# The development of Adaptive Neuro-Fuzzy Inference System model to diagnosis diabetes disease data set

**[1]Mamman Mamuda and [2]Saratha Sathasivam**

[1,2]School of Mathematical Sciences
Universiti Sains Malaysia, 11800 Gelugor, Penang, Malaysia
e-mail: [1]maanty123@gmail.com, [2]saratha@usm.my

**Abstract** Medical diagnosis is the extrapolation of the future course and outcome of a disease and a sign of the likelihood of recovery from that disease. Diagnosis is important because it is used to guide the type and intensity of the medication to be administered to patients. A hybrid intelligent system that combines the fuzzy logic qualitative approach and Adaptive Neural Networks (ANNs) with the capabilities of getting a better performance is required. In this paper, a method for modeling the survival of diabetes patient by utilizing the application of the Adaptive Neuro-Fuzzy Inference System (ANFIS) is introduced with the aim of turning data into knowledge that can be understood by people. The ANFIS approach implements the hybrid learning algorithm that combines the gradient descent algorithm and a recursive least square error algorithm to update the antecedent and consequent parameters. The combination of fuzzy inference that will represent knowledge in an interpretable manner and the learning ability of neural network that can adjust the membership functions of the parameters and linguistic rules from data will be considered. The proposed framework can be applied to estimate the risk and survival curve between different diagnostic factors and survival time with the explanation capabilities.

**Keywords** Hybrid intelligent system; neuro-fuzzy; fuzzy inference system; gradient descent algorithm; recursive least square error algorithm.

**2010 Mathematics Subject Classification** 92B05; 92B20

## 1 Introduction

Diabetes is a major health problem in both industrial and developing countries, and its incidence is rising. It is a disease that does not allowed the body to produce or properly use insulin, allowing glucose to enter and fuel them [1]. Diabetes is the leading cause of kidney failure, no traumatic lower limb amputations, and cases of blindness, nerve damage, and blood vessel damage among adults [2] which contributes to heart disease. More than 80 percent of people with diabetes die from some form of heart or blood vessel disease [3].

The cause of diabetes continues to be anonymous, although both genetics and environmental factors such as obesity and lack of exercise appear to play roles [4]. Although detection of diabetes is improving, about half of the patients with type 2 diabetes are undiagnosed and the delay from the disease onset to diagnosis may exceed 10 years [5]. Earlier detection of type 2 diabetes and treatment of related metabolic abnormalities is of vital importance. Type 2 diabetes is commonly found in ancient countries with Pima Indians of Arizona having the highest prevalence and incidence of any population in the world [6].

One of the central problems of the information age is dealing with the enormous amount of raw information that is available. More and more data is being collected and stored in databases or spreadsheets. As this increases, the gap between generating and collecting the data and actually being able to understand it is widening. In order to bridge this knowledge

gap, a variety of techniques known as data mining or knowledge discovery is being developed [7,8]. ANNs have been used as computational tools for pattern classification including diagnosis of diseases because of the belief that they have greater predictive power than signal analysis techniques. However, fuzzy set have attracted the growing attention and interest in modern information technology, production technique, decision making, pattern recognition, diagnostics, data analysis etc.[9]. Neuro-fuzzy systems are fuzzy systems that use ANNs theory in order to determine their properties, i.e. (fuzzy sets and fuzzy rules) by processing data samples. Neuro-fuzzy systems harness the power of the paradigms, i.e. fuzzy logic and ANNs by utilizing the mathematical properties of ANNs in turning the rule-based fuzzy systems that approximate the way human's process information. ANFIS is a specific approach in Neuro-Fuzzy development which has shown significant results in modeling nonlinear functions. The membership function parameters in ANFIS are extracted from a data set that describes the system behavior. The ANFIS learns features in the data set and adjusts the system parameters according to a given error criterion [10]. Implementations of ANFIS in biomedical engineering have been successfully reported for classification and data analysis [11].

In this paper, we have proposed a new approach based on ANFIS to model the survival of diabetes. The rest of the paper is organized as follows: Section 2 discusses the material and method. This has subsections, in each of the subsection, the detailed information was given. The result obtained in the applications was given in Section 3. We conclude the paper in Section 4 by summarizing the result.

## 2   Material and method

### 2.1   Dataset used for diabetes disease

The data set was obtained from Ref. [12]. The data set was selected from a larger data set held by the National Institutes of Diabetes and Digestive and Kidney Disease. All patients in this database are women from Pima-Indian living near phoenix, Arizona, USA. A total of 768 cases were considered from the data. The binary response variable that takes the values '0' or '1' were used, where '0' stands for a negative test for diabetes and '1' stands for a positive test for diabetes. There are 268 cases which represent (34.9%) in the class '1' and 500 cases which represent (65.1%) in the class '0'. There are eight (8) clinical findings from the data: (1) Number of times pregnant. (2) Plasma glucose concentration a 2-h in an oral glucose tolerance test. (3) Diastolic blood pressure (mm Hg). (4) Triceps skin fold thickness (mm). (5) 2-h serum insulin (mu U/ml). (6) Body mass index. (7) Diabetes pedigree function. (8) Age (years).

### 2.2   Adaptive neuro-fuzzy inference system

The architecture of ANFIS is a fuzzy Sugeno model put in the framework of adaptive systems to facilitate learning and adaptation. The framework makes the ANFIS modeling more systematic and less reliant on expert knowledge. The architecture of the ANFIS is shown in Figure 1.
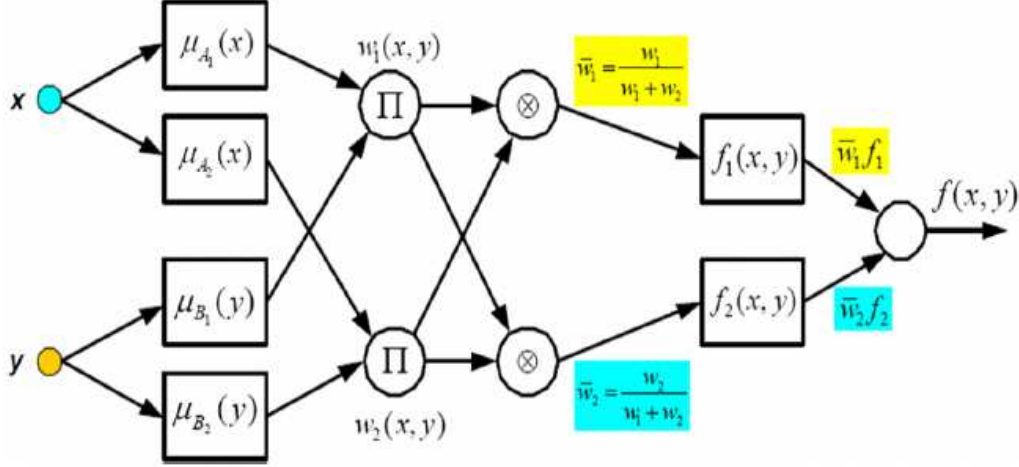
Figure 1: ANFIS architecture

### 2.2.1    Proposed method

The main objective of this paper is to develop an adaptive neuro-fuzzy inference system ANFIS classifier to diagnosis diabetes disease using the diabetes disease data set obtained from Ref [12] that has 8 attributes. We divided the data set into two subsets: the training data set which contains 512 patients' records and the testing data set which contains 256 patients' records. The training data set were used for learning the diabetes pattern.

Presenting the ANFIS architecture, two fuzzy if-then rules that are based on a first order Sugeno model are considered

Rule 1: If $(x$ is $A_1)$ and $(y$ is $B_1)$ then $(f_1 = p_1x + q_1y + r_1)$

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 the inputs $A_i$ and $B_i$ are the fuzzy sets, $f_i$ are the outputs within the fuzzy region specified by the fuzzy rule, $p_i$, $q_i$ and $r_i$ are the design parameters that are determined during the training process. The ANFIS architecture shown in Figure 1 above indicates a circle as a fixed node, whereas a square indicates an adaptive node. The architecture has five layers. All the nodes of the first layer are adaptive nodes. The outputs of the first layer are fuzzy membership grade of the inputs, given by:

$$O_i^1 = \mu A_i(x), i = 1, 2, \tag{1}$$

$$O_i^1 = \mu B_{i-2}(y), i = 3, 4, \tag{2}$$

where $\mu$ an obtained weight according to related fuzzy membership function is, $\mu A_i(x)$, $\mu B_{i-2}(y)$ can adopt any fuzzy membership function, for example, if we employed the bell

shaped membership function; the membership functions are given as:

$$\mu A_i\left(x\right) = \frac{1}{1 + \left\{\left(\left(x - c\right)/a_i\right)^2 \right\}^{b_i}}; \; i = 1, 2, \tag{3}$$

$$\mu B_{i-2}\left(y\right) = \frac{1}{1 + \left\{\left(\left(y - c_i\right)/a_i\right)^2 \right\}^{b_i}}; \; i = 3, 4, \tag{4}$$

where $a_i$, $b_i$, $c_i$ the parameters of the membership function, governing the bell shaped functions accordingly.

The nodes of the second layer are fixed; they are labeled with $\otimes$, indicating that they perform as a simple multiplier. The output of the second layer can be represented as:

$$w = \mu A_1\left(x\right)\mu B_1\left(y\right) + \mu A_2\left(x\right)\mu B_2\left(y\right) \tag{5}$$
$$O_i^2 = w_i = \mu A_i\left(x\right)\mu B_i\left(y\right). \tag{6}$$

These nodes are the firing strengths of the rules.

The third layer nodes are labeled with $\otimes$, indicating that they play a normalization role to the firing strengths from the first and second nodes respectively. The nodes of the third layer are also fixed nodes. The output of this layer can be represented as:

$$O_i^3 = \bar{w}_i = \frac{w_i}{w_1 + w_2}, \quad i = 1, 2 \tag{7}$$

which are also called the normalized firing strengths.

Next is the fourth layer, which is simply as the product of the normalized firing strength. Thus, the outputs of this layer is given by:

$$O_i^4 \, \bar{w}_i f_i = \bar{w}_i\left(p_i x + q_i y + r_i\right), \quad i = 1, 2. \tag{8}$$

The last layer is the fifth layer, this layer consist of only one single fixed node, and this node performs the summation of all the incoming signals. The overall output of the model is given in the fifth layer represented as:

$$O_i^5 = \sum_{i=1}^{2} \bar{w}_i f_i = \frac{\sum_{i=1}^{2} w_i f_i}{w_1 + w_2} \tag{9}$$

In the first layer, there are three modifiable parameters $\{a_i, b_i, c_i\}$, which are related to the input membership functions that are so-called the premise parameters. Also in the fourth layer, there are three modifiable parameters $\{p_i, q_i, r_i\}$, pertaining to the first order polynomial. Learning algorithm for this architecture has the task of tuning all the modifiable parameters of the premise $\{a_i, b_i, c_i\}$ and the consequent $\{p_i, q_i, r_i\}$ to make the out of the ANFIS match the training data. When the premise parameters $a_i$, $b_i$ and $c_i$ of the membership function are fixed, the output of the ANFIS model can be written as:

$$f = \frac{w_1}{w_1 + w_2}f_1 + \frac{w_2}{w_1 + w_2}f_2. \tag{10}$$

But,

$$\bar{w}_i = \frac{w_i}{w_1 + w_2} \tag{11}$$

Substituting (11) into (10) yields

$$f = \bar{w}_1 f_1 + \bar{w}_2 f_2. \tag{12}$$

Substituting the fuzzy if-then rules into (12), it becomes

$$f = \bar{w}_1 (p_1 x + q_1 y + r_1) + \bar{w}_2 (p_2 x + q_2 y + r_2). \tag{13}$$

After rearrangement, the output can be expressed as

$$f = (\bar{w}_1 x) p_1 + (\bar{w}_1 y) q_1 + (\bar{w}_1) r_1 + (\bar{w}_2 x) p_2 + (\bar{w}_2 y) q_2 + (\bar{w}_2) r_2. \tag{14}$$

This is a linear combination of the modifiable consequent parameters $p_1$, $q_1$, $r_1$, $p_2$, $q_2$ and $r_2$. The optimal values of these parameters can easily be identified using the least-square method. When the premise parameters are not fixed, the search space becomes larger and the convergence of the training becomes slower. To overcome this problem, a hybrid algorithm that combines the least-square method and the gradient descent method is adopted. The hybrid algorithm is composed of a forward pass and a backward pass. The forward pass uses the least square method to optimize the consequent parameters with the premise parameters fixed. Once the optimal consequent parameters are found, the backward pass starts immediately. The backward pass uses the gradient descent method to adjust optimally the premise parameters corresponding to the fuzzy sets in the input domain. The output of the ANFIS is calculated by employing the consequent parameters found in the forward pass. The output error is used to adapt the premise parameters by means of a standard back propagation algorithm.

## 3    Results and discussions

The five ANFIS classifiers were trained using the hybrid algorithm of back propagation gradient descent method and the least square method where the 8 attributes representing the diabetes dataset were used as inputs. The ninth ANFIS classifier was trained using the outputs of the five ANFIS classifiers as input data in order to improve classification accuracy. The generalized bell shaped membership function defined in (3) was used as the fuzzy rule architecture of the ANFIS classifiers. The ANFIS classifiers were implemented using the MATLAB software package with fuzzy logic toolbox. The data sets were randomly divided into two: The training data set which consists of 512 patients' records and the testing data set which consists of 256 patients records.

The training data set was used to train the ANFIS model, while the testing data set was used to verify the accuracy and the effectiveness of the trained ANFIS model for classification of the diabetes disease data set. Each ANFIS used 512 training data in 10 training epochs with the error tolerance of 0.01. At the end of 10 training epochs, the network error convergence curve is shown in Figure 2 and Figure 3:
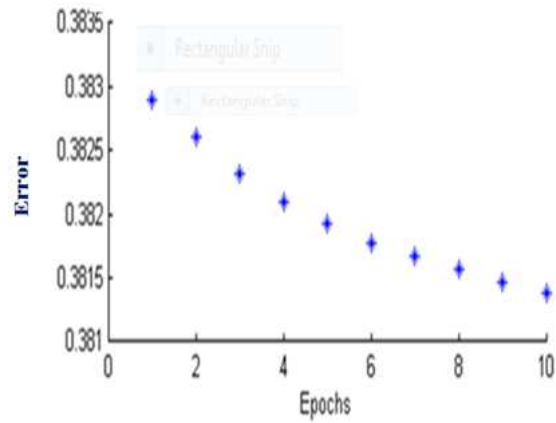
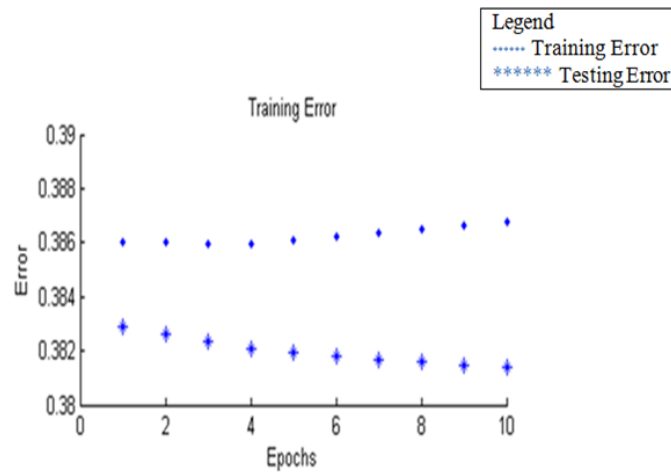Figure 2: The curve of network error convergence of ANFIS training data



Figure 3: The curve of network error convergence of ANFIS training and testing data

### 3.1   Classification accuracy

Classification results of the ANFIS model were displayed using a confusion matrix. In a confusion matrix, information about actual and predicted classification was done by classification system. The performance of the system is commonly evaluated using the data in the matrix. The confusion matrix showing the classification results of the ANFIS model is given in Table 1.

Table 1: Confusion matrix of the ANFIS model

| Output/desired | Result (normal) | Result (patient) |
|---|---|---|
| Result (normal) | 45 (TP) | 8 (FN) |
| Result (patient) | 8 (FP) | 80 (TN) |

The test performance of the classifiers can be determined by the computation of sensitivity, specificity and total accuracy. The sensitivity, specificity and the total accuracy are defined as:

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP + FN}}\,(\%) \tag{15}$$

$$\text{Specificity} = \frac{\text{TN}}{\text{FP + TN}}\,(\%) \tag{16}$$

$$\text{Accuracy} = \frac{\text{TP + TN}}{\text{TP + TN + FP + FN}}\,(\%) \tag{17}$$

where

$$\text{TP} = \text{True Positive},$$
$$\text{FN} = \text{False Negative},$$
$$\text{FP} = \text{False Positive},$$
$$\text{TN} = \text{True Negative}.$$

The results of the statistical parameters i.e. Sensitivity, specificity and accuracy of the ANFIS model were sensitivity (84.91%), specificity (90.1%) and accuracy (88.65%). The obtained classification result of this accuracy using ANFIS model is in line among classifier report from Ref [13]. Table 2 is shown to compare these classifiers with our method.

## 4   Conclusion

In this study, the diagnosis of the diabetes disease using a novel approach known as adaptive-neuro fuzzy inference system (ANFIS) is proposed. The ANFIS classifiers were trained using the hybrid learning algorithm of neural network and were implemented in MATLAB software package with fuzzy logic toolbox. 512 training data were used to train each of the ANFIS in 10 training epochs with tolerance error of 0.01. The research reported in this paper was applied with the task of diagnosing diabetes disease using the most accurate learning methods with a total accuracy of 88.65%. The results strongly agree with other

Table 2: Classification accuracies obtained with our result and other classifiers from literature

| Author (year) | Accuracy (%) | Method |
|---|---|---|
| Logdisc[13] | 77.7 | Statlog |
| BP[13] | 75.2 | Statlog |
| Yildirm *et al.* [13] | 77.08 | BFGS quasi Newton |
| Yildirm *et al.* [13] | 77.60 | Gradient descent |
| Yildirm *et al.* [13] | 80.21 | GRNN |
| Polat and Gunes [3] | 89.47 | PCA-ANFIS |
| Our study [2015] | 88.65 | ANFIS |

result reported from literature. It further suggests that ANFIS can aid in the diagnosis of diabetes disease.

## Acknowledgments

## References

[1] Mohamed, E. I., Linder, R., Perriell, O. G., Di Daniele, N. S., Pppl, J. and De Lorenzo, A. Predicting type 2 diabetes using an electronic nose base artificial neural network analysis. *Diabetes, Nutrition and Metabolism.* 2002. 15(4). 215-221.

[2] Center for disease control and prevention (CDC), National diabetes fact sheet, National estimates and general information on diabetes and prediabeties in the United States, Atlanta, GA: US Department of Health and Human Services, Centers for Disease Control and Prevention, 2011, 201.

[3] Polat, K. and Gnes, S. An expert system approach based on principal component analysis and adaptive neuro fuzzy inference system to diagnosis of diabetes disease. *Digital Signal Processing.* 2007. 17(4): 702-710

[4] Polat, K., Gnes, S. and Arslan, A. A cascade learning system for classification of diabetes disease, generalized discriminate analysis and least square support vector machine. *Expert Systems with Applications.* 2008. 34(1): 482-487.

[5] Temurtas, H., Yumusak, N. and Temurtas, F. A. Comparative study on diabetes disease diagnosis using neural networks. *Expert Systems with Applications.* 2009. 36(4): 8610-8615

[6] Baier, L. J. and Hanson, R. L. Genetics studies of the etiology of type 2 diabetes in Pima Indians. *Diabetes.* .2004. 53(5): 1181-1186.

[7] Knowler, W. C., Pettitt, D. J., Saad , M. F., Charles, M. A.,Nelson, R. G., Howard, B. V. and Bennett, P. H. Obesity in the Pina Indians: Its magnitude and relationship with diabetes. *The American Journal of Clinical Nutrition.* 1991. 53(6): 1543S-1551S.

[8] Michie, D. Methodologies from machine learning in data analysis and software. *The Computer Journal.* 1991. 34(6): 559-565.

[9] Dubois, D. and Prade, H. An introduction to fuzzy systems. *Clinica Chimica Acta.* 1998. 270(1): 3-29.

[10] Jan, J. S. R. ANFIS: Adaptive-network-based fuzzy inference system; Systems. *Man and Cybernetics, IEEE Transactions.* 1993. 23(3): 665-685.

[11] Übeyli, E. D. and Gler, I. Automatic detection of erythemato-squamoles diseases using adaptive neuro-fuzzy inference systems. *Computers in Biology and Medicine.* 2005. 35(5): 421-433.

[12] Blake, C. L. and Merz; C. J. UCI repository of machine learning data bases http://www.ics.uci.edu/MLrepository.html. CA: University of California. Department of Information and Computer Science. 1998. 55.

[13] http://www.phys.uni.torun.pl/kmk/projects/datasets.html. Datasets used for classification comparison of results.