

Connection Optimization of a Neural Emotion Classifier Using Hybrid Gravitational Search Algorithms

Mansour Sheikhan

Electrical Engineering Department
Islamic Azad University-South Tehran Branch
Tehran, Iran
msheikhn@azad.ac.ir

Mahdi Abbasnezhad Arabi

Electrical Engineering Department
Islamic Azad University-South Tehran Branch
Tehran, Iran
m_abasnejad@yahoo.com

Davood Gharavian

Electrical Engineering Department
Shahid Beheshti University
Tehran, Iran
d_gharavian@sbu.ac.ir

Received: May 17, 2014-Accepted: November 21, 2014

Abstract— Artificial neural network is an efficient model in pattern recognition applications, but its performance is heavily dependent on using suitable structure and connection weights. This paper presents a hybrid heuristic method for obtaining the optimal weight set and architecture of a feedforward neural emotion classifier based on Gravitational Search Algorithm (GSA) and its binary version (BGSA), respectively. By considering various features of speech signal and concatenating them to a principal feature vector, which includes frame-based Mel frequency cepstral coefficients and energy, a rich medium-size feature set is constructed. The performance of the proposed hybrid GSA-BGSA-neural model is compared with the hybrid of Particle Swarm Optimization (PSO) algorithm and its binary version (BPSO) used for such optimizations. In addition, other models such as GSA-neural hybrid and PSO-neural hybrid are also included in the performance comparisons. Experimental results show that the GSA-optimized models can obtain better results using a lighter network structure.

Keywords- emotion recognition; speech processing; neural network; connection optimization; structure optimization; gravitational search algorithm.

I. INTRODUCTION

Automatic human emotions recognition has attracted research efforts in the field of man-machine communication in recent years [1-5]. This task can be performed using speech or (and) image signals. It is noted that there are two important information sources in the speech signal: a) an explicit source which contains the linguistic content, b) an implicit source which carries the paralinguistic information

about the speaker. In the last four decades, several methods have been proposed for developing Automatic Speech Recognition (ASR) systems to extract linguistic information. Although decoding the paralinguistic information such as emotion needs more research efforts. The emotion recognizer has become an effective tool in human-computer interfacing applications such as computer tutorial [6], call-center [7], lie detecting [8], consumer

relationships and in-car boards [9], developing learning environments [10, 11].

In recent years, the research efforts on emotion recognition from speech have been focused on extracting reliable informative features [12-15], selecting appropriate feature set [16], and combining powerful classifiers to improve the performance of emotion detection systems in real-life applications [17, 18].

Several features have been extracted and experienced for emotion recognition from speech such as:

- Pitch frequency (F_0), Log Energy (LE), formant frequencies, and Mel-Frequency Cepstral Coefficients (MFCCs) [19, 20],
- F_0 , LE, formant frequencies, MFCCs, vocal tract cross-section areas (A_k), and speech rate [21, 22],
- Linear Prediction Coefficients (LPCs) and MFCCs [23],
- F_0 , LE, MFCCs, and LPCs [24],
- Zero Crossing Rate (ZCR), LE, F_0 , and Harmonics-to-Noise Ratio (HNR) [25],
- Harmony features which are based on the psychoacoustic harmony perception known from music theory [26],
- Statistics of MFCCs computed over three phoneme types (stressed vowels, unstressed vowels, and consonants) [27],
- Jitter, shimmer, LPCs, Linear Prediction Cepstral Coefficients (LPCCs), MFCCs, derivative of MFCCs (dMFCCs), second derivative of MFCCs (ddMFCCs), Log Frequency Power Coefficients (LFPCs), and Perceptual Linear Prediction Coefficients (PLPCs) [28],
- Modulation Spectral Features (MSFs) using an auditory filter-bank and a modulation filter-bank for speech analysis [29].

Similarly, different classification methods have been employed in this field such as K-Nearest Neighbor (KNN) [23, 30, 31], decision trees [32-34], Bayesian networks [34], Optimum Path Forest (OPF) [35], Hidden Markov Models (HMMs) [36], Gaussian Mixture Models (GMMs) [37], Support Vector Machines (SVMs) [38], Artificial Neural Networks (ANNs) [39-41], and hybrid approaches [42].

This paper presents a hybrid heuristic method to find the best weight set and architecture of a feedforward neural emotion classifier based on Gravitational Search Algorithm (GSA) and its binary version (BGSA), respectively. This hybrid model is called GSA-BGSA in this paper. By considering various supplementary features, based on the first three formants (F_1 , F_2 , and F_3) and F_0 , and concatenating them to a principal feature vector, which includes MFCCs (shown by c_i , $i=1,2,\dots,12$), LE, and their derivative (dc_i , dLE) and second derivative (ddc_i , $ddLE$), a rich medium-size feature vector was constructed in this study. So, a total of 55 features were extracted over Farsi sentences. Four classes of emotion were considered in this study: neutral, happiness, anger, and surprise.

The rest of paper is organized as follows: the background and related work on optimizing ANNs are reviewed in Section 2. Section 3 explains the detailed structure of the ANN model used in this study. The details of hybrid GSA-BGSA algorithm are presented in Section 4. The emotional speech dataset is introduced in Section 5. The experimental results are reported in Section 6 in which the performance of proposed method is compared with the hybrid of Particle Swarm Optimization (PSO) and its binary version (BPSO) used for such optimizations. In addition, other models such as GSA-neural hybrid, PSO-neural hybrid, and standard Error Back-Propagation (EBP) are also included in the performance comparisons. The paper is concluded in Section 7.

II. RELATED WORK ON OPTIMIZING ANNS

It is noted that ANN is a nature-based computing technique that has been developed as a parallel-distributed network model based on the biological learning process of human brain. The mostly used training algorithm for ANNs is the EBP algorithm, which is a gradient-based method. However, some inherent problems exist in the EBP algorithm. One of these problems is trapping in local minima, especially for nonlinearly separable pattern classification problems or complex function approximation problems [43]. In addition, the training performance is sensitive to the choice of algorithm's parameters and initial values of weights. In other words, selecting the appropriate network architecture and weight parameters strongly affect the convergent behavior of the EBP algorithm [44].

Several approaches have been proposed with the aim of introducing systematic and automatic ways for tuning the network structure and the training parameters of ANNs. These approaches can be categorized as follow:

- Statistical or empirical methods that have been used to study the role of an ANN's internal parameters and choosing appropriate values for it based on the model's performance [45-47]. For example, Salchenberger *et al.* [48] suggested the number of hidden node as follows:

$$n_{hidden_layer} = 0.75 \times n_{input} \quad (1)$$

and Subramanian *et al.* [49] recommended the number of nodes in a single hidden layer ANN as follows:

$$n_{hidden-layer} = n_{input} + n_{output} + 1 \quad (2)$$

- Constructive and/or pruning algorithms that trace the network performance by adding/removing neurons from an initial architecture using a previously specified criterion [50, 51]. The Dynamic Code Creation (DNC) and Cascade Correlation (CC) [52, 53] algorithms are the most well-known methods in this category.
- Computational intelligence algorithms such as Genetic Algorithm (GA) [54-56], fuzzy logic [57], Bayesian training using genetic programming [58], simulated annealing [59], immune algorithm [60], the



PSO algorithm and its variants [61-68], Tabu search [69], fish swarm algorithm [70], harmony search algorithms [71], and the GSA [72-75]. For example, the GA searches in a multi-dimensional space based on its global searching capability. The GA varies the number of hidden layers and hidden neurons through application of genetic operators and evaluation of the different architectures according to a fitness function [76-78].

The GSA is a heuristic algorithm that was introduced by Rashedi *et al.* [79] and is based on the gravitational law and laws of motion. The GSA has a flexible and well-balanced mechanism to enhance exploration and exploitation abilities. A hybrid GSA-BGSA algorithm is used in this paper to optimize the network structure (i.e., the number of hidden layer nodes in a feedforward neural network with single hidden layer using the BGSA and connection weights of this network using the GSA). The initial number of hidden nodes of the mentioned network was considered as 75% of the input features number and it varied by iterations to achieve the Minimum Mean Square Error (MMSE).

III. ANN MODEL EQUIPPED WITH SWITCHES IN HIDDEN LAYER

A multi-layer feedforward neural network has been used in this study. This network is characterized by the number of input nodes, number of hidden layers, number of hidden nodes, transfer function of neurons, and number of output nodes. The number of input and output nodes is equal to the number of input features and number of classes in a pattern recognition problem, respectively. Number of hidden layers is problem dependent. Chester [80] indicated that the appropriate number of hidden layers in most of the problems is one or two. The EBP, as a well known training algorithm of this neural model, is slowing convergent and its speed depends on initial value of the connection weights and the initial learning rate. So, the BGSA is used in this study to determine the optimum number of nodes in the single hidden layer of this model and the GSA is used to obtain the optimum value of weights.

Finding optimal number of hidden layer nodes is a critical task, because if a network is smaller than needed, it may be unable to provide good performance owing to its limited information processing power, and a network larger than needed has redundancy and also loses its generalization.

The structure of the three-layer feedforward ANN that will be optimized in this study is shown in Figure 1. As seen, $X=[x_1, \dots, x_{ni}]$ is the input feature vector to this model. The output vector of this model is shown as $Y=[y_1, \dots, y_{no}]$. So, “ ni ” denotes the number of inputs, “ nh ” denotes the number of hidden nodes, and “ no ” denotes the number of outputs.

$v_{i,j}$ denotes the weight link between the i th input node and the j th hidden node. $w_{j,k}$ denotes the weight link between the j th hidden node and the k th output node. “ S_i ” is the switch value of the i th hidden node ($i=1, \dots, nh$) that is equal to 0 or 1. This switching function is shown by small boxes in Figure 1.

In other words, the links between a hidden node and the input/output nodes will be established, if the switch value of that hidden node is 1. So, a zero value for this switch indicates that corresponding hidden node and its links to the input and output nodes are removed. The output of this model can be determined as follows:

$$y_k(t) = \sum_{j=1}^{nh} w_{j,k} S_j \text{Sigmoid} \left(\sum_{i=1}^{ni} v_{i,j} S_j x_i(t) \right) \quad (3)$$

In the proposed algorithm, we use the GSA for weight updating procedure and the BGSA for obtaining the optimum structure of the neural classifier:

$$y_k^{GSA-BGSA}(t) = g(x(t)) \quad (4)$$

where $Y^{GSA-BGSA}(t)=[y_1^{GSA-BGSA}(t), \dots, y_{no}^{GSA-BGSA}(t)]$ and $X(t)=[x_1(t), \dots, x_{ni}(t)]$ indicate the output and input of the unknown nonlinear function (g), respectively. The GSA-BGSA algorithm is used to minimize the Mean Square Error (MSE) that is considered as the fitness function and defined as follows:

$$MSE(x, y) = \frac{1}{no} \sum_{k=1}^{no} (y_k - y_{target})^2 \quad (5)$$

IV. GSA-BGSA HYBRID FOR OPTIMIZING ANN

A. Review of GSA

Rashedi *et al.* [79] introduced an optimization algorithm based on the law of gravity and mass interactions. In the GSA, a set of agents called masses were introduced to find the optimum solution by simulation of Newtonian laws of gravity and motion. The performance of objects were measured by their masses, and all these objects attracted each other by the gravity force, while this force caused a global movement of all objects towards the objects with heavier masses.

Based on [79], the mass of each agent is calculated after computing the current population fitness, as follows:

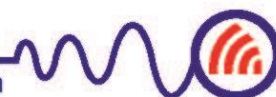
$$M_i(t) = \frac{q_i(t)}{\sum_{j=1}^N q_j(t)}; \quad q_i(t) = \frac{fit_i(t) - worst(t)}{best(t) - worst(t)} \quad (6)$$

where N , $M_i(t)$ and $fit_i(t)$ represent the population size, the mass, and the fitness value of agent i at t , respectively. The $worst(t)$ and $best(t)$ are defined for a minimization problem as follows:

$$best(t) = \min_{j \in \{1, \dots, N\}} fit_j(t) \quad (7)$$

$$worst(t) = \max_{j \in \{1, \dots, N\}} fit_j(t) \quad (8)$$

The acceleration of an agent is computed using (9) in which a_i^d presents the acceleration of agent i in dimension d .



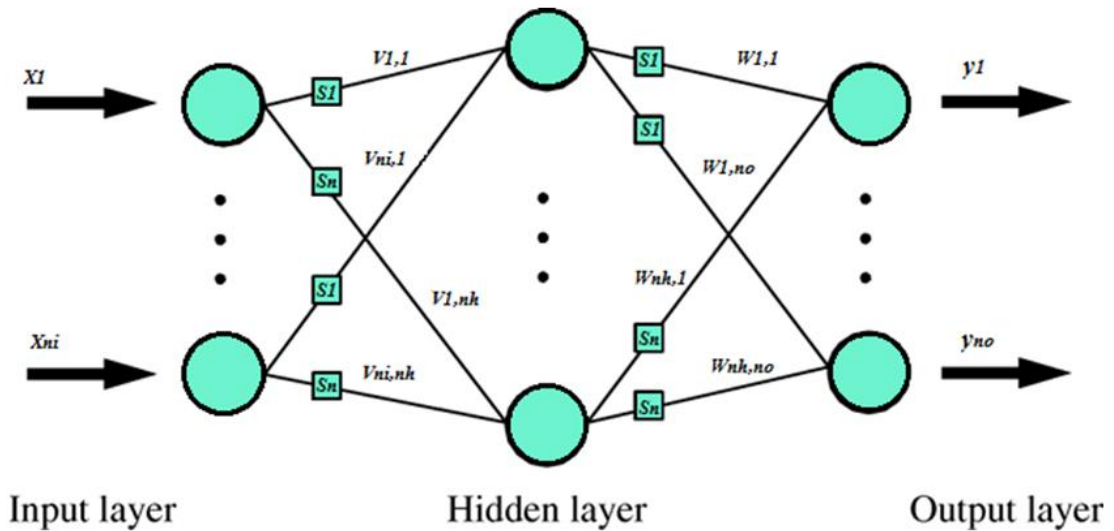


Fig. 1. Structure of a multi-layer feedforward neural network with switched links

$$a_i^d(t) = \sum_{j \in kbest, j \neq i} rand_j G(t) \frac{M_j(t)}{R_{i,j}(t) + e} (x_j^d(t) - x_i^d(t)), \quad (9)$$

$$d = 1, 2, \dots, n, \quad i = 1, 2, \dots, N$$

rand is a uniform random in the interval [0, 1], *e* is a small value, *n* is the dimension of the search space, and $R_{i,j}(t)$ is the Euclidean distance between two agents, *i* and *j*. *kbest* is the set of first *K* agents with the best fitness value and biggest mass, which is a function of time, initialized to K_0 at the beginning and decreased with time. Here K_0 is set to *N* and is decreased linearly to 1. *G(t)* is a decreasing function of time, which is set to G_0 at the beginning and decreases exponentially towards zero with lapse of time. The exponential reduction is given in (10).

$$G(t) = G_0 \exp\left(-\frac{gt}{t_{max}}\right) \quad (10)$$

where t_{max} is the total number of iterations. It is noted that $X_i = (x_i^1, x_i^2, \dots, x_i^N)$ indicates the position of agent *i* in the search space, which is a candidate solution.

The next velocity of an agent is calculated using (11) where v_i^d presents the velocity of agent *i* in dimension *d*:

$$v_i^d(t+1) = rand_i \times v_i^d(t) + a_i^d(t) \quad (11)$$

Then, the position of agent *i* in dimension *d* is calculated as follows:

$$x_i^d(t+1) = x_i^d(t) + v_i^d(t+1) \quad (12)$$

The steps of the GSA algorithm are as follows:

Step 1: Initialization of $X_i(t)$; $i=1,2,\dots,N$;

Step 2: Fitness evaluation of agents;

Step 3: Update of $G(t)$, $best(t)$, $worst(t)$, and $M_i(t)$; $i=1,2,\dots,N$;

Step 4: Calculation of acceleration and velocity;

Step 5: Update of agents' position to obtain $X_i(t+1)$; $i=1,2,\dots,N$;

Step 6: Repeat steps 2 to 5 until the stop criteria is reached.

The BGSA was introduced in [81] to extend the GSA algorithm to tackle binary problems effectively. In the BGSA, the position of agents has two values; 0 or 1, and the velocity of an agent represents the probability that a bit (position) takes on 0 or 1. The velocity updating formula remains unchanged, and the position updating formula is redefined as (13):

$$x_i^d(t+1) = \begin{cases} 1 - x_i^d(t); & rand_i < |\tanh(v_i^d(t+1))| \\ x_i^d(t); & otherwise \end{cases} \quad (13)$$

In the proposed GSA-BGSA algorithm, the agents of GSA and BGSA work together and evaluated simultaneously. Each agent is divided to two sub-agents which have been subjected to two independent and consecutive processes. The first one is a regular GSA, i.e. the traditional velocity and position update of neural network weights. The second one is a BGSA, which allows the agent to determine the number of nodes in single hidden layer of a feedforward neural network.

B. GSA-BGSA as a Tuning Algorithm

In the GSA section of this hybrid algorithm, each agent was encoded as a vector of floating numbers, including all the connected weights of a multi-layer perceptron (MLP) shown in Figure 1. So, the agent in the traditional GSA was encoded as follows:

$$Agent_j = [v_{1,1}, \dots, v_{ni,nh}, b_1^1, \dots, b_{nh}^1, w_{1,1}, \dots, w_{nh,no}, b_1^2, \dots, b_{no}^2]; j = 1, \dots, N \quad (14)$$



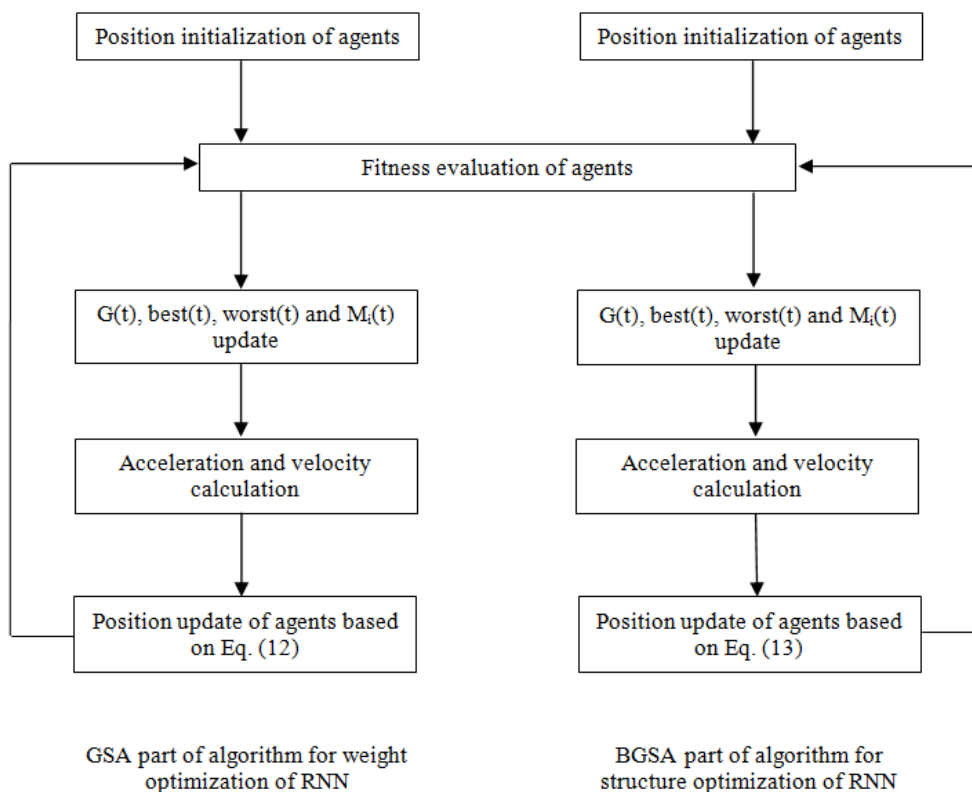


Fig. 2. GSA-BGSA algorithm for optimizing the structure and weights of a neural network

So, the dimension of agent was equal to $[(ni + 1) \times nh] + [(nh + 1) \times no]$ and N is the number of agents. The BGSA was used to update the neural network structure that was considered as switch link and was encoded by:

$$Binary - Agent_j = [\delta_1, \delta_2, \delta_3, \dots, \delta_{nh}]; j = 1, \dots, N \quad (15)$$

The search process of the GSA-BGSA hybrid algorithm for updating the network structure and connection weights is shown in Figure 2. The binary and real GSA algorithms were implemented independently to search the space for finding the best solution. The GSA and the BGSA shared their information by fitness function and mass calculation. The iteration process will repeat for a fixed number of iterations or will end when the search process converges to a pre-defined MSE. The agents' vector in the GSA and corresponding BGSA agents' vector resulting in minimum MSE were used to determine the optimized neural network model.

V. EMOTIONAL SPEECH DATASET

The proper preparation of an emotional speech database requires recording of emotional manifestations. However, real-life emotion data is hard to collect [36, 82]. The text of sentences of FARSDAT speech corpus [83] was used in forming emotional dataset. The FARSDAT is a continuous Farsi neutral speech corpus including 6000 utterances from 300 speakers with various accents. Using 30 non-professional speakers, the emotional speech corpus was recorded in this study. The non-professional speakers were graduate students and speech samples were recorded in a quiet room. The speakers were also directed to keep the degree of

expressiveness of each emotion almost constant. For this purpose, each speaker uttered 252 sentences in four emotional states: neutral, happiness, anger, and surprise.

The speakers were amateur and uttered each sentence several times from the template corpus. The emotional sentences with better quality were selected from the recorded sentences (4964 uttered sentences).

The basic features used for emotion classification were 12 MFCCs, LE, the first three formant frequencies, and F_0 . Each principal feature vector contained 39 components which were 12 MFCCs ($c_1 - c_{12}$), LE, and their derivative ($dc_1 - dc_{12}$ and dLE) and their second derivative ($ddc_1 - ddc_{12}$ and $ddLE$) coefficients of the 13 mentioned features. Also, using three formant frequencies and pitch frequency, 16 supplementary features were calculated. These features contained pitch and formant frequencies (F_0, F_1, F_2, F_3), derivative and logarithm of them (dF_0, dF_1, dF_2, dF_3 and LF_0, LF_1, LF_2, LF_3), and their normalized (zero-mean) values (zF_0, zF_1, zF_2, zF_3) at each frame. To compute zF_i ($i = 0, 1, 2, 3$), the mean value of F_i in each sentence is subtracted from the original value at each 25-ms frame.

The training dataset contained 3475 utterances corresponding to 70% of the corpus and the test dataset included 1489 utterances corresponding to 30% of the corpus.

VI. SIMULATION AND EXPERIMENTAL RESULTS

In this study, the neural emotion classifier was implemented using five methods: the standard EBP algorithm, the PSO algorithm, the GSA, the PSO-

Table 1. Emotion recognition rates of proposed neural classifier with 39-component input feature vector using different optimization methods (averaged over 10 runs)

Method	ANN training algorithm	ANN-structure optimization algorithm	Minimum recognition rate (%)	Maximum recognition rate (%)	Average recognition rate (%)	Number of hidden nodes
EBP	EBP	-	77.16	79.11	77.99	30
PSO	PSO	-	71.70	75.82	72.83	30
PSO-BPSO	PSO	BPSO	68.14	72.78	70.42	17
GSA	GSA	-	79.91	82.20	81.08	30
GSA-BGSA	GSA	BGSA	77.43	82.53	78.84	16

Table 2. Emotion recognition rate of proposed neural classifier with 55-component input feature vector using different optimization methods (averaged over 10 runs)

Method	ANN training algorithm	ANN-structure optimization algorithm	Minimum recognition rate (%)	Maximum recognition rate (%)	Average recognition rate (%)	Number of hidden nodes
EBP	EBP	-	79.71	84.08	81.79	42
PSO	PSO	-	73.27	77.70	75.30	42
PSO-BPSO	PSO	BPSO	68.17	73.35	72.29	31
GSA	GSA	-	82.47	86.43	84.63	42
GSA-BGSA	GSA	BGSA	80.59	85.49	82.60	27

BPSO hybrid algorithm, and the proposed GSA-BGSA hybrid algorithm. The structure of ANN was considered with fixed number of hidden layer nodes when using the EBP, the GSA and the PSO algorithms. This number was equal to about 75% of the input vector dimension [48]. Two feature sets were used in our simulations: a) 39-component feature vector (including c_1-c_{12} , dc_1-dc_{12} , ddc_1-ddc_{12} , LE, dLE, and ddLE), b) 55-component feature vector (including 39 mentioned components and 16 supplementary features introduced in Section 5). So, the number of hidden layer nodes for the first and second feature sets, was set to 30 and 42, respectively. However, in the GSA-BGSA and the PSO-BPSO algorithms, the initial number of hidden nodes was set same as the other methods and was changed with time to optimize the structure of neural network. Because of stochastic search in heuristic algorithms, 10 runs of each algorithm were performed. Initial value of weights was generated at random in the range of [-1, 1]. Other parameter settings of five mentioned methods were performed as follows:

- The EBP algorithm: learning rate was set to 0.001.
- The PSO and the PSO-BPSO algorithms: population size was set to 30 ($N = 30$). The acceleration constants were set to 2, and the inertia factor was decreasing linearly from 0.9 to 0.2 [84, 85].
- The GSA and the GSA-BGSA hybrid algorithm: population size was set to 30 ($N = 30$). $G(t)$ was decreasing linearly from 0.1 to 0.01.

The maximum number of iterations was set to 10000 in all methods. The recognition rate of the proposed neural classifier when employing the first feature set and using each of five methods is reported in Table 1. As seen in this table, the GSA-BGSA hybrid algorithm offers the best maximum recognition

rate as compared to other four algorithms. It is important that this performance is achieved using smaller number of hidden nodes. However, the GSA algorithm offers the best minimum and average recognition rates as compared to other four algorithms but by using more hidden nodes as compared to the hybrid algorithms. Figure 3 shows the result of fitness evaluation or MSE for each of five methods in 10 runs. Results indicate that the GSA-BGSA can achieve the least MSE (in the 4th run), however; the GSA-optimized ANN performs the best averaged over 10 runs.

The recognition rate of the proposed neural classifier when employing the second feature set and using each of five methods is reported in Table 2. As seen in this table, the GSA performs better than the other methods and the GSA-BGSA hybrid method is positioned in the second rank. Figure 4 shows the MSE for each of five methods in 10 runs when using the second feature set containing 55-component feature vectors. As seen, a near competence exists between the GSA and the GSA-BGSA algorithms; however, the performance of GSA-BGSA hybrid algorithm is achieved using smaller number of hidden nodes (27 nodes in the GSA-BGSA method as compared to 42 nodes in the GSA method). The PSO and PSO-BPSO methods trapped in local minima and faced with the main lack of PSO algorithm that is called premature convergence.

By comparing the results in Tables 1 and 2, it is seen that the recognition rates are improved when using supplementary features based on pitch and formant frequencies as compared to the system with 39 input features. The performance of the proposed system is compared with some other emotion recognition systems (Table 3).



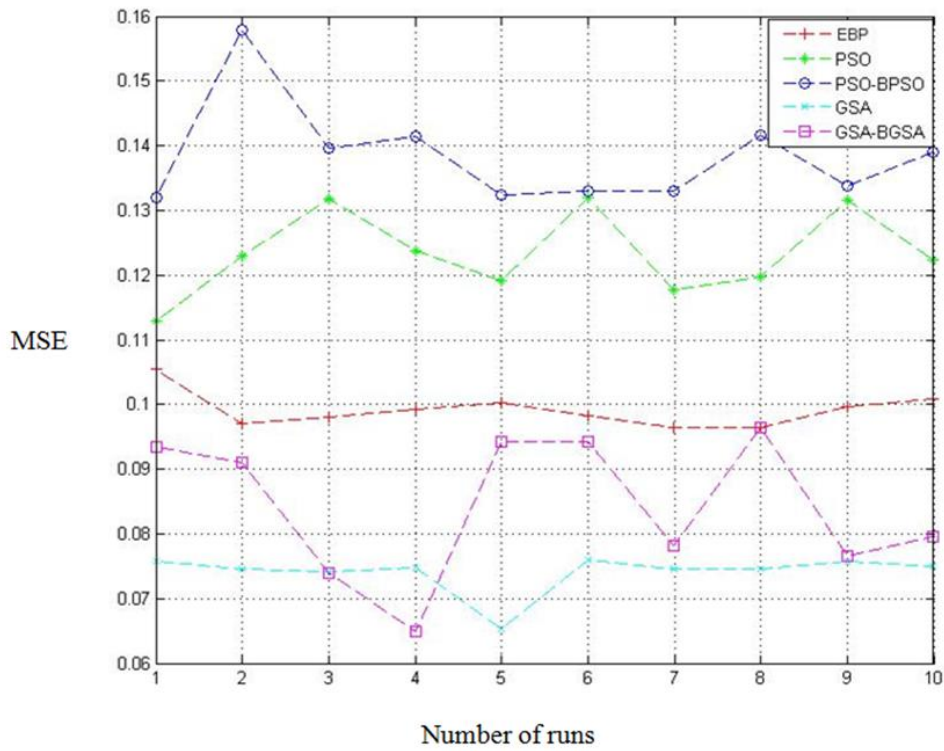


Fig. 3. MSE of the five methods for emotion recognition using 39 input features in 10 runs

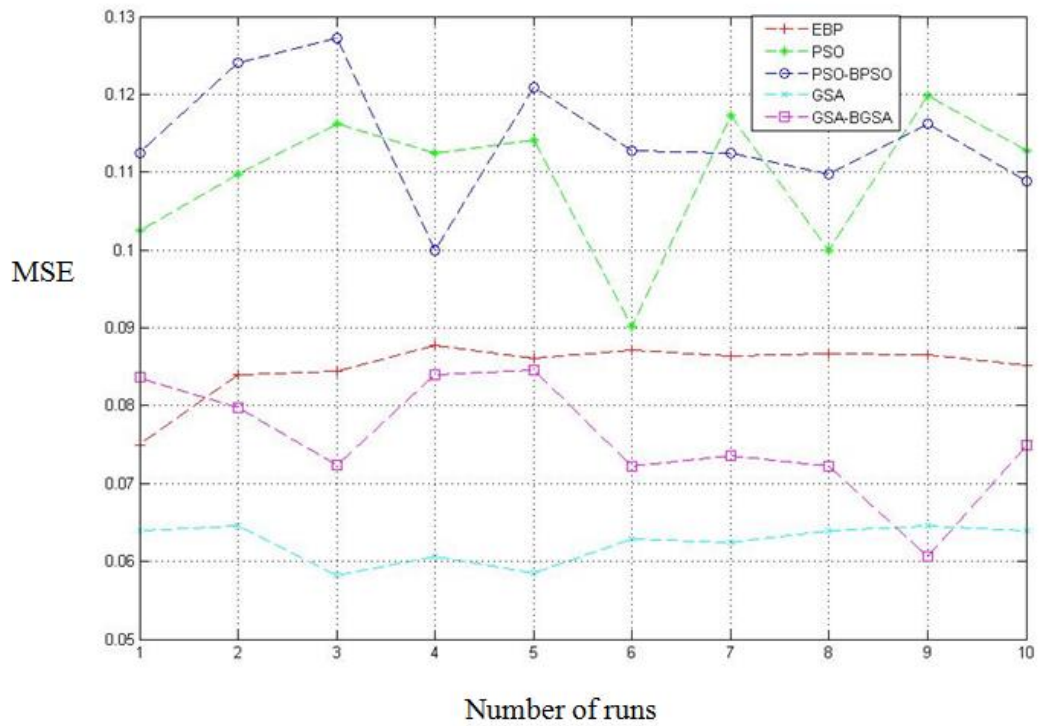


Fig. 4. MSE of the five methods for emotion recognition using 55 input features in 10 runs



Table 3. Performance comparison of proposed emotion recognition system and some similar researches

Emotional states	Input features	Type of classifier(s)	Recognition rate (%)
Happiness, anger, sadness, neutral [86]	F ₀ , dF ₀ , formants, MFCCs	SVM, ANN	71, 42
Happiness, anger, tiredness, sadness, neutral [87]	F ₀ , LE, formants, MFCCs and their first and second derivatives	Gaussian SVM	41
Happiness, anger, anxiety, fear, tiredness, disgust, neutral [88]	MFCCs, E, d _{c_i} , dE, dd _{c_i} , ddE	GMVAR ^a , ANN, HMM	76, 55, 71
Happiness, anger, neutral [16]	55 features introduced in this study	GMM (32 mixtures)	65.9
Happiness, anger, neutral [16]	55 features introduced in this study	C5.0	56.3
Happiness, anger, neutral [16]	55 features introduced in this study	MLP	68.3
Happiness, anger, tiredness, sadness, disgust, fear, neutral [89]	39 features introduced in this study	HMM	81
Happiness, anger, sadness, neutral [24]	F ₀ , sub-band energies, MFCCs, LPCs	Multi-class SVM	80
Happiness, anger, surprise, neutral (proposed model)	55 features introduced in this study	MLP trained by GSA	84.6
Happiness, anger, surprise, neutral (proposed model)	39 features introduced in this study	MLP trained by GSA	81.1
Happiness, anger, surprise, neutral (proposed model)	55 features introduced in this study	BGSA-optimized MLP trained by GSA	82.6
Happiness, anger, surprise, neutral (proposed model)	39 features introduced in this study	BGSA-optimized MLP trained by GSA	78.8

^a Gaussian Mixture Vector Autoregressive Model

Table 4. Number of connection weights in different GSA-BGSA models simulated in this study

Number of inputs to emotion classifier	Number of weights in the optimized-structure MLP	Number of weights in the non-optimized-structure MLP	Reduction rate in the number of weights (%)
55	1593	2478	35.7
39	688	1290	46.7

Because of the different target emotional states and also feature sets in some of these researches, selection of the most effective approach is impossible. However, as can be seen the performance of proposed model is superior to the reported systems. As seen in Table 3, the average recognition rate of proposed neural model when using the BGSA for structure optimization is about 2% lower than the model with non-optimized structure. However, the number of weights in the non-optimized model is increased by 55.6% and 78.5% when using 55 and 39 input features, respectively. So, this recognition result of the GSA-BGSA model was obtained by a considerable lighter neural model. The number of weights for the BGSA-optimized and the non-optimized neural models is reported in Table 4.

Using an Intel quad-core 2.69 GHz CPU and 3GB RAM, the run time of the GSA was 89.9% and 91.5% of the run time of the PSO algorithm when using 39 and 55 input features, respectively. Similarly, the run time of the GSA-BGSA was 86.7% and 83.3% of the run time of the PSO-BPSO algorithm when using 39 and 55 input features, respectively.

VII. CONCLUSION AND FUTURE WORK

In this study, the GSA-BGSA hybrid method was proposed to tune simultaneously the structure and weights of a neural emotion classifier. Two feature sets were employed in the simulations: a) 39-component feature vectors which included 12 MFCCs, and logarithm of energy (i.e., 13 components), the first and second derivatives of these



13 components (i.e., totally 39 components), b) 55-component feature vectors by considering various supplementary features of speech signal based on the first three formants and pitch frequency (i.e., 16 components) and concatenating them to the 39-component feature vector (i.e., totally 55 components). The performance of the proposed hybrid GSA-BGSA-neural model was compared with the hybrid of PSO-BPSO used for such optimizations. In addition, other models such as the GSA-neural hybrid, the PSO-neural hybrid, and the standard EBP algorithm were also included in the performance comparisons. Experimental results showed that the proposed method obtained better results using a lighter network structure as compared with other peer investigated methods.

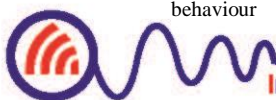
It is noted that the number of input features can also be decreased by employing feature reduction/selection algorithms. So, the proposed method can be modified in future works by inserting a feature selection unit such as ones used in similar works on emotion recognition from speech: least square bound [24], fast correlation-based filter [16], linear discriminate analysis [90], sequential floating forward selection [91, 92], mutual information-based feature selection [24], analysis of variations method [17], combination of the decision tree method and the random forest ensemble [33]. In this way, the number of weights can be decreased more and a very light neural model can be obtained by combining a feature selection method and optimizing strategies proposed in this study. In addition, the single neural classifier in this study can be replaced by a multiple-classifier scheme to improve the performance of system such as ones reported in [17, 93].

REFERENCES

- [1] A. Milton and S. Tamil Selvi, "Class-specific multiple classifiers scheme to recognize emotions from speech signals," *Computer Speech & Language*, vol. 28, pp. 727-742, May 2014.
- [2] A. Mencattini, E. Martinelli, G. Costantini, M. Todisco, B. Basile, M. Bozzali, and C. Di Natale, "Speech emotion recognition using amplitude modulation parameters and a combined feature selection procedure," *Knowledge-Based Systems*, vol. 63, pp. 68-81, June 2014.
- [3] S. Mariooryad and C. Busso, "Compensating for speaker or lexical variabilities in speech for emotion recognition," *Speech Communication*, vol. 57, pp. 1-12, Feb. 2014.
- [4] Y. Sun, G. Wen, and J. Wang, "Weighted spectral features based on local Hu moments for speech emotion recognition," *Biomedical Signal Processing and Control*, vol. 18, pp. 80-90, Apr. 2015.
- [5] H. Cao, R. Verma, and A. Nenkova, "Speaker-sensitive emotion recognition via ranking: Studies on acted and spontaneous speech," *Computer Speech & Language*, vol. 29, pp. 186-202, Jan. 2015.
- [6] H. Ai, D. J. Litman, K. Forbes-Riley, M. Rotaru, J. Tetreault, and A. Purandare, "Using system and user performance features to improve emotion detection in spoken tutoring systems," *Proceedings of Interspeech Conference*, pp. 797-800, 2006.
- [7] L. Devillers, and L. Vidrascu, "Real-life emotions detection with lexical and paralinguistic cues on human-human call center dialogs," *Proceedings of Interspeech Conference*, pp. 801-804, 2006.
- [8] K. B. Khanchandani and M. A. Hussain, "Emotion recognition using multilayer perceptron and generalized feed forward neural network," *Journal of Scientific & Industrial Research*, vol. 68, pp. 367-371, May 2009.
- [9] M. M. Javidi and E. Fazlizadahe Roshan, "Speech emotion recognition by using combinations of C5.0, neural network (NN), and support vector machines (SVM) classification methods," *Journal of Mathematics and Computer Science*, vol. 6, Iss. 3, pp. 191-200, 2013.
- [10] F. Tian, P. Gao, L. Li, W. Zhang, H. Liang, Y. Qian, and R. Zhao, "Recognizing and regulating e-learners' emotions based on interactive Chinese texts in e-learning systems," *Knowledge-Based Systems*, vol. 55, pp. 148-164, Jan. 2014.
- [11] M. Chatterjee, D. J. Zion, M. L. Deroche, B. A. Burianek, C. J. Limb, A. P. Goren, A. M. Kulkarni, and J. A. Christensen, "Voice emotion recognition by cochlear-implanted children and their normally-hearing peers," *Hearing Research*, vol. 322, pp. 151-162, Apr. 2015.
- [12] R. Fernandez and R. Picard, "Recognizing affect from speech prosody using hierarchical graphical models," *Speech Communication*, vol. 53, pp. 1088-1103, Nov. 2011.
- [13] L. Chen, X. Mao, Y. Xue, and L. L. Cheng, "Speech emotion recognition: Features and classification models," *Digital Signal Processing*, vol. 22, pp. 1154-1160, Dec. 2012.
- [14] S. Rigoulot and M. D. Pell, "Emotion in the voice influences the way we scan emotional faces," *Speech Communication*, vol. 65, pp. 36-49, Nov. 2014.
- [15] R. López-Cózar, J. Silovsky, and M. Kroul, "Enhancement of emotion detection in spoken dialogue systems by combining several information sources," *Speech Communication*, vol. 53, pp. 1210-1228, Nov. 2011.
- [16] D. Gharavian, M. Sheikhan, A. Nazerieh, and S. Garoucy, "Speech emotion recognition using FCBF feature selection method and GA-optimized fuzzy ARTMAP neural network," *Neural Computing and Applications*, vol. 21, pp. 2115-2126, Nov. 2012.
- [17] M. Sheikhan, M. Bejani, and D. Gharavian, "Modular neural-SVM scheme for speech emotion recognition using ANOVA feature selection method," *Neural Computing and Applications*, vol. 23, pp. 215-227, Jul. 2013.
- [18] B. Vlasenko, D. Prylipko, R. Böck, and A. Wendemuth, "Modeling phonetic pattern variability in favor of the creation of robust emotion classifiers for real-life applications," *Computer Speech & Language*, vol. 28, pp. 483-500, Mar. 2014.
- [19] Y. Kao and L. Lee, "Feature analysis for emotion recognition from Mandarin speech considering the special characteristics of Chinese language," *Proceedings of International Conference on Spoken Language Processing*, pp. 1814-1817, 2006.
- [20] J. P. Arias, C. Busso, and N. B. Yoma, "Shape-based modeling of the fundamental frequency contour for emotion detection in speech," *Computer Speech & Language*, vol. 28, pp. 278-294, Jan. 2014.
- [21] D. Ververidis and C. Kotropoulos, "Emotional speech recognition: Resources, features, and methods," *Speech Communication*, vol. 48, pp. 1162-1181, Sep. 2006.
- [22] D. Gharavian and M. Sheikhan, "Emotion recognition and emotion spotting improvement using formant-related features," *Majlesi Journal of Electrical Engineering*, vol. 4, pp. 1-8, Dec. 2010.
- [23] T. Pao, Y. Chen, J. Yeh, and Y. Chang, "Emotion recognition and evaluation of Mandarin speech using weighted D-KNN classification," *International Journal of Innovative Computing, Information and Control*, vol. 4, pp. 1695-1709, Jul. 2008.
- [24] H. Altun and G. Polat, "Boosting selection of speech related features to improve performance of multi-class SVMs in emotion detection," *Expert Systems with Applications*, vol. 36, pp. 8197-8203, May 2009.
- [25] R. Gajšek, V. Struc, and F. Mihelič, "Multi-modal emotion recognition using canonical correlations and acoustic features," *Proceedings of International Conference on Pattern Recognition*, pp. 4133-4136, 2010.



- [26] B. Yang and M. Lugger, "Emotion recognition from speech signals using new harmony features," *Signal Processing*, vol. 90, pp. 1415-1423, May 2010.
- [27] D. Bitouk, R. Verma, and A. Nenkova, "Class-level spectral features for emotion recognition," *Speech Communication*, vol. 52, pp. 613-625, Jul. 2010.
- [28] J. Yeh, T. Pao, C. Lin, Y. Tsai, and Y. Chen, "Segment-based emotion recognition from continuous Mandarin Chinese speech," *Computers in Human Behavior*, vol. 27, pp. 1545-1552, Sep. 2010.
- [29] S. Wu, T. H. Falk, and W-Y. Chan, "Automatic speech emotion recognition using modulation spectral features," *Speech Communication*, vol. 53, pp. 768-785, May 2011.
- [30] E. Väyrynen, J. Toivanen, and T. Seppänen, "Classification of emotion in spoken Finnish using vowel-length segments: Increasing reliability with a fusion technique," *Speech Communication*, vol. 53, pp. 269-282, Mar. 2011.
- [31] E. Fersini, E. Messina, and F. Archetti, "Emotional states in judicial courtrooms: An experimental investigation," *Speech Communication*, vol. 54, pp. 11-22, Jan. 2012.
- [32] J. Rong, G. Li, and Y. P. Chen, "Acoustic feature selection for automatic emotion recognition from speech," *Information Processing and Management*, vol. 45, pp. 315-328, May 2009.
- [33] E. Mower, C. Busso, S. Lee, S. Narayanan, and C. Lee, "Emotion recognition using a hierarchical binary decision tree approach," *Speech Communication*, vol. 53, pp. 1162-1171, Nov. 2011.
- [34] M. El Ayadi, M. S. Kamel, and F. Karray, "Survey on speech emotion recognition: Features, classification schemes, and databases," *Pattern Recognition*, vol. 44, pp. 572-587, Mar. 2011.
- [35] A. I. Iliev, M. S. Scordilis, J. P. Papa, and A. X. Falcão, "Spoken emotion recognition through optimum-path forest classification using glottal features," *Computer Speech & Language*, vol. 24, pp. 445-460, Jul. 2010.
- [36] M. Kockmann, L. Burget, and J. H. Černocký, "Application of speaker and language identification state-of-the-art techniques for emotion recognition," *Speech Communication*, vol. 53, pp. 1172-1185, Nov. 2011.
- [37] A. D. Dileep and C. Chandra Sekhar, "Class-specific GMM based intermediate matching kernel for classification of varying length patterns of long duration speech using support vector machines," *Speech Communication*, vol. 57, pp. 126-143, Feb. 2014.
- [38] S. Chandaka, A. Chatterjee, and S. Munshi, "Support vector machines employing cross-correlation for emotional speech recognition," *Measurement*, vol. 42, pp. 611-618, May 2009.
- [39] G. Caridakis, K. Karpouzis, and S. Kollias, "User and context adaptive neural networks for emotion recognition," *Neurocomputing*, vol. 71, pp. 2553-2562, Aug. 2008.
- [40] D. Kukolja, S. Popović, M. Horvat, B. Kovač, and K. Čosić, "Comparative analysis of emotion estimation methods based on physiological measurements for real-time applications," *International Journal of Human-Computer Studies*, vol. 72, pp. 717-727, Oct. 2014.
- [41] N. Ahmed Hendy and H. Farag, "Emotion recognition using neural network: A comparative study," *World Academy of Science, Engineering and Technology*, vol. 75, pp. 791-797, Mar. 2013.
- [42] D. Gharavian, M. Sheikhan, and F. Ashoftedel, "Emotion recognition improvement using normalized formant supplementary features by hybrid of DTW-MLP-GMM model," *Neural Computing and Applications*, vol. 22, pp. 1181-1191, May 2013.
- [43] M. Gori and A. Tesi, "On the problem of local minima in back-propagation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 14, pp. 76-86, Jan. 1992.
- [44] J. Ye, J. Qiao, M. Li, and X. Ruan, "A tabu based neural network learning algorithm," *Neurocomputing*, vol. 70, pp. 875-882, Jan. 2007.
- [45] H. R. Maiera and G. C. Dandya, "Understanding the behaviour and optimising the performance of back-propagation neural networks: An empirical study," *Environmental Modelling & Software*, vol. 13, pp. 179-191, Apr. 1998.
- [46] H. R. Maiera and G. C. Dandya, "The effect of internal parameters and geometry on the performance of back-propagation neural networks: An empirical study," *Environmental Modelling & Software*, vol. 13, pp. 193-209, Apr. 1998.
- [47] P. G. Bernardos and G. C. Vosniakos, "Optimizing feedforward artificial neural network architecture," *Engineering Applications of Artificial Intelligence*, vol. 20, pp. 365-382, Apr. 2007.
- [48] L. M. Salchenberger, E. M. Cinar, and N. A. Lash, "Neural networks: A new tool for predicting thrift failures," *Decision Sciences*, vol. 23, pp. 899-916, Jul. 1992.
- [49] V. Subramanian, M. S. Hung, and M. Y. Hu, "An experimental evaluation of neural networks for classification," *Computers & Operations Research*, vol. 20, pp. 769-782, Sep. 1993.
- [50] M. M. Islam and K. Murase, "A new algorithm to design compact two hidden-layer artificial neural networks," *Neural Networks*, vol. 14, pp. 1265-1278, Nov. 2001.
- [51] X. Jiang and A. H. K. SiewWah, "Constructing and training feedforward neural networks for pattern classification," *Pattern Recognition*, vol. 36, pp. 853-867, Apr. 2003.
- [52] T. Ash, "Dynamic node creation in backpropagation networks," *Connection Science*, vol. 1, pp. 365-375, Dec. 1989.
- [53] S. E. Fahlman and C. Lebiere, "The cascade-correlation learning architecture," In: *Advances in Neural Information Systems 2*, Morgan Kaufmann, Los Altos, CA, pp. 1-13, 1990.
- [54] Q. Zhao and T. Higuchi, "Efficient learning of NN-MLP based on individual evolutionary algorithm," *Neurocomputing*, vol. 13, pp. 201-215, Oct. 1996.
- [55] R. S. Sexton and R. E. Dorsey, "Reliable classification using neural network: A genetic algorithm and back propagation computation," *Decision Support Systems*, vol. 30, pp. 11-22, Dec. 2000.
- [56] M. Castellani and H. Rowlands, "Evolutionary artificial neural network design and training for wood veneer classification," *Engineering Applications of Artificial Intelligence*, vol. 22, pp. 732-741, Jun. 2009.
- [57] P. RoyChowdhury and K. K. Shukla, "Incorporating fuzzy concepts along with dynamic tunneling for fast and robust training of multilayer perceptrons," *Neurocomputing*, vol. 50, pp. 319-340, Jan. 2003.
- [58] T. Marwala, "Bayesian training of neural networks using genetic programming," *Pattern Recognition Letters*, vol. 28, pp. 1452-1458, Sep. 2007.
- [59] S. Amato, B. Apolloni, G. Caporali, U. Madesani, and A. Zanaboni, "Simulated annealing approach in backpropagation," *Neurocomputing*, vol. 3, pp. 207-220, Dec. 1991.
- [60] R. Pasti and L. N. De Castro, "The influence of diversity in an immune-based algorithm to train MLP networks," *Proceedings of International Conference on Artificial Immune Systems*, pp. 71-82, 2007.
- [61] J-R. Zhang, J. Zhang, T-M. Lok, and M. R. Lyu, "A hybrid particle swarm optimization-back-propagation algorithm for feedforward neural network training," *Applied Mathematics and Computation*, vol. 185, pp. 1026-1037, Feb. 2007.
- [62] J. Yu, S. Wang, and L. Xi, "Evolving artificial neural networks using an improved PSO and DPSO," *Neurocomputing*, vol. 71, pp. 1054-1060, Jan. 2008.
- [63] T. Ince, S. Kiranyaz, J. Pulkkinen, and M. Gabbouj, "Evaluation of global and local training techniques over feed-forward neural network architecture spaces for computer-aided medical diagnosis," *Expert Systems with Applications*, vol. 37, pp. 8450-8461, Dec. 2010.
- [64] S. N. Qasem and S. M. Shamsuddin, "Radial basis function network based on time variant multi-objective particle swarm



- optimization for medical diseases diagnosis," *Applied Soft Computing*, vol. 11, pp. 1427-1438, Jan. 2011.
- [65] L. Zhao and F. Qian, "Tuning the structure and parameters of a neural network using cooperative binary-real particle swarm optimization," *Expert Systems with Applications*, vol. 38, pp. 4972-4977, May 2011.
- [66] Z. Pian, S. Li, H. Zhang, and N. Zhang, "The application of the PSO based BP network in short-term load forecasting," *Physics Procedia*, vol. 24, Part A, pp. 626-632, 2012.
- [67] M. A. Cavuslu, C. Karakuzu, and F. Karakaya, "Neural identification of dynamic systems on FPGA with improved PSO learning," *Applied Soft Computing*, vol. 12, pp. 2707-2718, Sep. 2012.
- [68] M. Sheikhan and E. Hemmati, "PSO-optimized Hopfield neural network-based multipath routing for mobile ad-hoc networks," *International Journal of Computational Intelligence Systems*, vol. 5, pp. 568-581, May 2012.
- [69] C. H. Aladag, "A new architecture selection method based on tabu search for artificial neural networks," *Expert Systems with Applications*, vol. 38, pp. 3287-3293, Apr. 2011.
- [70] W. Shen, X. Guo, C. Wu, and D. Wu, "Forecasting stock indices using radial basis function neural networks optimized by artificial swarm algorithm," *Knowledge-Based Systems*, vol. 24, pp. 378-385, Apr. 2011.
- [71] S. Kulluk, L. Ozbakir, and A. Baykasoglu, "Training neural networks with harmony search algorithms for classification problems," *Engineering Applications of Artificial Intelligence*, vol. 25, pp. 11-19, Feb. 2012.
- [72] S. A. Mirjalili, S. Z. Mohd Hashim, and H. Moradian Sardroudi, "Training feedforward neural networks using hybrid particle swarm optimization and gravitational search algorithm," *Applied Mathematics and Computation*, vol. 218, pp. 11125-11137, Jul. 2012.
- [73] M. Sheikhan and M. Sharifi Rad, "Gravitational search algorithm-optimized neural misuse detector with selected features by fuzzy grids based association rules mining," *Neural Computing and Applications*, vol. 23, pp. 2451-2463, Dec. 2013.
- [74] M. Sheikhan and Z. Jadidi, "Flow-based anomaly detection in high-speed links using modified GSA-optimized neural network," *Neural Computing and Applications*, vol. 24, pp. 599-611, Mar. 2014.
- [75] M. Sheikhan and S. Ahmadluei, "An intelligent hybrid optimistic/pessimistic concurrency control algorithm for centralized database systems using modified GSA-optimized ART neural model," *Neural Computing and Applications*, vol. 23, pp. 1815-1829, Nov. 2013.
- [76] G. Bebis, M. Georgiopoulos, and T. Kasparis, "Coupling weight elimination with genetic algorithms to reduce network size and preserve generalization," *Neurocomputing*, vol. 17, pp. 167-194, Nov. 1997.
- [77] C. Jeong, J. H. Min, and M. S. Kim, "A tuning method for the architecture of neural network models incorporating GAM and GA as applied to bankruptcy prediction," *Expert Systems with Applications*, vol. 39, pp. 3650-3658, Feb. 2012.
- [78] F. H. F. Leung, H. K. Lam, S. H. Ling, and P. K. S. Tam, "Tuning of the structure and parameters of a neural network using an improved genetic algorithm," *IEEE Transactions on Neural Networks*, vol. 14, pp. 79-88, Jan. 2003.
- [79] E. Rashedi, H. Nezamabadi-pour, and S. Saryazdi, "GSA: A gravitational search algorithm," *Information Sciences*, vol. 179, pp. 2232-2248, Jun. 2009.
- [80] D. L. Chester, "Why two hidden layers are better than one?," *Proceedings of International Joint Conference on Neural Networks*, pp. 265-268, 1990.
- [81] E. Rashedi, H. Nezamabadi-pour, and S. Saryazdi, "BGSA: Binary gravitational search algorithm," *Natural Computing*, vol. 9, pp. 727-745, Sep. 2010.
- [82] C. Clavel, I. Vasilescu, and L. Devillers, "Fiction support for realistic portrayals of fear-type emotional manifestations," *Computer Speech and Language*, vol. 25, pp. 63-83, Jan. 2011.
- [83] M. Bijankhan, J. Sheikhzadegan, M. R. Roohani, Y. Samareh, C. Lucas, and M. Tebiani, "The speech database of Farsi spoken language," *Proceedings of International Conference on Speech Science and Technology*, pp. 826-831, 1994.
- [84] R. C. Eberhart and J. Kennedy, "Particle swarm optimization," *Proceedings of IEEE International Conference on Neural Networks*, vol. 4, pp. 1942-1948, 1995.
- [85] J. Kennedy and R. C. Eberhart, "A discrete binary version of the particle swarm algorithm," *Proceedings of IEEE International Conference on Systems, Man, and Cybernetics*, vol. 5, pp. 4104-4108, 1997.
- [86] F. Yu, E. Chang, Y. Xu, and H. Shum, "Emotion detection from speech to enrich multimedia content," *Proceedings of IEEE Pacific Rim Conference on Multimedia: Advances in Multimedia Information Processing*, pp. 550-557, 2001.
- [87] O. W. Kwon, K. Chan, J. Hao, and T. W. Lee, "Emotion recognition by speech signal," *Proceedings of European Conference on Speech Communication and Technology*, pp. 125-128, 2003.
- [88] M. Ayadi, S. Kamel, and F. Karray, "Speech emotion recognition using Gaussian mixture vector autoregressive models," *Proceedings of International Conference on Acoustics, Speech, and Signal Processing*, vol. 5, pp. 957-960, 2007.
- [89] B. Vlasenko and A. Wendemuth, "Tuning hidden Markov model for speech emotion recognition," *Proceedings of 33rd German Annual Conference on Acoustics*, pp. 317-320, 2007.
- [90] S. Haq, P. J. B. Jackson, and J. Edge, "Audio-visual feature selection and reduction for emotion classification," *Proceedings of International Conference on Auditory-Visual Speech Processing*, pp. 185-190, 2008.
- [91] D. Ververidis and C. Kotropoulos, "Fast and accurate sequential floating forward feature selection with the Bayes classifier applied to speech emotion recognition," *Signal Processing*, vol. 88, pp. 2956-2970, Dec. 2008.
- [92] A. Batliner, S. Steidl, B. Schuller, D. Seppi, T. Vogt, J. Wagner, L. Devillers, L. Vidrascu, V. Aharonson, L. Kessous, and N. Amir, "Whodunnit-searching for the most important feature types signalling emotion-related user states in speech," *Computer Speech & Language*, vol. 25, pp. 4-28, Jan. 2011.
- [93] E. M. Albornoz, D. H. Milone, and H. L. Rufiner, "Spoken emotion recognition using hierarchical classifiers," *Computer Speech & Language*, vol. 25, pp. 556-570, Jul. 2011.



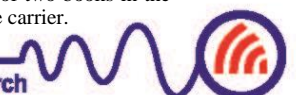
Mansour Sheikhan is currently an Associate Professor in Electrical Engineering Department of Islamic Azad University-South Tehran Branch. His research interests include speech signal processing, neural networks, and intelligent systems. He has published more than 80 journal papers, about 70 conference papers, three books in Farsi, and four book chapters for IET and Taylor & Francis.



Mahdi Abbasnezhad Arabi received his B.Sc. degree in electrical engineering from Noshirvani University, Babol, Iran in 2010 and his M.Sc. degree in electronic engineering from Islamic Azad University-South Tehran Branch in 2014. His research interests include optimization algorithms, neural networks, and intelligent systems.



Davood Gharavian is currently an Assistant Professor in Electrical Engineering Department of Shahid-Beheshti University. His research interests include digital signal processing, speech and image processing, digital signal processors and smart grid. Dr. Gharavian has published more than 20 journal papers and about 20 conference papers. He is the author of two books in the fields of communication systems and power line carrier.



IJICTR

This Page intentionally left blank.

