# Comparative Study of Diabetic Patient Data's Using Classification Algorithm in WEKA Tool

P.Yasodha

Pachiyappa's college for women,
Sri Chandrasekharendra Saraswathi Viswa Mahavidyalaya
Kanchipuram, India

N.R. Ananthanarayanan

Pachiyappa's college for women,
Sri Chandrasekharendra Saraswathi Viswa Mahavidyalaya
Kanchipuram, India

**Abstract**: Data mining refers to extracting knowledge from large amount of data. Real life data mining approaches are interesting because they often present a different set of problems for diabetic patient's data. The research area to solve various problems and classification is one of main problem in the field. The research describes algorithmic discussion of J48, J48 Graft, Random tree, REP, LAD. Here used to compare the performance of computing time, correctly classified instances, kappa statistics, MAE, RMSE, RAE, RRSE and to find the error rate  measurement for different classifiers in weka .In this paper the data classification is diabetic patients data set is developed by collecting data from hospital repository consists of 1865 instances with different attributes. The instances in the dataset are two categories of blood tests, urine tests. Weka tool is used to classify the data is evaluated using 10 fold cross validation and the results are compared. When the performance of algorithms, we found J48 is better algorithm in most of the cases.

**Keywords-** Data Mining, Diabetics data, Classification algorithm, Weka tool

## 1. INTRODUCTION

The main focus of this paper is the classification of different types of datasets that can be performed to determine if a person is diabetic. The solution for this problem will also include the cost of the different types of datasets. For this reason, the goal of this paper is classifier in order to correctly classify the datasets, so that a doctor can safely and cost effectively select the best datasets for the diagnosis of  the disease. The major motivation for this work is that diabetes affects a large number of the world population and it's a hard disease to diagnose. A diagnosis is a continuous process in which a doctor gathers information from a patient and other sources, like family and friends, and from physical datasets of the patient. The process of making a diagnosis begins with the identification of the patient's symptoms. The symptoms will be the basis of the hypothesis from which the doctor will start analyzing the patient. This is our main concern, to optimize the task of correctly selecting the set of medical tests that a patient must perform to have the best, the less expensive and time consuming diagnosis possible. A solution like this one, will not only assist doctors in making decisions, and make all this process more agile, it will also reduce health care costs and waiting times for the patients. This paper will focus on the analysis of data from a data set called Diabetes data set.

## 2. RELATED WORK

The few medical data mining applications as compared to other domains. [4] Reported their experience in trying to automatically acquire medical knowledge from clinical databases. They did some experiments on three medical databases and the rules induced are used to compare against a set of predefined clinical rules. Past research in dealing with this problem can be described with the following approaches: (a) Discover all rules first and then allow the user to query and retrieve those he/she is interested in. The representative approach is that of templates [3]. This approach lets the user to specify what rules he/she is interested as templates. The system then uses the templates to retrieve the rules that match the templates from the set of discovered rules. (b) Use constraints to constrain the mining process to generate only relevant rules. [12] Proposes an algorithm that can take item constraints specified by the user in the association rule mining processor that only those rules that satisfy the user specified item constraints are generated.
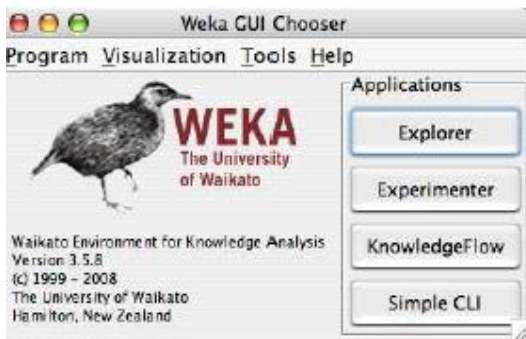
The study helps in predicting the state of diabetes i.e., whether it is in an initial stage or in an advanced stage based on the characteristic results and also helps in estimating the maximum number of women suffering from diabetes with specific characteristics. Thus patients can be given effective treatment by effectively diagnosing the characteristics.

Our research work based on the concept from Data Mining is the knowledge of finding out of data and producing it in a form that is easily understandable and comprehensible to humans in general. These further extended in this to make an easier use of the data's available with us in the field of Medicine.

The main use of this technique is the have a robust working model of this technology. The process of designing a model helps to identify the different blood groups with available Hospital Classification techniques for analysis of Blood group data sets. The ability to identify regular diabetic patients will enable to plan systematically for organizing in an effective manner. Development of data mining technologies to predict treatment errors in populations of patients represents a major advance in patient safety research.

## 3. MATERIALS AND METHODS

The **WEKA** (Waikato Environment for Knowledge Analysis) software was developed in the University of New Zealand. A number of data mining methods are implemented in the WEKA software. Some of them are based on decision trees like the J48 decision tree, some are rule-based like ZeroR and decision tables, and some of them are based on probability and regression, like the Naïve Bye's algorithm. The data that is used for WEKA should be made into the ARFF (Attribute Relation file format) format and the file should have the extension dot ARFF (.arff). WEKA is a collection of machine learning algorithms for solving real world data mining problems. It is written in Java; WEKA runs on almost any platform and is available on



the web at www.cs.waikato.ac.nz/ml/weka.

### 3.1. DATA PREPROCESSING

An important step in the data mining process is data preprocessing. One of the challenges that face the knowledge discovery process in medical database is poor data quality. For this reason we try to prepare our data carefully to obtain accurate and correct results. First we choose the most related attributes to our mining task.

### 3.2. DATA MINING STAGES

The data mining stage was divided into three phases. At each phase all the algorithms were used to analyze the health datasets. The testing method adopted for this research was parentage split that train on a percentage of the dataset, cross validate on it and test on the remaining percentage. Sixty six percent (66%) of the health dataset which were randomly selected was used to train the dataset using all the classifiers. The validation was carried out using ten folds of the training sets. The models were now applied to unseen or new dataset which was made up of thirty four percent (34%) of randomly selected records of the datasets. Thereafter interesting patterns representing knowledge were identified.

### 3.3 PATTERN EVALUATION

This is the stage where strictly interesting patterns representing knowledge are identified based on given metrics.

### 3.4 EVALUATION MATRICS

In selecting the appropriate algorithms and parameters that best model the diabetes forecasting variable, the following performance metrics were used:

**3.4.1**. **Time:** This is referred to as the time required to complete training or modeling of a dataset. It is represented in seconds

**3.4.2. Kappa Statistic:** A measure of the degree of nonrandom agreement between observers or measurements of the same categorical variable.

**3.4.3. Mean Absolute Error:** Mean absolute error is the average of the difference between predicted and the actual value in all test cases; it is the average prediction error.

**3.4.4. Mean Squared Error:** Mean-squared error is one of the most commonly used measures of success for numeric prediction. This value is computed by taking the average of the squared differences between each computed value and its corresponding correct value. The mean-squared error is simply the square root of the mean-squared-error. The mean-squared error gives the error value the same dimensionality as the actual and predicted values.

**3.4.5. Root relative squared error:** Relative squared error is the total squared error made relative to what the error would have been if the prediction had been the average of the absolute value. As with the root mean-squared error, the square root of the relative squared error is taken to give it the same dimensions as the predicted value.

**3.4.6. Relative Absolute Error:** Relative Absolute Error is the total absolute error made relative to what the error would have been if the prediction simply had been the average of the actual values.

## 4. METHODOLOGY

### 4.1. CLASSIFICATION

Classification is a data mining (machine learning) technique used to predict group membership for data instances. For example, you may wish to use classification to predict whether the weather on a particular day will be "sunny", "rainy" or "cloudy". Popular classification techniques include decision trees and neural networks.

### 4.2. J48 Pruned Tree

J48 is a module for generating a pruned or unpruned C4.5 decision tree. When we applied J48 onto refreshed data, we got the results shown as below on Figure .
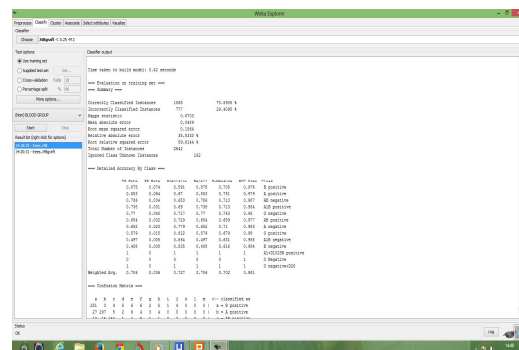
Fig- 1: J48 Tree

### 4.3. J48 graft

Perhaps C4.5 algorithm which was developed by Quinlan [13] is the most popular tree classifier till today. Weka classifier package has its own version of C4.5 known as J48 or J48graft
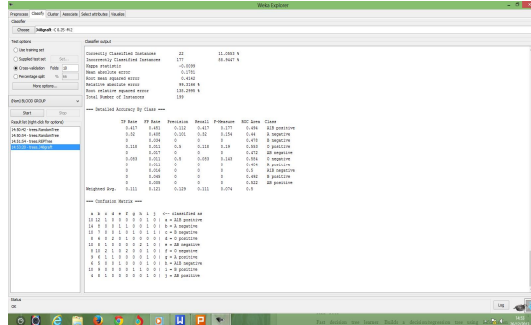

Fig-2: J48 Graft

### 4.4. LAD tree

LADTree is a class for generating a multiclass alternating decision tree using logistics strategy. LADTree produces a multi- class LADTree. It has the capability to have more than two class inputs. It performs additive logistic regression using the Logistics Strategy.
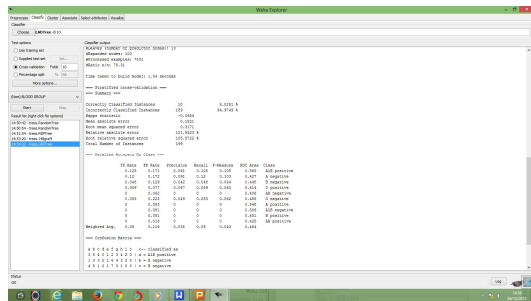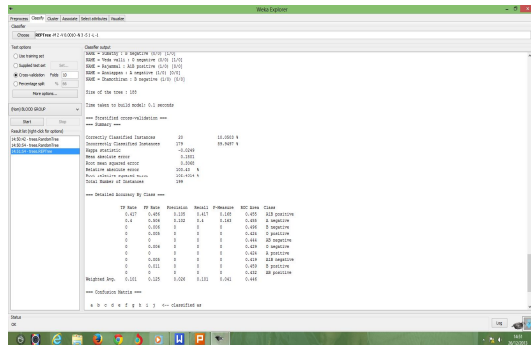

Fig-3: LAD Tree

### 4.5. REP Tree

Fast decision tree learner. Builds a decision/regression tree using information gain/variance and prunes it using reduced-error pruning (with back fitting). Only sorts values for numeric attributes once. Missing values are dealt with by splitting the corresponding instances into pieces (i.e. as in C4.5).



## 5. RESULT AND DISCUSSION

J48 algorithm was selected for the prediction because out of the five classifiers used to train the data, it had the best performance measures.

=== Run information ===

Scheme:     weka.classifiers.trees.J48 -C 0.25 -M 2
Relation:   py1
Instances:  2804
Attributes: 11
        NAME
        GENDER
        AGE
        HEIGHT
        BLOOD GROUP
        BLOOD SUGAR(F)
        BLOOD SUGAR (PP)
        BLOOD SUGAR (R)
        URINE SUGAR(F)
        URINE SUGAR(PP)
        URINE SUGAR (R)
Test mode:   evaluate on training data

=== Classifier model (full training set) ===

J48 pruned tree
------------------
J48 pruned tree
------------------

AGE <= 46
|   AGE <= 35
|   |   GENDER = Male
|   |   |   AGE <= 26: B positive (2.0/1.0)
|   |   |   AGE > 26: A positive (3.0/1.0)
|   |   GENDER = Female
|   |   |   AGE <= 34: O negative (2.0)
|   |   |   AGE > 34: A positive (2.0/1.0)
|   AGE > 35: B positive (7.0/4.0)
AGE > 46
|   GENDER = Male
|   |   AGE <= 60: O positive (5.0/3.0)
|   |   AGE > 60: AB positive (4.0/2.0)
|   GENDER = Female
|   |   AGE <= 63
|   |   |   AGE <= 55: AB positive (2.0/1.0)
|   |   |   AGE > 55: A1B positive (4.0/2.0)
|   |   AGE > 63: A negative (2.0/1.0)
Number of Leaves  :         10
Size of the tree :     19

Time taken to build model: 0.29 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly   Classified   Instances           1865
70.5905%
Incorrectly   Classified   Instances          777
29.4095%

Kappa statistic              0.6703
Mean absolute error          0.0489
Root mean squared error         0.1564
Relative absolute error      35.5333 %

Root relative squared error          59.6144%
Total Number of Instances            2642
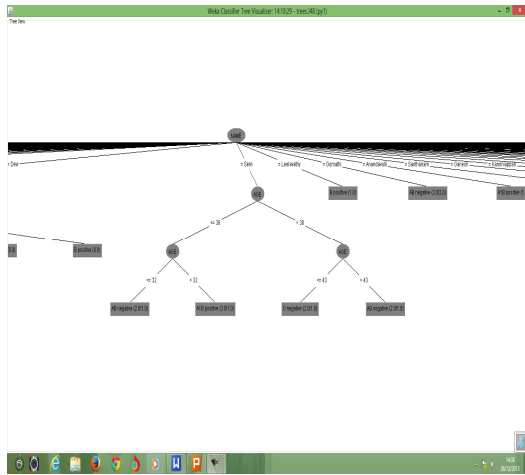Ignored Class Unknown Instances          162



Fig -6: **VISUALISE THE TREE**

| CLAS SIFIE R | CORRE CTLY CLASSI FIED INSTAN CES | TP RA TE | FP RA TE | PR ECI SIO N | RE CA LL | F-M E AS U R E | R O C A R E A |
|---|---|---|---|---|---|---|---|
| **J48** | 1865 (70.5%) | 0.70 6 | 0.03 6 | 0.72 7 | 0.70 6 | 0.7 02 | 0. 9 8 1 |
| **J48 GRAF T** | 1524 (57.6%) | 0.60 7 | 0.02 4 | 0.67 8 | 0.52 0 | 0.6 00 | 0. 7 8 1 |
| **LAD TREE** | 553 (20.9%) | 0.05 | 0.11 6 | 0.03 8 | 0.05 | 0.0 43 | 0. 4 6 4 |
| **RAND OM TREE** | 350 (13.2%) | 0.11 1 | 0.12 2 | 0.09 8 | 0.11 1 | 0.0 7 | 0. 4 6 4 |
| **REP TREE** | 348 (0.13%) | 0.13 2 | 0.13 2 | 0.01 7 | 0.13 2 | 0,0 31 | 0. 5 |

**Table-1: DIFFERENT PERFORMANCE METRICES RUNNING IN WEKA**

In this study, we examine the performance of different classification methods that could generate accuracy and some error to diagnosis the data set. According to above Table 1 , we can clearly see the highest accuracy is 70.5% belongs to J48 and lowest accuracy is 0.13% that belongs to REP. The total time required to build the model is also a crucial parameter in comparing the classification algorithm.

|  | **J48** | **J48GR AFT** | **RAND OM TREE** | **REP** | **LAD** |
|---|---|---|---|---|---|
| **TIME** | 0.29 | 0.42 | 0.02 | 0.05 | 1.85 |
| **CORRECTL Y CLASSIFIE D INSTANCES** | 1865 (70.5%) | 1524 (57.6%) | 350 (13.2% ) | 348 (0.13 %) | 553 (20.9%) |
| **KAPPA STATISTIC** | 0.011 | 0.6700 | 0.011 | 0.012 | 0.0654 |
| **MAE** | 0.0123 | 0.0480 | 0.1798 | 0.1377 | 0.1821 |
| **RMSE** | 0.1154 | 0.1560 | 0.3199 | 0.2624 | 0.3171 |
| **RAE%** | 12.53% | 35.50% | 100.24 % | 99.98 % | 101.55 % |
| **RRSE%** | 22.61% | 58.63% | 106.82 % | 100% | 105.87 % |

**Table- 2: ERRORS MEASUREMENT FOR DIFFERENT CLASSIFIERS IN WEKA**

Based on above table, we can compare errors among different classifiers in WEKA. We clearly find out that J48 is the best, second best is the j48 graft ,LAD, REP & random. An algorithm which has a lower error rate will be preferred as it has more powerful classification capability and ability in terms of medical and bio informatics fields.

# 6. CONCLUSION AND FUTURE WORK

The objective of this study is to evaluate and investigate FIVE selected classification algorithms based on WEKA. The best algorithm in WEKA is J48 classifier with an accuracy of 70.59% that takes 0.29 seconds for training. They are used in various healthcare units all over the world. In future to improve the performance of these classification.

I had been use the data mining classifiers to generate decision tree format. In this paper WEKA software for my experiment. Identify the diabetic patient's behavior using the classification algorithms of data mining. The analysis had been carried out using a standard blood group data set and using the J48 decision tree algorithm implemented in WEKA. The research work is used to classify the diabetic patient's based on the gender, age, height & weight, blood group, blood sugar(F), blood sugar(PP), urine sugar(F), urine sugar(PP). The J48 derived model along with the extended definition for identifying regular patients provided a good classification accuracy based model.

The distribution of blood groups in both positive and negative are shown in Table-1. Overall blood group A was the commonest (24.03 %), followed by B (18.77%), AB (19.11%), O (23.65) and A1B (17.14%).

| Blood group spectrum | Nos (%) | +ve (%) | –ve (%) |
|---|---|---|---|
| A | 635 (24.03) | 348 13.17 | 287 10.85 |
| B | 496 (18.77) | 289 (10.93) | 207 (7.83) |
| AB | 505 (19.11) | 196 (7.41) | 309 (11.69) |
| A1B | 453 (17.14) | 300 (11.35) | 153 (5.79) |
| O | 625 (23.65) | 345 (10.59) | 280 (13.05) |

**Table-3: Spectrum of Blood groups +ve and -ve in major population. (n-2642)**

In the present blood group-A was the predominant (24.03%) while A1B was the least common (17.14%). Blood group "A" was the most predominant (24.03%) in both positive and negative subjects, followed by blood group A, B,O,A1B and AB.

The future work will be focused on using the other classification algorithms of data mining. It is a known fact that the performance of an algorithm is dependent on the domain and the type of the data set. Hence, the usage of other classification algorithms like machine learning will be explored in future.

The future work can be applied to blood groups to identify the relationship that exits between diabetic, diagnosing cancer patients based on blood cells or predicting the cancer types on the blood groups, blood pressure, personality traits and medical diseases.

# 7. REFERENCES

[1] Mats Jontell, Oral medicine, Sahlgrenska Academy, Göteborg University (1998) "A Computerised Teaching Aid in Oral Medicine and Oral Pathology. " Olof Torgersson, department of Computing Science, Chalmers University of Technology, Göteborg.

[2] T. Mitchell, "Decision Tree Learning", in T. Mitchell, Machine Learning (1997) the McGraw- Hill Companies, Inc., pp. 52-78.

[3] Klemetinen, M., Mannila, H., Ronkainen, P., Toivonen, H., and Verkamo, A. I (1994) "Finding interesting rules from large sets of discovered association rules," CIKM.

[4] Tsumoto S., (1997)"Automated Discovery of Plausible Rules Based on Rough Sets and Rough Inclusion," Proceedings of the Third Pacific-Asia Conference (PAKDD), Beijing, China, pp 210-219.

[5] Liu B., Hsu W., (1996) "Post-analysis of learned rules," AAAI, pp. 828-834.

[6] Liu B., Hsu W., and Chen S., (1997) "Using general impressions to analyze discovered classification rules," Proceedings of the Third ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.

[7] Stutz J., P. Cheeseman. (1996) Bayesian classification (autoclass): Theory and results. In Advances in Knowledge Discovery and Data Mining. AAAI/MIT Press

[8] Witten Ian H., E. Frank, Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations, Ch. 8, © 2000 Morgan Kaufmann Publishers

[9] http://www.cs.waikato.ac.nz/ml/weka/, accessed 06/05/21.

[10] http://grb.mnsu.edu/grbts/doc/manual/ J48_Decision_T rees.html, accessed

[11] Wikipedia, ID3-algorithm (accessed 2007/12/09) (URL: http://en.wikipedia.org/wiki/ID3_algorithm)

[12] Srikant,R.,Vu,Q.andAgrawal,R.,(1997), "Mining association rules with item constraints," Proceedings of the Third International Conference on Knowledge Discovery and Data Mining, Newport Beach, USA, pp 67-73.