases such as exploration and exploitation stages. The exploration phase involves chasing, driving, and blocking the prey. The exploitation phase should attack the prey, and the chasing and driving are characterized as follows:

$$d = |c \cdot x_{\text{prey}}(t) - m \cdot x_{\text{chimp}}(t)|, \quad (17)$$

$$x_{\text{chimp}}(t+1) = x_{\text{prey}}(t) - a \cdot d,$$

where X_{prey} denotes the vector of prey location, x_{chimp} indicates the vector of chimp location, t denotes the amount of

present iterations, a , c , and m represent coefficient vectors and are calculated as follows:

$$\begin{aligned} a &= 2 \cdot f \cdot r_1 \cdot f, \\ c &= 2 \cdot r_2, \\ m &= \text{chaotic} - \text{value}, \end{aligned} \quad (18)$$

where f declined nonlinearly from 2.5 to 0, r_1 and r_2 denote the random value within $[0, 1]$, and m represents the chaotic vector. The dynamic coefficient f is selected for distinct slopes and curves; therefore, chimps employ distinct capabilities for searching the prey. Chimps upgrade the position according to the other chimps, and the arithmetical expression can be given in the following equation:

$$\begin{aligned} d_{\text{Attacker}} &= |c_1 x_{\text{Attacker}} - m_{1x}|, \\ d_{\text{Barrier}} &= |c_2 x_{\text{Barrier}} - m_{2x}|, \\ d_{\text{Chaser}} &= |c_3 x_{\text{Chaser}} - m_3 x| \\ d_{\text{Driver}} &= |c_4 x_{\text{Driver}} - m_4 x|, \\ x_1 &= x_{\text{Attacker}} - a_1 (d_{\text{Attacker}}), \\ x_2 &= x_{\text{Barrier}} - a_2 (d_{\text{Barrier}}), \\ x_3 &= x_{\text{Chaser}} - a_3 (d_{\text{Chaser}}), \\ x_4 &= x_{\text{Driver}} - a_4 (d_{\text{Driver}}), \\ x(t+1) &= \frac{x + x_2 + x_3 + x_4}{4}. \end{aligned} \quad (19)$$

4. Experimental Validation

In this section, the experimental validation of the FSS-OANFIS model has been performed using four benchmark datasets [23–26]. The details of the dataset are given in Table 1. The results are investigated and the outcomes are assessed in terms of different measures. For experimental validation, a 10-fold cross-validation process is utilized.

4.1. Results Analysis of Proposed Model. Figure 3 illustrates a set of confusion matrices offered by the FSS-OANFIS model on test datasets. The figure reported that the FSS-OANFIS model has properly recognized the class labels on all datasets. For instance, on the Leukemia dataset, the FSS-OANFIS model has identified 8 samples in class 0 and 2 samples in class 1. In addition, on the Prostate dataset, the FSS-OANFIS system has identified 16 samples in class 0 and 9 samples in class 1. Also, on the DLBCL Stanford dataset, the FSS-OANFIS approach has identified 2 samples in class 0 and 9 samples in class 1. Besides, on the Colon Cancer dataset, the FSS-OANFIS algorithm has identified 12 samples in class 0 and 5 samples in class 1.

Table 2 provides the overall classification outcomes of the FSS-OANFIS model on the test datasets. The experimental outcomes pointed out that the FSS-OANFIS model

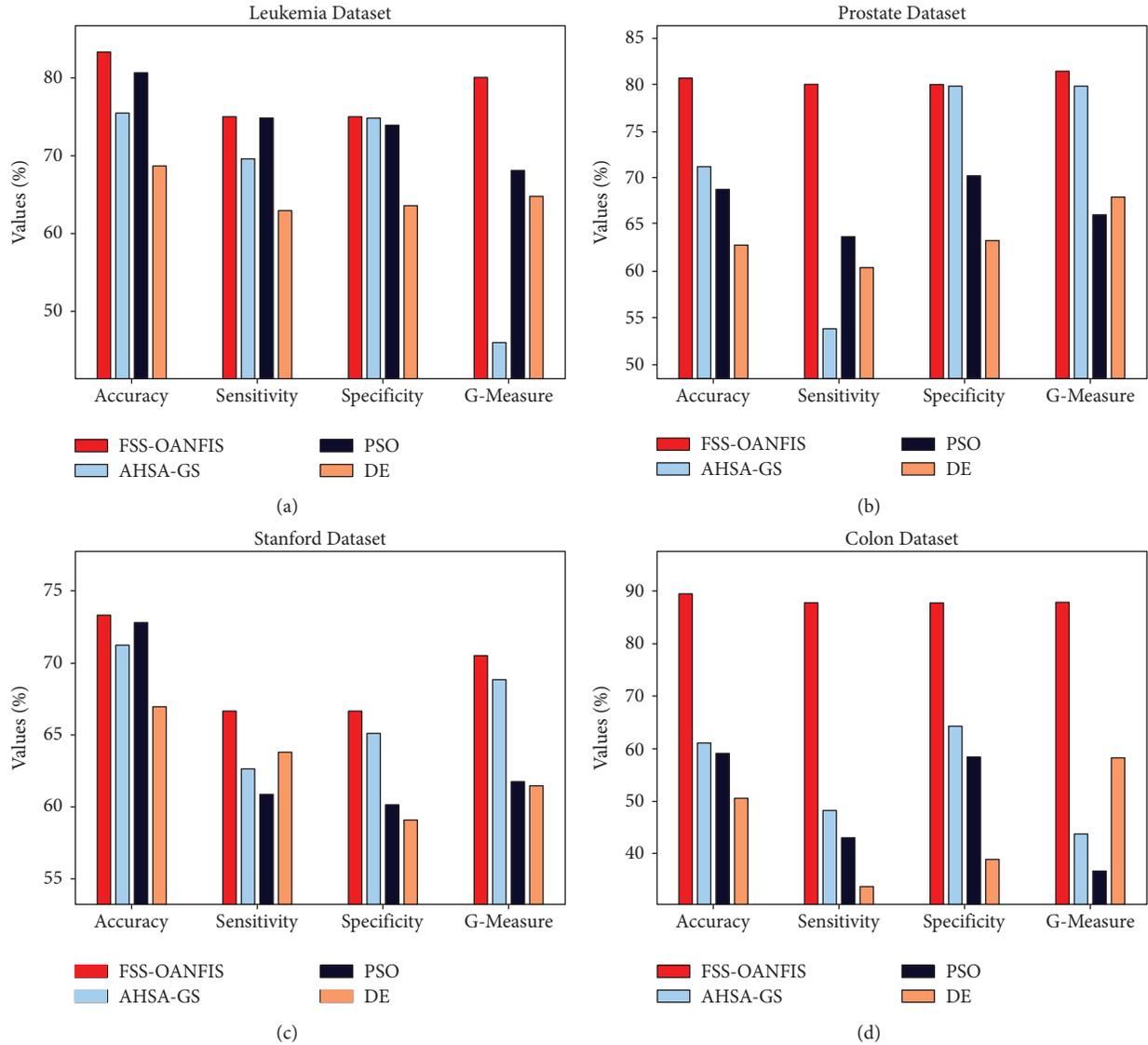


FIGURE 10: Comparative analysis of FSS-OANFIS technique with existing approaches.

has offered effective outcomes on all the datasets applied. For instance, with the Leukemia dataset, the FSS-OANFIS model has resulted in average acc_y , $recal$, $spec_y$, F_{score} , and $G_{measure}$ of 83.33%, 75%, 75%, 77.78%, and 80.08%, respectively. Following this, with the Prostate dataset, the FSS-OANFIS methodology has resulted in average acc_y , $recal$, $spec_y$, F_{score} , and $G_{measure}$ of 80.65%, 80%, 80%, 79.61%, and 81.37%, respectively. Meanwhile, with the DLBCL Stanford dataset, the FSS-OANFIS algorithm has resulted in average acc_y , $recal$, $spec_y$, F_{score} , and $G_{measure}$ of 73.33%, 66.67%, 66.67%, 65.91%, and 70.47%, respectively. Eventually, with the Colon Cancer dataset, the FSS-OANFIS technique has resulted in average acc_y , $recal$, $spec_y$, F_{score} , and $G_{measure}$ of 89.47%, 87.82%, 87.82%, 87.82%, and 87.82%, respectively.

Figure 4 shows the precision-recall curves offered by the FSS-OANFIS model on the test Leukemia dataset. The figure indicated that the FSS-OANFIS model has depicted effective precision-recall values on the classification of two classes, namely, class 0 and class 1 on the test Leukemia dataset.

Next, Figure 5 shows the precision-recall curves offered by the FSS-OANFIS model on the test Prostate dataset. The figure revealed that the FSS-OANFIS technique has depicted effective precision-recall values on the classification of two classes, namely, class 0 and class 1 on the test Prostate dataset.

Similarly, Figure 6 shows the precision-recall curves offered by the FSS-OANFIS system on the test DLBCL Stanford dataset. The figure exposed that the FSS-OANFIS model has depicted effective precision-recall values on the classification of two classes, namely, class 0 and class 1 on the test DLBCL Stanford dataset.

Figure 7 shows the precision-recall curves offered by the FSS-OANFIS method on the test Colon Cancer dataset. The figure indicated that the FSS-OANFIS approach has depicted effective precision-recall values on the classification of two classes, namely, class 0 and class 1 on the test Colon Cancer dataset.

A brief ROC investigation of the FSS-OANFIS model on the distinct four datasets is described in Figure 8. The results

indicated that the FSS-OANFIS technique has exhibited its ability in categorizing two different classes such as class 0 and 1 on the test four datasets.

Figure 9 shows the accuracy and loss graph analysis of the ODBN-IDS technique under four datasets. The results show that the accuracy value tends to increase and the loss value tends to decrease with an increase in epoch count. It is also observed that the training loss is low and validation accuracy is high under four datasets.

4.2. Discussion. Finally, a detailed comparative study of the FSS-OANFIS model with recent methods on distinct datasets is shown in Table 3 and Figure 10 [27]. The experimental results indicated that the FSS-OANFIS model has shown effectual outcomes under all datasets. For instance, with the Leukemia dataset, the DE and AHSA-GS models have depicted lower performance over the other methods. Though the PSO algorithm has resulted in slightly reasonable performance with $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 80.59%, 74.95%, 73.96%, and 68.07%, the FSS-OANFIS model has resulted in higher $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 83.33%, 75%, 75%, and 80.08%, respectively.

At the same time, with the Prostate dataset, the DE and AHSA-GS models have depicted lower performance over the other methods. Likewise, the PSO algorithm has resulted in somewhat reasonable performance with $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 68.78%, 63.63%, 70.15%, and 66.01% and the FSS-OANFIS methodology has resulted in superior $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 80.65%, 80%, 80%, and 81.37%, respectively. In addition, with the DLBCL Stanford dataset, the DE and AHSA-GS techniques have showcased lesser performance over the other methods. Though the PSO algorithm has resulted in slightly reasonable performance with $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 72.80%, 60.82%, 60.17%, and 61.74%, the FSS-OANFIS approach has resulted in higher $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 73.33%, 66.67%, 66.67%, and 70.47%, respectively.

Along with that, with the Colon Cancer dataset, the DE and AHSA-GS models have portrayed lower performance over the other methods. But, the PSO approach has resulted in slightly reasonable performance with $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 59%, 43.02%, 58.34%, and 36.66%, and the FSS-OANFIS system has resulted in superior $accu_y$, $sens_y$, $spec_y$, and $G_{measure}$ of 89.47%, 87.80%, 87.82%, and 87.82%, respectively. After examining the results and discussion, it is apparent that the FSS-OANFIS model has accomplished maximum performance in the microarray gene expression classification process.

5. Conclusion

In this study, a new FSS-OANFIS model has been developed for microarray gene expression data classification. The presented FSS-OANFIS model encompasses a series of processes, namely, data pre-processing, IGWO-FS-based election of features, ANFIS classification, and COA-based parameter optimization. The application of IGWO-FS and

COA techniques helps in accomplishing enhanced microarray gene expression classification outcomes. For examining the enhanced outcomes of the FSS-OANFIS model, a wide range of simulations were performed on distinct datasets. The experimental results indicated that the FSS-OANFIS model has resulted in enhanced performance over the recent approaches. In future, the feature reduction and clustering approaches can be integrated to enhance gene expression classification outcomes.

Data Availability

Data are available and can be provided upon direct request to the corresponding author.

Ethical Approval

This article does not contain any studies with human participants performed by any of the authors.

Consent

Not applicable.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

The authors extend their appreciation to the Deanship of Scientific Research at King Khalid University for funding this work under grant number RGP 2/180/43, Princess Nourah bint Abdulrahman University supporting project number PNURSP2022R151, Princess Nourah bint Abdulrahman University, Riyadh, Saudi Arabia. The authors would like to thank the Deanship of Scientific Research at Umm Al-Qura University for supporting this work under grant code: 22UQU4340237DSR15.

References

- [1] M. Maniruzzaman, M. Jahanur Rahman, B. Ahammed et al., "Statistical characterization and classification of colon microarray gene expression data using multiple machine learning paradigms," *Computer Methods and Programs in Biomedicine*, vol. 176, pp. 173–193, 2019.
- [2] R. Tabares-Soto, S. Orozco-Arias, V. Romero-Cano, V. Segovia Bucheli, J. L. Rodríguez-Sotelo, and C. F. Jiménez-Varón, "A comparative study of machine learning and deep learning algorithms to classify cancer types based on microarray gene expression data," *PeerJ Computer Science*, vol. 6e270 pages, 2020.
- [3] J. M. Franks, G. Cai, and M. L. Whitfield, "Feature specific quantile normalization enables cross-platform classification of molecular subtypes using gene expression data," *Bioinformatics*, vol. 34, no. 11, pp. 1868–1874, 2018.
- [4] E. H. Mahood, L. H. Kruse, and G. D. Moghe, "Machine learning: a powerful tool for gene function prediction in plants," *Applications in Plant Sciences*, vol. 8, no. 7, Article ID e11376, 2020.

- [5] N. Borisov, V. Tkachev, M. Suntsova et al., "A method of gene expression data transfer from cell lines to cancer patients for machine-learning prediction of drug efficiency," *Cell Cycle*, vol. 17, no. 4, pp. 486–491, 2018.
- [6] E. a. Alhenawi, R. Al-Sayyed, A. Hudaib, and S. Mirjalili, "Feature selection methods on gene expression microarray data for cancer classification: a systematic review," *Computers in Biology and Medicine*, vol. 140, Article ID 105051, 2022.
- [7] L. Rukhsar, W. H. Bangyal, M. S. Ali Khan, A. A. Ag Ibrahim, K. Nisar, and D. B. Rawat, "Analyzing RNA-seq gene expression data using deep learning approaches for cancer classification," *Applied Sciences*, vol. 12, no. 4, 1850 pages, 2022.
- [8] T. Thakur, I. Batra, M. Luthra et al., "Gene expression-assisted cancer prediction techniques," *Journal of Healthcare Engineering*, vol. 2021, Article ID 4242646, 10 pages, 2021.
- [9] M. Abd-Elnaby, M. Alfonse, and M. Roushdy, "Classification of breast cancer using microarray gene expression data: a survey," *Journal of Biomedical Informatics*, vol. 117, Article ID 103764, 2021.
- [10] K. Showalter, R. Spiera, C. Magro et al., "Machine learning integration of scleroderma histology and gene expression identifies fibroblast polarisation as a hallmark of clinical severity and improvement," *Annals of the Rheumatic Diseases*, vol. 80, no. 2, pp. 228–237, 2021.
- [11] J. P. Sarkar, I. Saha, A. Sarkar, and U. Maulik, "Machine learning integrated ensemble of feature selection methods followed by survival analysis for predicting breast cancer subtype specific miRNA biomarkers," *Computers in Biology and Medicine*, vol. 131, Article ID 104244, 2021.
- [12] S. M. Ayyad, A. I. Saleh, and L. M. Labib, "Gene expression cancer classification using modified K-Nearest Neighbors technique," *Biosystems*, vol. 176, pp. 41–51, 2019.
- [13] Y. Kong and T. Yu, "A deep neural network model using random forest to extract feature representation for gene expression data classification," *Scientific Reports*, vol. 8, no. 1, pp. 16477–16479, 2018.
- [14] A. K. Dwivedi, "Artificial neural network model for effective cancer classification using microarray gene expression data," *Neural Computing & Applications*, vol. 29, no. 12, pp. 1545–1554, 2018.
- [15] A. K. Shukla, "Feature selection inspired by human intelligence for improving classification accuracy of cancer types," *Computational Intelligence*, vol. 37, no. 4, pp. 1571–1598, 2021.
- [16] M. Rostami, S. Forouzandeh, K. Berahmand, M. Soltani, M. Shahsavari, and M. Oussalah, "Gene selection for microarray data classification via multi-objective graph theoretic-based method," *Artificial Intelligence in Medicine*, vol. 123, Article ID 102228, 2022.
- [17] N. Bhui, "Ensemble of deep learning approach for the feature selection from high-dimensional microarray data," in *Proceedings of the International Conference on Paradigms of Communication, Computing and Data Sciences*, pp. 591–600, Springer, Singapore, January 2022.
- [18] A. Forestiero, G. Papuzzo, R. De Simone, and R. Varchera, "A microarray analysis technique using a self-organizing multi-agent approach," *Methods in Molecular Biology*, vol. 2401, pp. 39–50, 2022.
- [19] D. S. Wankhede and R. Selvarani, "Dynamic architecture based deep learning approach for glioblastoma brain tumor survival prediction," *Neuroscience Informatics*, vol. 2, no. 4, Article ID 100062, 2022.
- [20] Y. Xiao, X. Sun, Y. Guo, S. Li, Y. Zhang, and Y. Wang, "An improved gorilla troops optimizer based on lens opposition-based learning and adaptive 13-hill climbing for global optimization," *Cmes-Computer Modeling In Engineering & Sciences*, 2022.
- [21] M. Khishe and M. R. Mosavi, "Chimp optimization algorithm," *Expert Systems with Applications*, vol. 149, Article ID 113338, 2020.
- [22] M. Hosseini and M. Zekri, "Review of medical image classification using the adaptive neuro-fuzzy inference system," *Journal of Medical Signals & Sensors*, vol. 2, no. 1, 49 pages, 2012.
- [23] A. A. Alizadeh, M. B. Eisen, R. E. Davis et al., "Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling," *Nature*, vol. 403, no. 6769, pp. 503–511, 2000.
- [24] U. Alon, N. Barkai, D. A. Notterman et al., "Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays," *Proceedings of the National Academy of Sciences*, vol. 96, no. 12, pp. 6745–6750, 1999.
- [25] D. Singh, P. G. Febbo, K. Ross et al., "Gene expression correlates of clinical prostate cancer behavior," *Cancer Cell*, vol. 1, no. 2, pp. 203–209, 2002.
- [26] T. R. Golub, D. K. Slonim, P. Tamayo et al., "Molecular classification of cancer: class discovery and class prediction by gene expression monitoring," *Science (New York, N.Y.)*, vol. 286, no. 5439, pp. 531–537, 1999.
- [27] R. Dash, "An adaptive harmony search approach for gene selection and classification of high dimensional medical data," *Journal of King Saud University-Computer and Information Sciences*, vol. 33, no. 2, pp. 195–207, 2021.