

Modelling the prevalence of diabetes mellitus risk factors based on artificial neural network and multiple regression

Kamal Gholipour,^{1,2} Mohammad Asghari-Jafarabadi,^{3,4} Shabnam Iezadi,⁵ Ali Jannati^{1,2} and Sina Keshavarz⁶

¹Iranian Center of Excellence in Health Management, School of Management and Medical Informatics, Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran. ²Tabriz Health Services Management Research Center, Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran. ³Road Traffic Injury Research Center, Health Management and Safety Promotion Research Institute, Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran. ⁴Department of Statistics and Epidemiology, Faculty of Health, Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran. ⁵Social Determinants of Health Research Center, Health Management and Safety Promotion Research Institute, Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran (Correspondence to: S. Iezadi: sh_iezadi@yahoo.com). ⁶Public Health and Preventive Medicine, University of Social Welfare and Rehabilitation Sciences, Tehran, Islamic Republic of Iran.

Abstract

Background: Type 2 diabetes mellitus (T2DM) is a metabolic disease with complex causes, manifestations, complications and management. Understanding the wide range of risk factors for T2DM can facilitate diagnosis, proper classification and cost-effective management of the disease.

Aims: To compare the power of an artificial neural network (ANN) and logistic regression in identifying T2DM risk factors.

Methods: This descriptive and analytical study was conducted in 2013. The study samples were all residents aged 15–64 years of rural and urban areas in East Azerbaijan, Islamic Republic of Iran, who consented to participate (n = 990). The latest data available were collected from the Noncommunicable Disease Surveillance System of East Azerbaijan Province (2007). Data were analysed using SPSS version 19.

Results: Based on multiple logistic regression, age, family history of T2DM and residence were the most important risk factors for T2DM. Based on ANN, age, body mass index and current smoking were most important. To test for generalization, ANN and logistic regression were evaluated using the area under the receiver operating characteristic curve (AUC). The AUC was 0.726 (SE = 0.025) and 0.717 (SE = 0.026) for logistic regression and ANN, respectively ($P < 0.001$).

Conclusions: The logistic regression model is better than ANN and it is clinically more comprehensible.

Keywords: artificial neural network, diabetes mellitus, multiple regression, risk factors.

<https://doi.org/10.26719/emhj.18.012>

Received: 08/09/15; **accepted:** 12/06/17

Copyright © World Health Organization (WHO) 2018. Some rights reserved. This work is available under the CC BY-NC-SA 3.0 IGO license (<https://creativecommons.org/licenses/by-nc-sa/3.0/igo>).

Introduction

Type 2 diabetes mellitus (T2DM) is a complex metabolic disease with complex causes, manifestations, complications and management (1,2). The chronic complications of T2DM include accelerated development of cardiovascular disease, end-stage renal failure, blindness and lower limb amputations, which can result in excess morbidity and mortality (3). These chronic complications not only have a major impact on patients and their families but also consume an increasing share of health system resources (4). Three hundred and forty-seven million people worldwide suffer from this serious and costly disease. Diabetes now affects both high- and low-income countries but > 80% of people with diabetes live in low- and middle-income countries (5). However, based on World Health Organization (WHO) reports, diabetes mortality will have doubled between 2005 and 2030 and the prevalence of T2DM is increasing worldwide (3,5,6).

Understanding the wide range of risk factors for T2DM can facilitate diagnosis, proper classification and cost-effective management of the disease (6,7). Recent intervention studies have indicated that T2DM can be prevented or delayed by lifestyle changes in high-risk individuals (2,6). Therefore, identifying such risk factors using the right models is the first stage in successful intervention.

Logistic regression is an efficient and powerful tool to assess independent variable contributions to a binary outcome and it is used to analyse the relationship between 1 or more predictors and a dichotomous outcome (8–10). Simultaneous analysis of multiple explanatory variables and reducing the effect of confounding factors are some important advantages of logistic regression (9). However, its accuracy strongly depends on careful variable selection with satisfaction of basic assumptions, as well as appropriate choice of model-building strategy and validation of results (10). Important considerations when conducting logistic regression include adopting independent variables, ensuring that relevant assumptions are met, and selection of the right modelling strategy (10). Logistic regression can be used to study the factors that predict improvement after an intervention (8).

An artificial neural network (ANN) is a nonlinear, computational and complex mathematical model that is constructed to simulate processes of the central nervous system of higher animals, distantly based on the human neuronal structure (11–15). ANNs represents a new method for predictive modelling in medical sciences and they are useful to predict complex, nonlinear and time-dependent relationships. ANNs also can be used when the measures influencing an event are not completely known (14,15). In contrast with traditional statistical techniques, ANNs are capable of automatically resolving these relationships without the need for a priori assumptions about the nature of the interactions between variables. ANNs use data to model and find relationships between factors (11). Another important difference in comparison with traditional statistical methods such as logistic regression is the learning ability of an ANN. A trained network has pooled regulations that are represented by the matrix of the weights between the neurons. This characteristic allows the ANN to forecast cases that have never been presented to the network before and it is called generalization (16).

When predicting and prioritizing risk factors of T2DM, it is questionable which one of these models is better. To respond to this question, we compared the power of these 2 models in terms of sensitivity, specificity and accuracy. For this purpose, we used receiver operating characteristic (ROC) analysis, which included sensitivity, specificity and accuracy of models to indicate the predictive power of models. The ROC curve is a graphical plot that illustrates the performance of a binary classifier system as its discrimination threshold is varied (17). In this study, we investigated the power of logistic regression and ANNs to identify T2DM risk factors and compared them to establish which one was better.

Methods

This descriptive and analytical study was conducted in 2013 to determine the risk factors for T2DM using 2 separate statistical methods. The study sample comprised all residents aged 15–64 years of rural and urban areas of East Azerbaijan Province, Islamic Republic of Iran who were willing to participate. We used a clustered randomized sampling method.

Neighbourhoods and parishes were considered as clusters. In urban settings, a cluster contained 1 or more or parts of a neighbourhood. In rural settings, a cluster contained 1 or more or parts of a village. Cluster heads were selected based on the last digit of the postal code. Each cluster had 20 individuals; 10 males and 10 females living in neighbouring households. From each cluster, we selected 2 men and 2 women from each age group (15–24, 25–34, 35–44, 45–54 and 55–64 years). Every individual in each cluster was selected randomly based on the postal address. We included the nearest right side neighbours to the cluster heads, who were eligible based on age group. Participants gave full informed consent after the study objectives and process were explained. Ethical approval was obtained from the Center for Disease Control of Iran. T2DM was defined as having a diagnosis or receiving a prescription for antidiabetic drugs. New T2DM was considered if fasting plasma glucose (FPG) level was ≥ 126 mg/dl. Impaired fasting glucose was defined by $\text{FPG} \geq 100$ mg/dl (5.6 mmol/l) but < 126 mg/dl (7.0 mmol/l).

At the first stage, during a home visit, health centre staff collected information about sociodemographic, lifestyle and health status through interview. A structured questionnaire (18) was used to explore demographic and ecological characteristics of participants, nutritional status,

diabetes risk factors such as high blood pressure and family history of T2DM, and patients' physical activities, based on WHO guidelines.

Anthropometric measurements were conducted by proficient and skilled healthcare staff of Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran. Body height and weight were measured using a portable electronic weighing scale and portable height-measuring instrument. Participants were asked to remove their shoes and any bulky clothing. Waist circumference was measured at the midpoint between the lower part of the lowest rib, and blood pressure was measured with a calibrated sphygmomanometer. The average of 3 measurements, with a mean time of 5 minutes, was used for analysis. Finally, blood samples (10 ml from every participant) were collected in 4 tubes and centrifuged immediately for measurement of FPG (≥ 126 mg/dl), total cholesterol (TC), high-density lipoprotein cholesterol and triglycerides. A cold chain was preserved while transferring blood samples to the Central Reference Laboratory in Tabriz.

ANN modelling and logistic regression were used to analyse the data. Variables associated with T2DM in the univariate analysis were included in multiple logistic regression models. The P values for entry and removal of variables in the logistic regression model were 0.05 and 0.1, respectively. The significant variables in univariate analyses along with confirmatory factors were used to calculate the individual T2DM risk with the ANN. Uncontrolled hypertension, gender and raised TC > 200 mg/dl were considered as confirmatory factors. The variable importance in the logistic regression analysis was calculated based on standardized coefficient (Wald). The data were divided into a training set (67.1%) and test set (32.9%). Automatic architecture selection was used to determine hidden layers. One hidden layer with 7 units was determined. A scaled conjugate gradient option was used to optimize the algorithm. Modelling was continued until the relative error of testing was less than that of training. Description and diagram network structures were used as the network structure.

We compared the importance of T2DM predictors revealed by logistic regression and ANN. For the ANN and logistic regression models, the area under the ROC curve (AUC) was calculated in the test set. ROC curve is a technique for visualizing, organizing and selecting possibly optimal

models based on their performance. This technique illustrates the performance of a binary classifier by considering sensitivity, specificity and accuracy of models (17). Data were analysed using SPSS version 19 (IBM Corp., Armonk, NY, USA). $P < 0.05$ was considered to be statistically significant.

Results

We used data from 990 participants to identify the T2DM risk factors using 2 separate methods of logistic regression and ANN. Selected risk factors were prioritized and the 2 methods were compared to determine how they differed. Table 1 shows the importance of T2DM predictors according to their priority using the ANN method. Age had the highest score of 0.34, which means that age can predict 34% of T2DM. Raised TC had the lowest score of 0.02. Figure 1 shows the sequence of predictors based on their importance.

In multiple logistic regression, after adjusting for other factors, there was a significant association between T2DM and age [odds ratio (OR): 1.05, 95% confidence interval (CI): 1.03–1.08; $P < 0.001$]. People living in urban compared with rural areas were more likely to develop T2DM (OR: 2.06, 95% CI: 1.26–3.37; $P = 0.004$). According to the association between having T2DM and a positive family history of the disease (OR: 2.56, 95% CI: 1.52–4.31; $P < 0.001$), people with a diabetes patient in their family had greater odds of developing the disease (Table 2). Considering the results of univariate and multiple logistic regression, age, positive family history of T2DM and residence, by adjusting for body mass index (BMI), gender, uncontrolled hypertension, raised TC, controlled hypertension and current smoking were significant and independent risk factors of T2DM (Tables 2 and 3).

We compared the importance of predictors of T2DM based on ANN and logistic regression modelling. For ANN, age, BMI and current smoking were the 3 most important predictors of T2DM, followed by residence, controlled hypertension, uncontrolled hypertension, family history of T2DM, gender and raised TC (> 200 mg/dl). The estimated errors of testing and training were 11% and 15%, respectively, so the goodness of the model was confirmed. For logistic regression modelling, age, family history of T2DM and residence were the most important predictors of T2DM, followed by current smoking, controlled hypertension,

uncontrolled hypertension, BMI, raised TC (> 200 mg/dl) and gender. The latter 5 factors were not significant.

To test the generalization of the results, we evaluated ANN and logistic regression in the test set using AUC values (Table 4). The AUC values were 0.726 (standard error 0.025) and 0.717 (standard error 0.026) for logistic regression and the ANN, respectively. So, the ability of the logistic regression model to predict those with and without T2DM was significantly greater than that of the ANN model ($P < 0.001$).

Discussion

Comparison of the power of an ANN and logistic regression indicated that the latter is a statistically better predictor. In both methods, age was predicted as the most important risk factor in East Azerbaijan Province. Therefore, we suggest paying attention to aged people in the diagnosis and management of T2DM. According to the results of both models together, people who smoke or live in rural areas and those with a family history of T2DM are more at risk of developing T2DM. Also, the risk may increase with BMI.

Logistic regression is easier than ANN to apply and understand. In contrast, ANN can be applied without assumptions used in logistic regression (such as residual normality, homogeneity of residual variances, residual independence and collinearity). Several studies have shown that ANN models have several advantages over conventional statistical methods (19,20). Such models can rapidly recognize linear patterns, categorical and stepwise linear patterns, nonlinear patterns with threshold impacts, and contingency effects. ANN analyses do not need to be started with a hypothesis or preselected key variables. Therefore, undocumented or quantified potential predictors may be specified if they already exist in the various datasets, although they may have been neglected in the past (19,20). Logistic regression as a recognized approach is able to predict clinically relevant dichotomous outcomes. It has some advantages over more traditional approaches to analyse such data (e.g., *t* test and regression), and it is better explored in this context than newer data analysis procedures (e.g., neural nets) (8).

It should be noted that due to the simplicity of interpretation of the variables in the logistic regression model, applying it clinically is more comprehensible. Rahman et al. compared the accuracy of ANN and binary logistic regression models for predicting glucose status (21). They showed a significantly better performance of ANN for detection of impaired glucose tolerance and T2DM patients from disease-free ones (21). Omurlu et al. compared performance of logistic regression and ANN for prediction of albuminuria in T2DM and demonstrated that multilayer perceptron had the highest predictive capability for the presence of albuminuria (22). Zandkarim et al. suggested that logistic regression was more powerful than discriminant analysis for distinguishing T2DM and prediabetes (23). In communities where there is high dependency among case and control groups, recognizing the differences needs stronger methods. Kazemnejad et al. demonstrated that there was no performance difference between models based on logistic regression and ANN in differentiating impaired glucose tolerance and diabetes patients from disease-free patients (24). Zandkarim and Safavi recommended ANNs for medical research in comparison to logistic regression (25).

The present study suggests that statistical analysis of the importance of T2DM risk factors differs using 2 separate models; however, age was the most important predictor in both models. Raised TC and sex had less importance in comparison with other risk factors. Rezaei et al. showed that age, FPG, BMI and mobility variables in their logistic regression model were significant, and FPG, glucose tolerance, BMI and mobility variables indicated the highest predictive power in the neural network model (26). A national survey in 2009 in the Islamic Republic of Iran showed that sex, age and residence were significant predictors of diabetes (27). Logistic regression analysis in a survey in Qatar showed that smoking and family history of DM had a significant association with DM (28).

Our study had 2 major limitations. The first was the small sample size. The second was that a glucose tolerance test was not done and a single FPG test was used. Nevertheless, our study had some strengths, such as inclusion of broad age groups, and it is believed to be the first study to compare 2 models in identifying and prioritizing T2DM risk factors in the Islamic Republic of Iran.

Conclusion

Comparison of the power of ANN and logistic regression models indicated that the latter is better than ANN and is clinically more comprehensible. Logistic regression can provide coefficients such as probability ratio to express the impact of each independent variable on the model and it is better to be used in medicine. However, we should bear in mind that ANNs can easily be used and analysed. It is possible to enter a large number of variables into an ANN and there is no need for assumptions such as normality. Thus, if there is no assumption, we recommend using an ANN model. Our results also showed that age, BMI, family history of T2DM, current smoking and residence are the most important predictors of T2DM in East Azerbaijan Province. A comprehensive programme of diagnosis and management of T2DM, as well as providing consultation for high-risk individuals, which is based on prioritizing people, can be an appropriate initiative to decrease the prevalence of T2DM in East Azerbaijan.

Acknowledgements

We would like to acknowledge the co-operation of health centre staff, in East Azerbaijan Province.

Funding: Student Research Committee, Tabriz University of Medical Sciences, Tabriz, Islamic Republic of Iran.

Competing interests: None declared.

References

1. Reinehr T. Type 2 diabetes mellitus in children and adolescents. *World J Diabetes*. 2013 Dec 15;4(6):270–81. <https://doi.org/10.4239/wjd.v4.i6.270> PMID:24379917
2. Narayan KMV, Benjamin E, Gregg EW, Norris SL, Engelgau MM. Diabetes translation research: where are we and where do we want to be? *Ann Intern Med*. 2004 Jun 1;140(11):958–63. <https://doi.org/10.7326/0003-4819-140-11-200406010-00037> PMID:15172921

3. Kahn R; American Diabetes Association. Type 2 diabetes in children and adolescents. *Diabetes Care*. 2000 Mar;23(3):381–9. <https://doi.org/10.2337/diacare.23.3.381>
PMID:10868870
4. Dyck R, Karunanayake C, Pahwa P, Hagel L, Lawson J, Rennie D, et al.; Saskatchewan Rural Health Study Group. Prevalence, risk factors and co-morbidities of diabetes among adults in rural Saskatchewan: the influence of farm residence and agriculture-related exposures. *BMC Public Health*. 2013 01 5;13(7):7. <https://doi.org/10.1186/1471-2458-13-7>
PMID:23289729
5. Screening for type 2 diabetes. Report of a World Health Organization and International Diabetes Federation meeting. Geneva: Department of Noncommunicable Disease Management, World Health Organization; 2003
(http://www.who.int/diabetes/publications/en/screening_mnc03.pdf, accessed 6 February 2018).
6. Lindström J, Tuomilehto J. The diabetes risk score: a practical tool to predict type 2 diabetes risk. *Diabetes Care*. 2003 Mar;26(3):725–31. <https://doi.org/10.2337/diacare.26.3.725>
PMID:12610029
7. Copeland KC, Becker D, Gottschalk M, Hale D. Type 2 diabetes in children and adolescents: risk factors, diagnosis, and treatment. *Clin Diabetes*. 2005;23(4):181–5.
<https://doi.org/10.2337/diaclin.23.4.181>
8. Reed P, Wu Y. Logistic regression for risk factor modelling in stuttering research. *J Fluency Disord*. 2013 Jun;38(2):88–101. <https://doi.org/10.1016/j.jfludis.2012.09.003>
PMID:23773663
9. Sperandei S. Understanding logistic regression analysis. *Biochem Med (Zagreb)*. 2014 02 15;24(1):12–8. <https://doi.org/10.11613/BM.2014.003> PMID:24627710
10. Stoltzfus JC. Logistic regression: a brief primer. *Acad Emerg Med*. 2011 Oct;18(10):1099–104. <https://doi.org/10.1111/j.1553-2712.2011.01185.x> PMID:21996075
11. Wei JT, Zhang Z, Barnhill SD, Madyastha KR, Zhang H, Oesterling JE. Understanding artificial neural networks and exploring their potential applications for the practicing urologist. *Urology*. 1998;52(2):161–72. PMID:9697777

12. Agatonovic-Kustrin S, Beresford R. Basic concepts of artificial neural network (ANN) modeling and its application in pharmaceutical research. *J Pharm Biomed Anal.* 2000 Jun;22(5):717–27. [https://doi.org/10.1016/S0731-7085\(99\)00272-1](https://doi.org/10.1016/S0731-7085(99)00272-1) PMID:10815714
13. Rodvold DM, McLeod DG, Brandt JM, Snow PB, Murphy GP. Introduction to artificial neural networks for physicians: taking the lid off the black box. *Prostate.* 2001 Jan 1;46(1):39–44. [https://doi.org/10.1002/1097-0045\(200101\)46:1<39::AID-PROS1006>3.0.CO;2-M](https://doi.org/10.1002/1097-0045(200101)46:1<39::AID-PROS1006>3.0.CO;2-M) PMID:11170130
14. Traeger M1 EA, Geldner G, Morin AM, Putzke C, Wulf H, Eberhart LH. Artificial neural networks. Theory and applications in anesthesia, intensive care and emergency medicine. *Anaesthesist.* 2003 52(11):1055–61.
15. Gamito EJ, Crawford ED. Artificial neural networks for predictive modeling in prostate cancer. *Curr Oncol Rep.* 2004 May;6(3):216–21. <https://doi.org/10.1007/s11912-004-0052-z> PMID:15066233
16. Traeger M, Eberhart A, Geldner G, Morin A, Putzke C, Wulf H, et al. Artificial neural networks. Theory and applications in anesthesia, intensive care and emergency medicine. *Anaesthesist.* 2003;52(11):1055-61.
17. Fawcett T. An introduction to ROC analysis. *Pattern Recognit Lett.* 2006;27(8):861–74. <https://doi.org/10.1016/j.patrec.2005.10.010>
18. Bonita R, Winkelmann R, Douglas KA, de Courten M. The WHO stepwise approach to surveillance (steps) of non-communicable disease risk factors. In: McQueen DV, Puska P, editors. *Global behavioral risk factor surveillance.* Boston: Springer; 2003:9–22.
19. Zhu L, Luo W, Su M, Wei H, Wei J, Zhang X, et al. Comparison between artificial neural network and Cox regression model in predicting the survival rate of gastric cancer patients. *Biomed Rep.* 2013 Sep;1(5):757–60. <https://doi.org/10.3892/br.2013.140> PMID:24649024
20. Levine RF. Clinical problems, computational solutions: a vision for a collaborative future. *Cancer.* 2001 Apr 15;91(8 Suppl):1595–602. [https://doi.org/10.1002/1097-0142\(20010415\)91:8+<1595::AID-CNCR1172>3.0.CO;2-P](https://doi.org/10.1002/1097-0142(20010415)91:8+<1595::AID-CNCR1172>3.0.CO;2-P) PMID:11309757

21. Rahman A, Nesha K, Akter M, Uddin SG. Application of artificial neural network and binary logistic regression in detection of diabetes status. *Sci J Public Health*. 2013;1(1):39–43. <https://doi.org/10.11648/j.sjph.20130101.16>
22. Omurlua IK, Tureb M, Unubolc M, Katrancid M, Guney E. Comparing performances of logistic regression, classification & regression trees and artificial neural networks for predicting albuminuria in type 2 diabetes mellitus. *Int J Sci Basic Appl Res*. 2014;16(1):173–87.
23. Zandkarimi E, Safavi AA, Rezaei M, Rajabi G. References comparison logistic regression and discriminant analysis in identifying the determinants of type 2 diabetes among prediabetes of Kermanshah rural areas. *J Kermanshah Univ Med Sci*. 2013;17(5):300–8.
24. Kazemnejad A, Batvandi Z, Faradmal J. Comparison of artificial neural network and binary logistic regression for determination of impaired glucose tolerance/diabetes. *East Mediterr Health J*. 2010 Jun;16(6):615–20. PMID:20799588
25. Zandkarim EI, Safavi AA. Comparison of artificial neural network predictive power with multiple logistic regressions to determine patients with and without diabetic retinopathy. *Razi J Med Sci*. 2014;21(124):79–90.
26. Rezaei M, Zandkarimi e, Hashemian A. Comparison of artificial neural network, logistic regression and discriminant analysis efficiency in determining risk factors of type 2 diabetes. *World Appl Sci J*. 2013;23(11):1522–9.
27. Esteghamati A, Meysamie A, Khalilzadeh O, Rashidi A, Haghazali M, Asgari F, et al. Third National Surveillance of Risk Factors of Non-Communicable Diseases (SuRFNCD-2007) in Iran: methods and results on prevalence of diabetes, hypertension, obesity, central obesity, and dyslipidemia. *BMC Public Health*. 2009 May 29;9:167. <https://doi.org/10.1186/1471-2458-9-167> PMID:19480675
28. Bener A, Zirir M, Janahi IM, Al-Hamaq AO, Musallam M, Wareham NJ. Prevalence of diagnosed and undiagnosed diabetes mellitus and its risk factors in a population-based study of Qatar. *Diabetes Res Clin Pract*. 2009 Apr;84(1):99–106. <https://doi.org/10.1016/j.diabres.2009.02.003> PMID:19261345

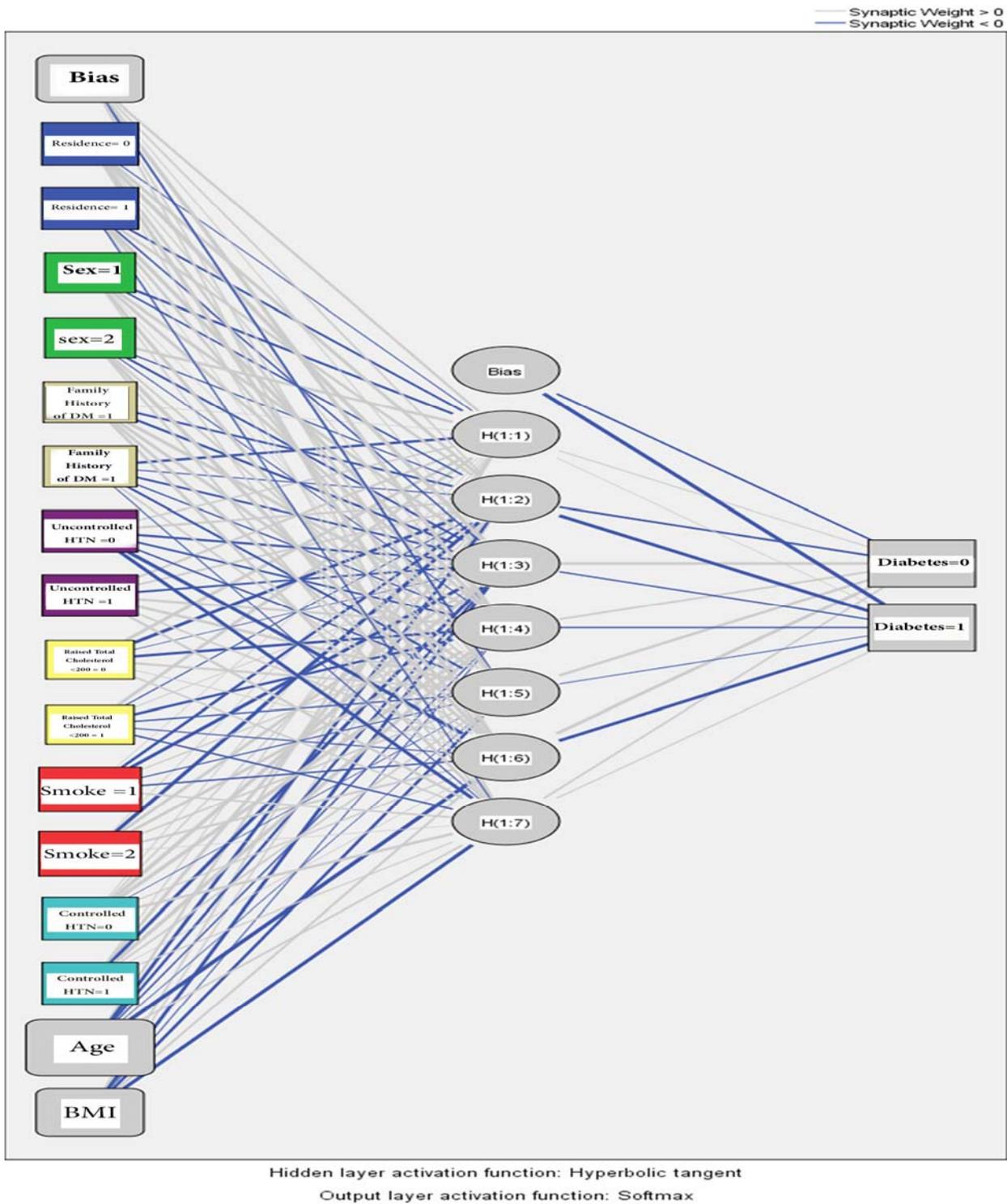


Figure 1 Artificial neural network diagram of relationship between selected risk factors and having had type 2 diabetes mellitus

Table 1 Importance of independent variables (artificial neural network)

Parameter	Importance	Normalized importance (0–100)
Age	0.34	100.0%
BMI	0.17	51.8%
Current smoking	0.13	38.3%
Residence	0.09	26.0%
Controlled HTN	0.08	23.2%
Uncontrolled HTN	0.08	22.6%
Family history of DM	0.06	17.8%
Sex	0.04	11.9%
Raised TC > 200 mg/dl	0.02	5.7%

BMI = body mass index; DM = diabetes mellitus; HTN = hypertension; TC = total cholesterol.

Table 2 Univariate logistic regression analysis for association between selected risk factors and having had DM

Parameter	DM yes (+)		Univariate logistic regression	
	n	N (%)	OR (95% CI)	P value
Age			1.04 1.06)	(1.02– < 0.001
BMI	990	105 (10.6)	1.05 1.10)	(1.01– 0.018
Gender				
Female	49 6	48 (9.67)	1	0.230
Male	49 4	57 (11.54)	0.77 1.18)	(0.51–
Residence				
Urban	59 3	75 (12.65)	2.21 3.47)	(1.40– 0.001
Rural ^a	39 7	30 (7.56)	1	
Uncontrolled				
HTN				
Yes	98	18 (18.4)	1.29 2.25)	(0.74– 0.376
No ^a	58 5	87 (14.9)	1	
Raised TC				
>200 mg/dl				
Yes	20 5	36 (17.56)	1.37 2.13)	(0.87– 0.169
No ^a	47 6	64 (13.44)	1	

Positive family history of DM				
Yes	16	34	2.67	(1.67– < 0.001
	0	(21.25)	4.26)	
No ^a	83	71 (8.55)	1	
	0			
Controlled HTN				
Yes	15	30 (20.0)	1.94	(1.20– 0.006
	0		3.11)	
No ^a	84	75 (8.93)	1	
	0			
Current smoking				
Yes	16	22	1.60	(0.95– 0.076
	0	(13.75)	2.72)	
No ^a	83	83 (10.0)	1	
	0			

^aReference category.

BMI = body mass index; CI = confidence interval; DM = diabetes mellitus; HTN = hypertension; OR = odds ratio; TC = total cholesterol.

Table 3 Multiple logistic regression analysis for association between risk factors and having had DM

Parameter	Multiple logistic regression		
	OR (95% CI)	β coefficient	<i>P</i> value
Age	1.05 (1.03–1.08)	20.66	< 0.001
BMI	1.04 (0.98–1.09)	1.88	0.170
Gender			
Female	0.90 (0.53–1.54)	0.136	0.712
Male	1		
Residence			
Urban	2.06 (1.26–3.37)	7.10	0.004
Rural ^a	1		
Uncontrolled			
HTN			
Yes	0.54 (0.25–1.14)	2.62	0.105
No ^a	1		
Raised TC			
>200 mg/dl			
Yes	1.33 (0.82–2.16)	1.37	0.242
No ^a	1		
Positive family history of DM			
Yes	2.56 (1.52–4.31)	12.48	< 0.001
No ^a	1		
Controlled			
HTN			

Yes	1.87	(0.99– 3.70	0.054
	3.54)		
No ^a	1		
<hr/>			
Current smoking			
Yes	2.05	(1.06– 4.51	0.034
	3.96)		
No ^a	1		

^aReference category.

The Hosmer–Lemeshow goodness-of-fit test: $\chi^2 = 14.17$, degrees of freedom = 8, significance = 0.077.

BMI = body mass index; CI = confidence interval; DM = diabetes mellitus; HTN = hypertension; OR = odds ratio; TC = total cholesterol.

Table 4 Comparison of logistic regression and ANN by area under the ROC curve

Models	Sensitivity	Specificity	Accuracy	Area	SE	Asymptotic significance	Asymptotic 95% CI	
							Lower boundary	Upper boundary
ANN	3.9%	99.5%	83.9%	0.717	0.026	< 0.001	0.666	0.768
Logistic regression	7.1%	99.1%	85.4	0.726	0.025	< 0.001	0.676	0.776

ANN = artificial neural network; CI = confidence interval; ROC = receiver operating characteristic; SE = standard error.