

AN EVOLUTIONARY ALGORITHM BASED FEATURE EXTRACTION AND SELECTION TO PERSIAN AND ARABIC HANDWRITTEN RECOGNITION

Hooman Kashanian, Samira Arabi Yazdi, Fariba Mahdavi, Masoomeh Esmaelnia
Department of Computer Engineering, Ferdows branch, Islamic Azad University, Ferdows,Iran



This is essentially identical to the following earlier work.

Mohammad Javad Aranian, Monireh Houshmand, Moein Sarvaghad Moghaddam

Feature dimensionality reduction for recognition of Persian handwritten letters using a combination of quantum genetic algorithm and neural network

ICEEE2015 (7th Iranian Conference on Electrical and Electronics Engineering) (2015)

http://www.civilica.com/Paper-ICEEE07-ICEEE07_373.html

<https://www.researchgate.net/publication/314094549>

The first page of the above source is being appended to the end of this document so that the reader can make his/her own judgment.

The material include below was obtained from

<https://doi.org/10.15520/ajcsit.v6i8.50>

<http://innovativejournal.in/ajcsit/index.php/ajcsit/article/view/50>

AN EVOLUTIONARY ALGORITHM BASED FEATURE EXTRACTION AND SELECTION TO PERSIAN AND ARABIC HANDWRITTEN RECOGNITION

Hooman Kashanian, Samira Arabi Yazdi, Fariba Mahdavi, Masoomeh Esmaelnia
Department of Computer Engineering, Ferdows branch, Islamic Azad University, Ferdows,Iran



ARTICLE INFO

Corresponding Author:

Samira Arabi Yazdi
Department of Computer
Engineering, Ferdows branch,
Islamic Azad University,
Ferdows,Iran
s_a.yazdi@yahoo.com



DOI:<http://dx.doi.org/10.15520/ajcsit.v6i8.50>

Keywords: Feature extraction;
pattern recognition; Feature
selection; Persian and Arabic
Handwritten Character recognition;
evolutionary algorithm

ABSTRACT

There are many feature extraction methods for handwritten letters. And selecting an effective subset of features is an important point in analyzing correlation rate in handwritten recognition. Feature selection is needed to select a subset of features that gives good recognition accuracy and has low computational overhead. In this article a methodology for feature selection in unsupervised learning is proposed. The main purpose of this article is enhancing characters recognition and classification, creating quick and low-cost classes, and eventually recognizing Persian and Arabic handwritten characters more accurately and faster. In this paper, to reduce feature dimensionality of datasets a hybrid approach using artificial neural network evolutionary algorithms algorithm is proposed that can be used to distinguish handwritten letters. A key property of our approach is that it does not require any a priori knowledge about the number of features to be used in the feature subset. Implementation results show that evolutionary algorithm are applied here to improve the recognition speed as well as the recognition accuracy.

©2016, AJCSIT, All Right Reserved.

INTRODUCTION

Classification is a process in which machine learns to assign new inputs to pre-defined classes [Kulkarni (1997)]. Different methods are used to classify patterns. Artificial neural networks and learning algorithms such as K-Nearest Neighbors and Support Vector Machine are some of the existing and practical methods that are used for patterns classification [Kulkarni (1998)].

One practical application of algorithms for classification tools is their use in large datasets to recognize one pattern with high number of features to a group of patterns. In this situation to achieve required accuracy and to speed up training and experimenting process, methods for reducing dimensionality of features are needed [Zanganeh (2009)], [Jarmulak (1999)], [Mika (2000)]. A categorization of algorithms used for features' selection and their comparison are also presented in [Ahmad (2005)].

“Optical Character Recognition” (OCR) [Liana (2006)] is a method in which a computer system can recognize texts available in digital images and convert them to text files. Nowadays, OCR is widely used in many contexts. Handwriting recognition, car plate recognition, extracting keywords from image, and indexing image based on content are some of these applications [Liana (2006)]. One of the challenges of OCR is reduction of classification speed and accuracy as a result of imposing lots of features to classifier algorithm. Hence, to achieve high accuracy and to increase training and testing speed, dimensionality reduction of problem is needed [Kheirkhah (2007)].

In [Kheirkhah (2007)], a genetic algorithm-based method for selecting subset of features in Persian OCR has been presented. For dimensionality reduction of problem, Bayesian Classification and genetic algorithm have been used in this

paper. To recognize printed Persian character a combination of genetic algorithm and simulated annealing has been used in [Saber (2005)]. In [Azmi (2010)], Particle Swarm Optimization and genetic algorithm have been used for better recognition of Persian handwritten digits. In [ghanbari (2012)], a hybrid method including neural networks and ant colony optimization has been presented in which neural network has been used as a classifier function used in ant colony optimization (ACO). In this paper, firstly, a hybrid approach for reducing features dimensionality of datasets using neural network and genetic algorithm has been presented to recognize Persian handwritten letters. Then, to enhance proposed method performance, Quantum Genetic Algorithm has been used instead of genetic algorithm. Presented hybrid method has been tested using heady dataset that includes 70000 images of scanned handwritten letters [Pourhabibi (2011)]. Results show that hybrid genetic algorithm-neural network strategy can reduce features dimensionality by 19% and quantum genetic algorithm can reduce number of features by 49%. It also has affected on accuracy and recognition rate of Persian handwritten letters on related dataset by 10%.

This paper is organized as follows: in section 2 how the features are extracted has been expressed. Different methods for selecting subset of features have been discussed in section 3. In section 4, short explanations about quantum genetic algorithm have been presented. The hybrid proposed method have been discussed in section 5 and datasets and how they are used have been discussed in section 6. Then, proposed method has been evaluated in section 7 and finally conclusion has been presented in the last section.

1. Features extraction

Features extraction is an important phase in OCR systems that may affect recognition phase quality. In recognition phase, a code or feature vector is assigned to each character or word input pattern which is that pattern indicator in features space and distinguishes it from other patterns units.

Due to segmentation and recognition phases, there are two main differences among Persian and Latin OCR systems. Because of major differences between Persian and Latin method of writing, it is not possible to apply Latin OCR's segmentation and recognition methods to Persian texts. Complexities available in Persian writing increase commercial OCR systems complexities. That is why most OCR software packages are not able to support Persian and Arabic languages [Sivagaminathan (2007)].

Features extraction techniques can be searched for, in methods related to four general groups of pattern recognition.

- Template matching
- Statistical methods
- Structural methods
- Neural networks

In this paper, 63 features have been extracted for each character using statistical methods. Lots of these features have been extracted based on character segmentation to 9 areas. Feature extraction using segmentation method eliminates context's or character's language limitation, remarkably. Therefore, extracted features can be used for many languages.

2. Different methods for selecting features subset

Algorithms for selecting features' subset are divided into two main categories: Filter method and wrapper method [Zanganeh (2009)].

2.1 Filter method

In this technique, no classification function is used. In other words, no feedback from applied learning algorithm will be used. This is a pre-selected method which is independent from applied machine learning algorithm. Features subset are evaluated using other concepts.

Filter technique works in a way in which, firstly, a weight is calculated for each feature. Then, these weights will be sorted and features with lowest weights are eliminated. A threshold is used for features weights. Then results on features subset are employed to a classifier system as input. Fig. 1 shows how filter technique works.

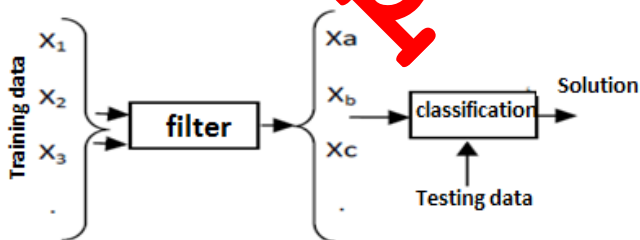


Fig. 1. Filter technique mechanism

2.2 Wrapper method

Wrapper method is known as a black box. In this method a classification function is used for evaluating fitness of feature subsets. This technique uses the feedback that has been applied to learning algorithm. A genetic algorithm has also been used to search for valid features. The main reason for using genetic algorithm is that this algorithm can establish a random search and is not prone to stuck in local minimum [Mika (2000)].

Anyway, crossover operator used in this algorithm works as hybrid solutions, while keeps successful selections of previous feature. In other words, this technique is a feedback method

that uses machine learning algorithm in feature selection procedure. Evaluation is done in each feature selection by execution of inductive algorithm during learning and testing phases. Fig 2 shows how this technique works [Mika (2000)].

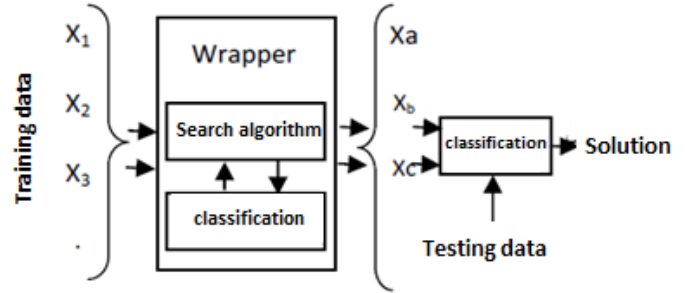


Fig. 2. Wrapper technique mechanism [8]

3. Quantum genetic algorithm

The smallest unit of information which is saved in a quantum computer named as q-bit¹ [Han (13)]. A q-bit can take values 0 or 1. A q-bit state can be represented as Eq. 1.

$$|\psi\rangle = \alpha|0\rangle + \beta|1\rangle \tag{1}$$

While α and β are complex numbers that specify probability amplitude of corresponding states. $|\alpha|^2$ specifies the probability of q-bit to be in state "0" and $|\beta|^2$ specifies the probability of q-bit to be in state "1". Considering the fact that q-bit will take either value "0" or "1", we will have Eq. 2:

$$|\alpha|^2 + |\beta|^2 = 1 \tag{2}$$

Quantum genetic algorithms present replies in a probabilistic format. In Eq. 3, each quantum-chromosome in a space is represented with n q-bits:

$$q = \begin{bmatrix} \alpha_{j1} & \alpha_{j2} & \dots & \alpha_{jn} \\ \beta_{j1} & \beta_{j2} & \dots & \beta_{jn} \end{bmatrix} \tag{3}$$

In which $j = (1,2,\dots,m)$ shows quantum-chromosome number in solution space, m shows total number of chromosomes in solution space, n shows number of qubits available in quantum-chromosome or optimization problem dimension, and t shows generation number of evolutionary algorithm. In fact, $[\alpha^t j i \ \beta^t j i]^T$ shows i -Th qubit from j -Th chromosome in t -Th generation. The main benefit of representing data using qubit is that a chromosome can have different values in solution space. E.g. consider following quantum-chromosome:

$$\begin{bmatrix} 1 & 1 & 1 \\ \sqrt{2} & \sqrt{2} & 2 \\ 1 & -1 & \sqrt{3} \\ \sqrt{2} & \sqrt{2} & 2 \end{bmatrix} \tag{4}$$

The states of this quantum-chromosome can be represented as follows:

$$\frac{1}{4}|000\rangle + \frac{\sqrt{3}}{4}|001\rangle - \frac{1}{4}|010\rangle - \frac{\sqrt{3}}{4}|011\rangle + \frac{1}{4}|100\rangle + \frac{\sqrt{3}}{4}|101\rangle - \frac{1}{4}|110\rangle - \frac{\sqrt{3}}{4}|111\rangle$$

In which the probabilities of states:

$$|111\rangle, |110\rangle, |101\rangle, |100\rangle, |011\rangle, |010\rangle, |001\rangle, |000\rangle$$

Are equal to $\frac{1}{16}, \frac{3}{16}, \frac{1}{16}, \frac{3}{16}, \frac{1}{16}, \frac{3}{16}, \frac{1}{16}, \frac{3}{16}$, respectively.

Obviously, mentioned quantum-chromosome can form eight different states. This means each chromosome in

¹ Quantum bit

quantum genetic algorithm comprises information embedded in many solutions, so only one quantum-chromosome is enough to represent eight different states.

4. Proposed hybrid method

Proposed method of this paper is based on evolutionary theory which says only populations will remain in sequential generations that are fitter than others. Operators such as crossover and mutation will also be applied on populations of a generation for diversity purposes. In other words, we are looking for best answer using population evolution and choosing the bests. In this paper, artificial neural network has been used as classifier function for evaluating generated population by genetic algorithm and quantum genetic algorithm. Fig 3 shows hybrid algorithm mechanism that has been used for obtaining smallest features set with maximum efficiency and desired effectiveness on output.

In Fig 3 dataset S and its subsets include n features for each character that most efficient features are specified by applying these features to hybrid algorithm.

5.1 Appropriate features selection using GA

As mentioned in section 2, 63 features have been extracted from each handwritten character. In this method using genetic algorithm appropriate features for letters classification are selected among 63 extracted features from characters images. For encoding these 63 features in genetic algorithm, a binary chromosome length in 63 is defined. If a bit in that chromosome is 1, that feature will be used in letters classification and if a bit is 0, that feature will not be used in letters classification (Fig 4).

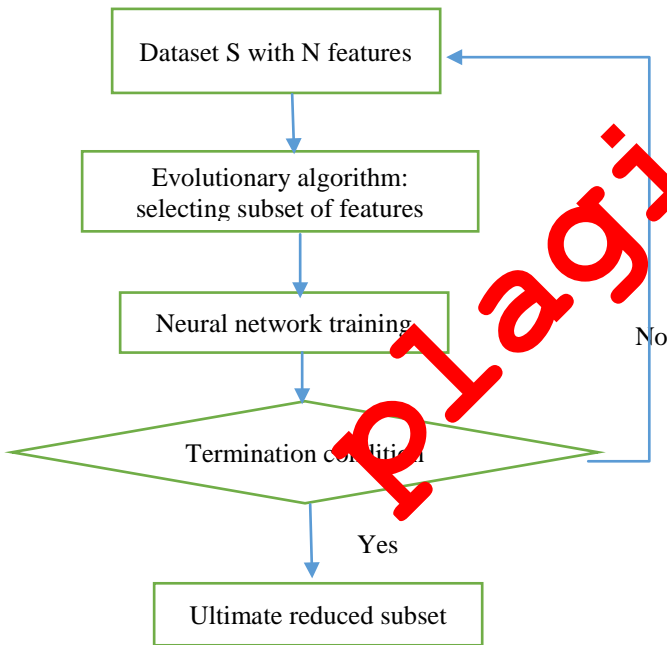


Fig. 3. Hybrid algorithm structure

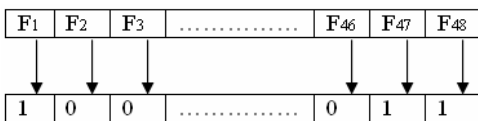


Fig. 4. Presentation of a chromosome and how to select features

Binary chromosomes are generated randomly to form initial population in genetic algorithm. Then, fitness value is computed for each chromosome using fitness function that artificial neural network method has been used for this purpose in this paper. Appropriate parents will also be selected using Roulette wheel technique.

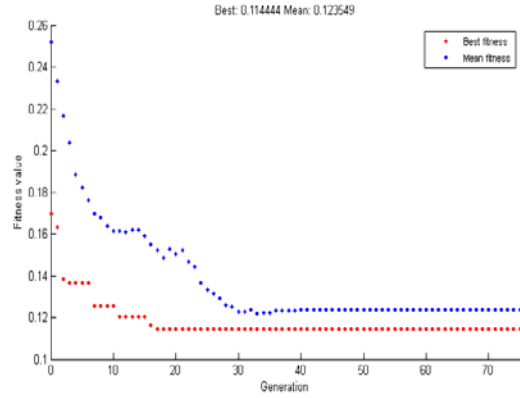


Fig. 5. Genetic algorithm convergence

5.2 Selecting effective features using QGA

Procedure QEA

```

begin
  t=0
  i. initialize Q(0)
  ii. make X(0) by observing Q(0)
  iii. evaluate X(0)
  iv. store the best solutions among X(0) into b
  v. while termination-condition do
    begin
      vi. make X(t) by observing the states of Q(t-1)
      vii. evaluate X(t)
      viii. update Q(t) using Q-gate
      ix. store the best solution among X(t) into b
    end
  end
end
  
```

Fig. 6. Quantum genetic algorithm pseudo code [Kulkarni (1998)]

Quantum genetic algorithm pseudo code is as Fig 6:

^ Above algorithm steps are as follows:

- i. The first step is initializing $Q(t)$.
In this step, initial value $\frac{1}{\sqrt{2}}$ is assigned to all q-bits α_i^0 and $\beta_i^0 (i = 1, 2, \dots, n)$ for all quantum-chromosomes $q_j^0 (j = 1, 2, \dots, m)$. This means the probability of observing “0” and “1” is equal. Here, n is length of chromosome vector and m is number of chromosomes in solution space.
- ii. This step makes a set of binary chromosomes, $P(0)$, from quantum-chromosomes $Q(0)$. In this step, depending on $|\alpha_i|^2$ and $|\beta_i|^2$ values where i is $(i = 1, 2, \dots, n)$, binary chromosomes $p(0) = \{x_1^0, x_2^0, \dots, x_m^0\}$ are made in generation (0). A binary chromosome $x_j^0, (j = 1, 2, \dots, n)$, is a binary solution length in n . Making a binary chromosome based on quantum-chromosome is done by observation. To make a x_i bit based on a q-bit $\begin{bmatrix} \alpha_i \\ \beta_i \end{bmatrix}$, the following equation is used where $U(\cdot)$ is the function that generates uniform random numbers (see Eq. 5).

$$x_i = \begin{cases} 0 & u(0,1) < \alpha_i^2 \\ 1 & \text{Otherwise} \end{cases} \tag{5}$$

- iii. Set of binary chromosomes resulted from second step are evaluated using fitness function.
- iv. The best solution among $P(\mathbf{0})$ is stored into b.
- v. Algorithm will be executed while termination condition is not satisfied.
- vi. In while loop P (t) binary values are made by observing the states of Q (t-1).
- vii. P (t) is evaluated.
- viii. Qubits are updated using quantum-gates. Quantum-gate is an operator that is applied on qubit where $|\alpha|^2 + |\beta|^2 = 1$. Here α' and β' are qubit's updated values. For updating quantum-gate $[\alpha \ \beta]$, following equation is used:

$$[\alpha'_i \ \beta'_i]^T = \begin{cases} U(\Delta\theta_i)[\alpha_i \ \beta_i]^T & \text{if } \alpha\beta > 0 \\ U(-\Delta\theta_i)[\alpha_i \ \beta_i]^T & \text{Otherwise} \end{cases} \quad (6)$$

In which $\Delta\theta$ is extracted from table 1.

In this paper, the following equation has been used for quantum-gate operator:

$$U(\Delta\theta_i) = \begin{bmatrix} \cos(\Delta\theta_i) & -\sin(\Delta\theta_i) \\ \sin(\Delta\theta_i) & \cos(\Delta\theta_i) \end{bmatrix} \quad (7)$$

- ix. The best found solution is stored into b.

C. Calculating features' set fitness using neural network

Table 1. $\Delta\theta$ values for updating q-bit [Laboudi (2012)].

| xi | bi | f(x) > f(b) | $\Delta\theta_i$ |
|----|----|-------------|------------------|
| 0 | 0 | 0 | 0.001π |
| 0 | 0 | 1 | 0.001π |
| 0 | 1 | 0 | 0.08π |
| 0 | 1 | 1 | 0.001π |
| 1 | 0 | 0 | 0.08π |
| 1 | 0 | 1 | 0.001π |
| 1 | 1 | 0 | 0.001π |
| 1 | 1 | 1 | 0.001π |

As mentioned in previous section, in this paper neural network has been used as classifier function for determining fitness value of features' subset. In this phase, features subset has been given to neural network with constant number of neurons in hidden layer for training purpose, which number of neurons in input layer depends on number of features in related subset. Each subset of data has been divided into two sections training and testing on a 70/30 ratio.

After training, neural network will be tested using some new data. Wrong number of classified samples is considered as related subset's error. Error inverse is considered as a criterion for that subset fitness. Hence, each subset of features has an estimation error that helps to determine best subset. Therefore, neural network helps individuals in the *population* of genetic algorithm and quantum genetic algorithm to find best solution fig (7).

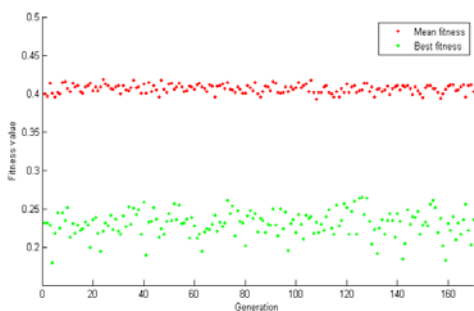


Fig. 7. Quantum genetic algorithm convergence

5. Data Sets

One of the problems with Persian OCR is lack of standard datasets for comparing different applied methods. In this paper the heady dataset [<http://farsiocr.ir>] including 70 thousand and 17 thousand training and testing data from Persian handwritten character images with a resolution of 200 dpi (dot per inch) has been used. At first, all characters' images were resized to a constant size using MATLAB software and then data pre-processing functions like noise removal were applied on them. In the next phase, by combining statistical and structural techniques, 63 features were extracted from all characters. Therefore, a resulted dataset including 63 features was created for 17000 Persian letters to be used in proposed method.

6. Evaluation of proposed method

In this research firstly all extracted features including 63 features were analyzed using SMO² algorithm in RapidMiner 5 data mining software. Achieving an acceptable accuracy in recognition of Persian handwritten letters was the result of this work. Then, firstly using proposed hybrid genetic algorithm-neural network method, 63 selected features were applied to proposed algorithm as input. Proposed hybrid technique could reduce number of features by 19.04%. Once more features resulted from this algorithm tested using SMO technique in RapidMiner software. The result showed recognition precision of about 90.49% for Persian handwritten letters. In the next step, 63 features were applied on quantum genetic algorithm. Quantum genetic algorithm could reduce feature dimensionality by 49.20%. The result of applying these features to SMO algorithm in RapidMiner software was achieving 91.39% accuracy for recognition of handwritten letters. Here, an important issue is use of features that have been extracted by authors of this research to achieve high accuracy in recognition of handwritten character. That is the reason why results did not change considerably after applying proposed algorithm. It is predicated that by applying this algorithm on a dataset with more variety of features and greater number of *inappropriate* features better results can be achieved. It is worth mentioning that one of the objectives of this paper is comparing functionality of genetic algorithm with that of quantum genetic algorithm regarding convergence rate and finding global optimum of problem. Obtained results can be found in table 2. Considering the fact that evolutionary algorithms are heuristic and time-consuming techniques and the fact that proposed algorithm needs chromosome fitness to be specified by neural network in each phase, the algorithm mechanism is time-consuming comparing to other methods. However, the results obtained from features reduction may cause an increase in speed and accuracy of character recognition during test phase.

7. Conclusion and future works

All OCR systems in all existing languages are in need of high accuracy and speed for characters recognition. High dimensionality of problem and great number of features are main reasons for speed and accuracy reduction in these systems. Numerous methods have been proposed for features reduction. In this paper a hybrid approach combining neural network with genetic algorithm and quantum genetic algorithm is used to select an effective subset of features. Since genetic and neural network learning algorithms are heuristic algorithms, runtime prolonged in training phase. However, the purpose is reducing the number of features in order to boost speed and accuracy in recognition of Persian handwritten characters in test phase. Positive results obtained from dimensionality reduction compensate training phase costs in test phase. In this paper, by combining genetic

² sequential minimal optimization

algorithm with neural network classifier, number of features fell from 63 to 51. Additionally, combining quantum genetic algorithm network with neural network decreased number of features by 49%. This reduction in number of features means selecting the features that are more influential in characters classification and leads to the result to be obtained more

accurately and faster. As mentioned in section 7, features used in this paper are extracted by its own authors. Since extreme efforts have been made for appropriate feature to be extracted, selecting a useful subset among good features seems delicate and difficult.

Table 2. Results obtained prior to the execution of algorithm and after it.

| DATA SET | TRAINING DATA | NO. OF TESTED DATA | NO. OF FEATURES | INITIAL ACCURACY WITH SMO ALGORITHM | NO. OF FEATURES WITH GA-ANN ALGORITHM | DIMENSION REDUCTION PERCENTAGE WITH GA-ANN | RECOGNITION ACCURACY WITH GA-ANN | NO. OF FEATURES WITH QGA-ANN ALGORITHM | DIMENSION REDUCTION PERCENTAGE WITH QGA-ANN | RECOGNITION ACCURACY WITH QGA-ANN |
|----------|---------------|--------------------|-----------------|-------------------------------------|---------------------------------------|--|----------------------------------|--|---|-----------------------------------|
| HODA | 70.645 | 17.706 | 63 | 85.11 | 51 | 19.04 | 90.49 | 32 | 49.20 | 91.83 |

As mentioned before, comparing effectiveness of genetic algorithm and quantum genetic algorithm is another purpose of this research. Results show that quantum genetic algorithm is more capable to find optimized answer comparing with genetic algorithm. This point encouraged this paper's authors to further research on solving other optimization problems and features dimensionality reduction using quantum genetic algorithm.

REFERENCES

- ahmad, a; l. dey, (2005): a feature selection technique for classificatory analysis, pattern recognition letters 2(5), pp. 321- 332.
- ahmadreza kheirkhah, esmaeil rahmanian, (2007): optimization of recognition of farsi handwriting characters based on effective feature selection by ga, 8th conference on intelligent systems, Ferdowsi University of Mashhad, (in farsi).
- azmi, reza, boshra pishgoo, narges norozi, marziyeh koohzadi, fahimeh baesi (2010): a hybrid genetic algorithm for feature selection in recognition of hand-printed farsi characters, 12(3), pp. 41-54.
- jarmulak, j; s. craw, (1999): genetic algorithms for feature selection and weighting", appears in proceedings of the ijcai'99 workshop on automatic construction of case based reasoners.
- karthik, rahul sivagaminathan, preeram ramakrishnan, (2007): a hybrid approach for feature subset selection using neural networks and ant colony optimization", expert systems with applications 33 (4), pp. 49-60.
- kuk-hyun han, jong-hwan kim (2002): genetic quantum algorithm and its application to combinatorial optimization problem, iee transaction on pattern analysis and machine intelligence 12(2), pp. 87-95
- kulkarni, r. s.; vidyasagar, m, (1997): learning decision rules for pattern classification under a family of probability measures, iee transactions on information theory, 43(1), pp.154-166.
- kulkarni, r. s.; lugosi, g.; santosh, v. s., (1998): learning pattern classification – a survey, iee transaction on information theory, 44(6),
- kulkarni, r. s.; lugosi, g.; santosh, v. s., (1998): learning pattern classification – a survey, iee transaction on information theory, 44(6), pp. 203-212.
- laboudi, zakaria and salim chikhi, (2012): comparison of genetic algorithm and quantum genetic algorithm, scal group of the misc laboratory, university mentouri, algeria.
- liana, m. lorigo and venu govindaraju, (2006): offline arabic handwriting recognition: a survey" iee transaction on pattern analysis and machine intelligence.
- mika, s; g. ratsch; j. weston; b. scholkopf; a. j. smola; k. r. muller, (2000): invariant feature extraction and classification in kernel spaces, advances in neural information processing systems, massachusetts, usa: mit press, 32(4), pp. 50-
- najme ghanbani seyed mohammad razavi, sedighe ghanbani (2012): optimizing recognition system of persian handwritten digits, majlesi journal of multimedia processing vol. 1, no. 2, june.
- oliveira, l. s, n. benhamed, r. sabourin3, f. bortolozzi, c.v.suen(2001): feature subset selection using genetic algorithms for handwritten digit recognition, proc. xiv brazilian symposium on computer graphics and image processing (sibgrapi'01), 32(8), pp. 362-380.
- oh, i. s. j. s. lee; b. r. moon, (2004): hybrid genetic algorithms for feature selection, iee transaction on pattern analysis and machine intelligence, 26(11),pp. 43-54.
- pourhabibi, tahereh, maryam bahojb imani, saman haratizadeh, (2011): feature selection on persian fonts: a comparative analysis on gaa, gesa and ga, procedia computer science 3(15), pp.14-29.
- saberi, m; d. safaai, (2005): feature selection method using genetic algorithm for the classification of small and high dimension data, iee transaction on pattern analysis and machine intelligence, 23(11), pp. 103-114
- singhi, s. k; h. liu, (2006): feature subset selection bias for classification learning, appearing in proceedings of the 23rd international conference on machine learning, pittsburgh, pa, 31(2),pp. 187-199.
- soryani, m; n.rafat,(2008): application of genetic algorithms to feature subset selection in a farsi ocr", world academy of science, engineering and technology 5(8), pp. 409-417.
- zanganeh, sakineh; reza javanmard; mohamad mahdi ebadzadeh, (2009): a hybrid approach for features dimension reduction of datasets using hybrid algorithm artificial neural network and genetic algorithm-in medical diagnosis, 3rd data mining conference,(in farsi). <http://farsiocr.ir> 2015



Feature dimensionality reduction for recognition of Persian handwritten letters using a combination of quantum genetic algorithm and neural network

Mohammad Javad Aranian

Department of Electrical and Computer Engineering
Imam Reza International University
Mashhad, Iran
Mj.aranian@gmail.com

Monireh Houshmand

Department of Electrical and Computer Engineering
Imam Reza International University
Mashhad, Iran
M_houshmand61@yahoo.com

Moein Sarvaghad Moghaddam

Department of Electrical and Computer Engineering
Semnan University
Semnan, Iran
moeinsarvaghad@yahoo.com

Abstract- Curse of dimensionality is one of the biggest challenges in classification problems. High dimensionality of problem increases classification rate and brings about classification error. Selecting an effective subset of features is an important point in analyzing correlation rate in classification issues. The main purpose of this paper is enhancing characters recognition and classification, creating quick and low-cost classes, and eventually recognizing Persian handwritten characters more accurately and faster. In this paper, to reduce feature dimensionality of datasets a hybrid approach using artificial neural network, genetic algorithm and quantum genetic algorithm is proposed that can be used to distinguish Persian handwritten letters. Implementation results show that proposed algorithms are able to reduce number of features by 19% to 49%. They also show that recognition and classification accuracy of resulted subset of features has risen, by 7/31%, comparing to primitive dataset.

Keywords- *dimensionality reduction of features; recognition of Persian handwritten letters; genetic algorithm; quantum genetic algorithm; neural networks*

I. INTRODUCTION

Classification is a process in which machine learns to assign new inputs to pre-defined classes [1]. Different methods are used to classify patterns. Artificial neural networks and learning algorithms such as K Nearest Neighbors and Support Vector Machine are some of the existing and practical methods that are used in patterns classification [2].

One practical application of algorithms and classification tools is their use in large datasets to designate one pattern with high number of features to a group of patterns. In this situation to achieve required accuracy and to speed up training and experimenting process, methods for reducing dimensionality of features are needed [6, 7, 8]. A categorization of algorithms used for features' selection and their comparison are also presented in [3].

“Optical Character Recognition” (OCR) [4] is a method in which a computer system can recognize texts available in digital images and convert them to text files. Nowadays, OCR is widely used in many contexts. Handwriting recognition, car plate recognition, extracting keywords from image, and indexing image based on content are some of these applications [4].

One of the challenges of OCR is reduction of classification speed and accuracy as a result of imposing lots of features to classifier algorithm. Hence, to achieve high accuracy and to increase training and testing speed, dimensionality reduction of problem is needed [5].

In [5], a genetic algorithm-based method for selecting subset of features in Persian OCR has been presented. For dimensionality reduction of problem, Bayesian Classification and genetic algorithm have been used in this paper. To recognize printed Persian character a combination of genetic algorithm and simulated annealing has been used in [9]. In [10], Particle Swarm Optimization and genetic algorithm have been used for better recognition of Persian handwritten digits. In [11], a hybrid method including neural networks and ant