29677

Amir Amiri and Vahid Rafe/ Elixir Comp. Engg. 78 (2015) 29677-29680

Available online at www.elixirpublishers.com (Elixir International Journal)

## **Computer Engineering**



Elixir Comp. Engg. 78 (2015) 29677-29680

# Diagnosing diabetes using data mining algorithms and artificial intelligence

systems

Amir Amiri<sup>1\*</sup> and Vahid Rafe<sup>2</sup>

<sup>1</sup>Department of Computer Engineering, Malayer branch, Islamic Azad University, Malayer, Iran. <sup>2</sup>Department of Computer Engineering, Faculty of Engineering, Arak University, Arak 38156-8-8349, Iran.

#### **ARTICLE INFO**

Article history: Received: 29 November 2014; Received in revised form: 21 December 2014; Accepted: 6 January 2015;

## Keywords

Data mining, Diagnostic, Diabetes, C5.0, Decision trees, Artificial neural networks, Feature selection.

#### ABSTRACT

The purpose of this study is to investigate the role and the scope of the application of data mining predictive science and medicine in a bid to build, evaluate, and the exploitation of data mining models in this regard. In this study to examine the related works in the field of prediction data mining in medicine, who recently published and is trying to highlight important issues and summarize methods and algorithms applied in the form of a series of training. According to the study of used in this study, in most cases to explore knowledge in the medical data from a combination of such as smart algorithms artificial neural network and the decision in the direction of the optimal - the former methods have been used.

© 2015 Elixir All rights reserved.

## Introduction

Today, in medical data are collected in the case of the importance of different diseases. Medical centers with intentions to collect data. A study of the data and to gain the results and useful patterns in connection with the disease is one of the goals of the use of the data. The high volume data and confusion resulting from the problems that impede reaching remarkable results. So of data mining to overcome the problem and to obtain good relations between risk factors, according to the outbreak of the disease and the contribution in mortality humans have be [1].



Figure 1: The process of implementing the proposed system

Data mining and knowledge useful patterns relations hidden in a large volume is given. Such studies and searches can be located along the same duration of the Old Creatures and comprehensive statistics. The major difference in the scale of the extent and variety of fields and applications and data size is contemporary approaches to modeling and machine learning the will seek [2]. In this thesis is trying to provide a solution for the new data mining diagnosis diabetes, according to a similar cases in this case should offer a solution has been forecast accuracy than investigations and secondly could factors affecting the disease diagnosis.

#### The proposed algorithm

The proposed algorithm consists of three stages in the following form is also shown. In the first stage of the database, which is about to be fully in the previous section.

Initially the data from the mining database. Then the data explore, by the process makingof decision - tree (D - T). The best the research and the way of the input characteristics of patients (to - output), the relative avoidance rate ill or healthy (by a symbolic depiction of a tree. Artificial Neural Network (neural networks) to estimate, education, adaptability, machine learning, in order to identify and detect the disease has been used.

Knowledge discovery stages in the database 1959 machine learning year term for the first time by Samuel.

1 - data cleansing: " completely different data

2 - integrating data: a combination of multiple sources, dispersed, and heterogeneous data

3- data choice: retrieve data related to discover the knowledge

4 - data: To apply in different ways

5 - data mining: necessary step in the discovery of knowledge and wisdom of data base in which the statistical methods and machine learning proper model is used to extract.

## Data mining:

Data mining is adoption or exploitation of knowledge of a very large amounts of data, in other words, data mining process.

Tele: E-mail addresses: Amiryamir57@gmail.com With the use of Smart techniques, knowledge of a set of data mining statistical analysis is simply not be able to do it.Data mining very complicated mathematical algorithms to divide data and prediction events is used [3].With the advancement of science and technology and technology tools, the ability to review and storing important data provided with large - scale. Scientific need to search in the data and receive necessary beneficial results.

Data mining, automatic search large data sources, to find patterns and affiliation statistical analysis simply could not do [4].Data mining include information and analysis tools to explore reliable patterns and unknown among a lot of data.Data mining algorithms in various professional methods are used. Data mining, exploitation of knowledge of the data, and learning deductive are called.

Data mining techniques can be in the range of types of data, including style Text, databases, the location information, when information, and other sophisticated data used [5].

Data mining is credible information extraction process, unknown, understandable and reliable of large databases and its use in decision making on major commercial activities. [6]

Data Mining is a process that data mining techniques, intelligent, knowledge of a set of data. [7]

Data Mining i. e., a search in a database to find patterns between the data. [7]

In fact, the discovery of data mining structures, interesting and valuable through a vast collection of data.Data MiningIs an activity that basically with detailed analysis of the data.[6]

- Data mining, extract hidden patterns of relations between the data in a large amount of data and summarize data with innovative ways. [5]

- data mining Recognition is reliable patterns, new, essentially good sense of data.

- Data mining methods in a series of knowledge discovery process can be used to detect patterns and vague relations in the data used.

1 - Data mining operations (categories, clustering, forecast, determining dependence, etc.)

2 - data mining methods (neural networks, the decision - making, the genetic algorithm, etc.)

3-Data - to find a suitable model.

4 - Assessing schemes: best suit pattern.

5-To provide knowledge - knowledge derived using data techniques.

### Data mining techniques:

Various data mining techniques can be based on the types of operations are carried out in two predictors and descriptive divide. Predictive techniques to build a model to the database, predicted task unknown cases. While describing techniques, understandable patterns of data for humans. Descriptive techniques include:

- Summarize
- clustering
- community laws
- consecutive patterns
- Divination techniques include:
- classification
- then extremism
- time series analysis

#### Database:

In the relevant variables in this study as a database of PID medical system has been used as a round - of - the - 8 as the property. The variables used in this research - which includes the following characteristics: [8] of the 867, 500 and 267 people

healthy people infected with the disease, diabetes. 8 registered feature of individuals based on the definition of the World Health Organization (WHO)May - in this case.

The way it works is the first of a series of included in any field is that includes Beat - to the number of features. In this paper each include 8 Bit could be - that every bit represents one of the features. Being a zero Bit shows a lack of property and it is a sign of the property in [9]. The features were randomly selected from among all existing properties. The suffering of the variables and changes in the table below.

100101		bailering sji	iptoms of insolutory
Attribute number	Mean	Standard Deviation	Unit
times_pregnant	3.8	3.4	Number of times pregnant
glucose_tol	120.9	32.0	Plasma glucose concentration a 2 hours in an oral glucose tolerance test
diastolic_pb	69.1	19.4	Diastolic blood pressure (mm Hg)
triceps	20.5	16.0	Triceps skin fold thickness (mm)
insulin	79.8	115.2	2-Hour serum insulin (mu U/ml)
mass_index	32.0	7.9	Body mass index (weight in kg/(height in m)^2)
pedigree	0.5	0.3	Diabetes pedigree function
age	33.2	11.8	Age (years)

Table 1 - Normal suffering symptoms of laboratory

#### Normalization data

Normalization data scale change is that they are to a narrow range and defined as the distance between the 1 to - 1 map. Normalization of various methods that are described below. Normalization causes large - scale data to divert his side. Several ways to normalization exists in this study of normalization z - score is used. In this way, using data from the mean and standard deviation, are normal feature in the relationship is shown.

 $V' = (v - A_v)/\sigma_v \quad or(x - mean(x))/std(x)$ (1)

 $A_{\rm v}$  mean and standard deviation  $\sigma_{\rm v}$  here. Feature technique selection technical feature selection is to reduce the number of features before applying data mining algorithm is used. Characteristics of irrelevant or sub - may have a negative impact on the task to be expected, or complicated calculations. Feature selection consists of three stages:

Selection: and cases of non - important records. To do this, there are several techniques and methods that one of these techniques using genetic algorithms to do that.
Ratings remaining records on the basis of their significance and scored rating.

🔶 diabets **×** <u>≣</u>Eile 🏷 <u>G</u>enerate 8 Rank - 🔺 🛅 🗞 Field 2 glucose\_ 2 mass\_in. 3 pedigree 5 insulin 6 triceps 7 diastolic... Importance Important Important Important Important Important Marginal A Range Range Range Range Range Range Range Range Rank / Туре Value Val 1.0 1.0 1.0 1.0 0.962 0.928  $\mathbf{Y}$ Total fields available:7 Selected fields:5 ★ > 0.95 + <= 0.95 • < 0.9 0 Screened Fields Field  $\nabla$ Туре Reason Model Summary Annotations Apply <u>R</u>eset OK Cancel

**Figure 2. Model Feature Selection** 

- Selection: a record subsets, to maintain the inputs are more important than all the cases, the rest of the filter. To concentrate on the fields and get recordings that more important processing, faster and easier and less number of computation is reduced significantly.

## Model PCA:

PCA to apply the technique is also necessary to reduce the dimensions of the features. This technique must be carefully and that of First of all, the accuracy of the data is not so much diminished when data without this technique are used. When this technique are examined acceptable results and the reduction of the dimensions of the attribute to a good result seems to be necessary. Besides, with regard to reduce by reducing the time and the complexity of the algorithm. From this perspective is important, too. With regard to the aforementioned points only for 5 after new data was considered to be cases mentioned in the figure below.



Figure 3 - impose PCA technique on data and the formation of a new components

In high - tech sector, including the warnings, information than any of the variance, the total variance, solidarity, partial solidarity between attributes between attributes and solidarity be simple. In Table resulting output of PCA technique five major factor resulting from the output - observed.

glucose_tol	triceps	insulin	mass_index	pedigree	age	diabets	\$F-Factor-1	\$F-Factor-2	\$F-Factor-3	\$F-Factor-4	\$F-Factor-5
148.000	35.000	0.000	33.600	0.627	50	yes	0.609	0.993	0.080	0.559	-0.11
85.000	29.000	0.000	26.600	0.351	31	no	-0.682	-0.513	-0.693	0.271	0.72
183.000	0.000	0.000	23.300	0.672	32	yes	-0.445	1.048	1.713	0.026	-0.64
89.000	23.000	94.000	28.100	0.167	21	no	-0.633	-1.002	-0.638	-0.680	0.54
137.000	35.000	168.0	43.100	2.288	33	yes	1.844	-1.312	2.949	4.293	-1.30
116.000	0.000	0.000	25.600	0.201	30	no	-1.067	0.541	-0.231	-0.409	0.08
78.000	32.000	88.000	31.000	0.248	26	yes	-0.546	-1.312	-0.490	-0.360	0.24
115.000	0.000	0.000	35.300	0.134	29	no	-1.627	-0.695	1.283	-0.980	-2.71
197.000	45.000	543.0	30.500	0.158	53	yes	2.311	0.560	1.820	-3.268	1.90
125.000	0.000	0.000	0.000	0.232	54	ves	-1.573	2.348	0.607	-0.298	3.76
110.000	0.000	0.000	37.600	0.191	30	ne	-0.416	0.714	-1.373	-0.142	-0.71
168.000	0.000	0.000	38.000	0.537	34	ves	0.085	1.136	0.394	0.050	-1.77
139.000	0.000	0.000	27.100	1.441	57	ne	0.183	1.962	1.609	2.980	0.47
189.000	23.000	846.0	30.100	0.398	59	ves	2.682	0.509	3.541	-3.599	2.75
166.000	19.000	175.0	25.800	0.587	51	yes	0.578	1.300	1.375	-0.424	0.92
100.000	0.000	0.000	30.000	0.484	32	ves	-1.735	-0.759	1.783	0.119	-1.85
118.000	47.000	230.0	45.800	0.551	31	yes	1.709	-0.837	-0.919	-0.140	-0.07
107.000	0.000	0.000	29.600	0.254	31	yes	-0.939	0.473	-0.465	-0.117	-0.21
103.000	38.000	83.000	43.300	0.183	33	ne	0.034	-1.152	-0.286	-0.770	-1.61
115.000	30.000	96.000	34,600	0.529	32	yes	0.333	-0.371	-0.242	0.166	0.02
126.000	41.000	235.0	39.300	0.704	27	ne	1.511	-0.802	-0.347	0.105	0.45
99.000	0.000	0.000	35.400	0.388	50	no	-0.431	1.500	-0.845	0.629	-0.08
196.000	0.000	0.000	39.800	0.451	41	yes	0.579	2.055	0.185	-0.321	-1.79
119.000	35.000	0.000	29.000	0.263	29	yes	-0.083	-0.175	-0.923	-0.274	0.43
143.000	33.000	146.0	36.600	0.254	51	yes	1.107	1.207	-0.762	-0.739	0.66
125.000	26.000	115.0	31.100	0.205	41	ves	0.116	0.363	-0.207	-0.837	0.44
147.000	0.000	0.000	39.400	0.257	43	yes	-0.105	1.505	-0.337	-0.358	-1.45
97.000	15.000	140.0	23.200	0.487	22	no	-0.607	-0.844	0.364	-0.095	1.06
145,000	19.000	110.0	22.200	0.245	57	no	0.044	1,789	0.361	-0.785	1.64

Figure 4 - new components of the execution of the algorithm PCA

## Decision tree algorithm:

According to the algorithm, it could be two or more branches. With the survey a tree at the root of the decision to a lower level than a little or no. Each node of the data to decide which split. Decision trees through consecutive separation of data to separate groups are made. The goal in the process of increasing the distance between groups in each. One of the differences between the tree construction methods, the decision is how to measure the distance. Decision trees for predicting a bunch of variables used classification trees. Because these samples in the categories. Every path in the tree usually decided to a leaf is understandable. Of the terms of a tree decision could explain its projections, which is an important advantage.

Analysis	of [diabets]	#6			
<u>F</u> ile	🗐 <u>E</u> dit 🛛 🛛	I 🖻 🗞			×
8 Colla	apse All	🌳 Expand All			
Results- Con⊡⊡Con	for output fin nparing \$C-0 'Partition'	eld diabets diabets with diabets 1 Training		2 Testing	
	Correct	496	81.44%	113	71.07%
	Wrong	113	18.56%	46	28.93%
	Total	609		159	

#### Figure 5. Decision tree [10]

Decision trees frequency of some of the data for each level crossing and tree with variables large anticipated work well. As a result, models quickly, that they are made for the most suitable data sets. If allow unfettered grow spent more time to build a non - intelligent, but it's more important is that fit the data over. **Model C5.0:** 

C5.0 is the decision to a tree, or a series of laws. The model the Field, the most important information, and classification. Each subsample by the first Turk, or by using different fields, the main categories. This procedure is repeated until other subassemblies fails to other smaller subset or be divided. Finally, the lowest side (where isleaf) test again view was, in fact, that leaves a very important not detected and greed.

The concept of neural networks, related to simulate the learning in humans and the implementation of the computer algorithms. The learning in humans is learning about patterns by iterations. The brain as a system of information processing parallel with the structure of the millions of neuron is formed. That are called nerve tissue of social are neurons that information and messages from one part of the other part of the body. Learning System of neural networks complex, which consists of a series of brain neurons inspired. Although there is a neuron alone is a simple structure, a network of neuron that are connected. Can learn a complex tasks. Artificial neural networks provide the junior level of a nonlinear training are naturally neuron there is real.

Of artificial neural networks would offer to use circuitry hardware and software (algorithms) to build intelligent machines, capabilities of the robots, programs, and so on. This method are able to learn during the process. Neuron biological of four parts:

1 - Dendrite: Information and cells.

2 - Cell body received information.

3-Exxon Information, the cells to another neuron. 4- The confluence of a synapse Exxon from a laptop to dendrite of another neuron say synapse.

In Figure under artificial neuron similarity with biological neuron is shown. Neuron false information processing the smallest units. Neuron inlets by a form of communication in the name of the weight to enter neuron[11].

First, a combined function (typically gathered  $\Sigma$ ) is a linear combination of the nodes (x), along with the weight of the imposed related to each node (W) are together and become a scalar value.

$$net_{j} = \sum_{i} W_{ij} x_{ij} = W_{0j} x_{0j} + W_{1j} x_{1j+\ldots} + W_{Ij} x_{Ij}$$
(2)

#### Amir Amiri and Vahid Rafe/ Elixir Comp. Engg. 78 (2015) 29677-29680

data mining technique	Prediction accuracy of diabetes
SVM	77/34
SSVM	76/73
Navies Bayesian	76/30
AD Tree	72/91
Decision Table	71/22
Kstar	69/14
ANN	73/40
FLANN	71/84
MLP	76/89
ANFIS	80.60%
KNN	75.55%
The proposed method	88.02%

Table 2.	The perf	formance o	of the algo	orithm	proposed	in com	parison t	<u>o other</u>	algorithms
		1.4.		n	1. 4.		. C. 1 L		



Figure 6. Similarity with artificial neuron biological neurons [12]

To the value of the transfer function is used as input. Neuron in biological, when the combination of neuron to the greatest extent. Signals are sent between the neurons. This is a non - linear behavior. Artificial neural networks. This behavior with regard to the modeling of nonlinear transfer function.

In the table below the results of the implementation of the model with neural network technique is shown. In the model proposed in this paper, we use multiple neural networks. After using neural network in the project, the proposed method of a multilayer perceptron's after the publication of the algorithm for training.

E-Coincidence Matrix for \$N-diabets (rows show actuals)

'Partition' = 1_Training	no	yes
no	351	55
yes	71	132
'Partition' = 2_Testing	no	yes
no	82	12
yes	24	41

Figure 7. The results of the model with neural network multiple techniques

⊡⊡Results	for output field	diabets						
⊟…Indiv	/idual Models							
	Comparing \$C-0	diabets with dial	bets					
	'Partition'	1_Traini	ng		2_Test	ing		
	Correct	4	87 79.	97%	1	12	70.4	44%
	Wrong	1	22 20.	03%		47	29.5	56%
	Total	6	09	159				
	Comparing \$N-0	diabets with dial	bets					
	'Partition'	1_Traini	ng		2_Test	ing		
	Correct	4	74 77.	83%	1	25	78.6	32%
	Wrong		35 22.	17%		34	21.3	38%
	Total	6	609			159		
Ė…Agre	ement betweer	⊨\$C-diabets \$N	l-diabets					
	'Partition'	1_Training			2_Testing			]
	Agree	512	84.079	6	130 8		1.76%	
	Disagree	97	15.939	6	29	1	8.24%	,
	Total	609	609		159			
	Comparing Agre	ement with diab	pets					
	'Partition'	1_Traini	ng		2_Test	ing		]
	Correct	4	32 84.	38%	1	04	80%	ف
	Wrong		80 15.	62%		26	20%	6
	Total	5	12		1	30		

Figure 8. The results of the implementation of the proposed model

#### **Comparison of the results**

For comparison, the proposed method with other existing methods has been included in the table. In which all of the methods discussed in this dissertation with accuracy and classification criteria.

Went on ways to implement the first compared and then each performance compared with each other in our schedules. In the model implemented all specimens were used, and all the default settings were considered software, for example, including input layer with 4 neuron, a hidden layer with the number of 19 and output layer with the number of the cases, including one neuron. Finally, after the implementation of the following results business model that is shown below.

In Figure above the results of the implementation of natural birth for training and test data is shown separately. First is also the performance of the decision and neural network separately, and then to the proposed model.

In Table up as well as you can see the performance of the algorithm proposed in comparison higher accuracy with other techniques. This shows that for reducing the dimensions of feature in diagnosis is better algorithms use playoffs. Or indeed all the way from a combination of functions.

#### References

[1]. Morgan Kaufmann, "Data Mining - Concepts and Techniques - Second Edition Second Edition [2006], Jiawei Han University of Illinois at Urbana-Champaign "MichelineKamber" [2]. Ian H. Eibe Frank Mark A. Hall, "Data Mining- Practical Machine Learning Tools and Techniques Third Edition" [2011] Shahrabi J, Shakoorniaz V. Concepts of data mining in; Tehran Shahrabi J, ZolghadrShojaei A. Advanced data mining: Concepts and algorithms.

[5] Kantardzic Mehmed, "*Data Mining: Concepts, Models, Methods, and Algorithms*" Second edition. John Wiley & Sons press, 2011.

[6] Larose, Daniel T., "Discovering knowledge in data An Introduction to Data Mining", John Wiley & Sons press, 2005.

[7] Jiawei Han, Micheline Kamber, Jian Pei , "Data Mining: Concepts and Techniques", Third Edition, elsivier,2006.
[8]. ZhaoHui, "Tang" and "Jamie MacLennan" [2005] Data Mining with SQL Server [2005] Bellazzi R, Zupan B. Predictive data mining in clinical medicine: Current issues and guidelines. International Journal of Medical Informatics,

[10] Data mining Mhmd Kamtardzyk Translation: AmirSamiraSummer1385Publisher: Computer Science

[11] Ozbakır Lale, Baykasoglu Adil and Kulluk Sinem. "A soft computing-based approach for integrated training and rule extraction from artificial neural networks:DIFACONN-miner", Applied Soft Coputing, 2010, Vol.10,pp: 304-317.

[12] Martin T.Hagan, Howard B.Demuth,"Neural Network Design", PWS PublishingCompany, 1996.