

## Random Forest v/s Scaled Conjugate Gradient To Predict Diabetes Mellitus

**Neha Shukla and Dr. Meena Arora**

*Student, Associate Professor  
CSE Department  
JSS Academy Of Technical Education*

### Abstract

Diabetes Mellitus is one of the real wellbeing challenges everywhere throughout the world. The pervasiveness of diabetes is expanding at a quick pace, falling apart human, financial and social fabric. Aversion and expectation of diabetes mellitus is progressively picking up enthusiasm for social insurance group. Albeit a few clinical choice emotionally supportive networks have been commended that fuse a few information digging methods for diabetes forecast and course of movement. These ordinary frameworks are ordinarily based either just on a solitary classifier or a plain mix thereof. As of late broad attempts are being made for enhancing the exactness of such frameworks utilizing gathering classifiers. This study takes after the procedures utilizing random forest tree as a base learner alongside standalone information mining procedure scaled conjugate gradient to characterize patients with diabetes mellitus utilizing diabetes hazard variables. This characterization is done crosswise over three diverse ordinal grown-ups bunches in PIMA indian dataset. Test result demonstrates that, general execution of adaboost group strategy is superior to anything sacking and in addition standalone random forest tree.

**Keywords:** diabetes mellitus, random forest tree, classification, prediction, scaled conjugate gradient

### 1 INTRODUCTION

Information mining is a standout amongst the most alluring interdisciplinary subfield of software engineering. It includes some computational procedure, factual methods, grouping, bunching and finding designs in huge information sets. The general objective of the information mining strategy is to extricate the valuable data from the

extensive information set and to change it into a reasonable configuration so it can be utilized for the future use. Diabetes Mellitus is a gathering of metabolic ailment in which the measure of sugar substance can't be controlled. There are principally four sorts of Diabetes Mellitus. They are Type1, Type2, Gestational diabetes, Congenital diabetes. Sort 1 likewise called as "Insulin ward Diabetes Mellitus" or "Adolescent Onset Diabetes Mellitus" happens when the human body disappointments to deliver insulin. They are portrayed by the loss of insulin creating beta cells. Sort 2 is additionally called as "Non Insulin subordinate Diabetes Mellitus" or "Grown-up onset diabetes". Non Insulin subordinate Diabetes Mellitus is described by the insulin resistance and is found in individual above age 40 i.e., human body can't adequately utilize the insulin that is created. Eating routine, practice and glucose level control ought to be kept up to manage the Type 2 diabetes. Gestational diabetes mellitus (GDM) are provisional diabetes that looks like Type 2 diabetes in a few viewpoints. It is the condition in which the pregnant ladies, without already analyzed diabetes show an expansion level of glucose in the blood. GDM is treatable under cautious therapeutic supervision and determines totally once the child is conceived. The Congenital diabetes is brought on because of hereditary deformities of insulin discharge, cystic fibrosis-related diabetes, steroid diabetes affected by high dosages of glucocorticoids. Diabetes mellitus causes genuine entanglements, for example, coronary illness, stroke, visual deficiency, kidney disappointment and malignancy. As indicated by World Health Organization (WHO), 37 crores of individuals are experiencing diabetes around the globe and it copies before the year 2030[1]. Around 48 lakhs of individuals were passed on in the year 2012. The vast majority of them have a place with lower and white collar class families [1].

Diabetes can be controlled by utilizing distinctive measures like insulin and eating regimen. For this it ought to be recognized as right on time as would be prudent and thusly give suitable treatment. The vast majority of the grouping, recognizing and diagnosing medications depend on synthetic and physical tests. Taking into account the deduction acquired from these outcomes, a specific infection can be anticipated. Forecast may have mistakes. This is because of various vulnerability of different parameters utilized for testing [2]. Such instabilities make the expectations wrong and keep the odds of curing the sickness. The registering office has been advanced with extraordinary headways. These headways gave by data innovation, orders the information, anticipate the results and finding of numerous sicknesses all the more precisely. The primary point of preference of data innovation is that a gigantic information stockpiling of past patient's records are kept up and checked by healing centers persistently for different references [2]. These restorative information helps the specialists to analyze distinctive examples in the information set. The examples found in information sets might be utilized for grouping, expectation and determination of the sicknesses [2].

## **2. LITERATURE REVIEW**

Number of information mining calculations has been proposed to characterize, anticipate and analyze diabetes. For this, information preprocessing ought to be

finished. It is a strategy that includes changing crude information into reasonable arrangement. This fills the missing qualities in the middle of the information. By investigating the information utilizing the qualities, it is feasible for a specialist to discover values that are startling and incorrect. Examines the performance of the Levenberg-Marquardt (LM) algorithm on a single dataset, the Pima Indian Diabetes dataset, attempting to minimize error in classifying the patients as diabetes positive or negative. The learning algorithm is applied on dynamically constructed neural network to minimize the error by continuously training the network until the optimum efficiency level is obtained. The performance of the approach is verified by performing a comparison study.[5]

Prevention and prediction of diabetes mellitus is increasingly gaining interest in healthcare community. Although several clinical decision support systems have been proposed that incorporate several data mining techniques for diabetes prediction and course of progression. These conventional systems are typically based either just on a single classifier or a plain combination thereof. Recently extensive endeavors are being made for improving the accuracy of such systems using ensemble classifiers. This study follows the adaboost and bagging ensemble techniques using J48 (c4.5) decision tree as a base learner along with standalone data mining technique J48 to classify patients with diabetes mellitus using diabetes risk factors.[11]

It proposed the use of a rule extraction algorithm, Re-RX with J48graft,combined with sampling selection techniques(sampling Re- RX withJ48graft)to achieve highly accurate, concise, and interpretable classification rules The use of this algorithm resulted in an average accuracy of 83.83% after 10 runs of 10-fold cross validation. Sampling Re-RX with J48 graft achieved substantially better accuracy and provided a considerably fewer average number of rules and antecedents than the original Re-RX algorithm [14]

Table 2: Comparison of supervised Algorithms based on performance

no	Algo	CTime	TP	FN	FP	TN	Acc %	Spec %	Sen %	CV Erate	P (Prec)	N (Prec)	BVErate
	C4.5	550ms	31	23	19	77	72	80%	57%	0.28	0.38	0.23	0.3196
	SVM	546ms	24	30	14	82	70.6667	85%	44%	0.29	0.368	0.2679	0.2929
	K-NN	640ms	20	34	18	78	65.3333	81%	37%	0.34	0.474	0.3036	0.3532
	PNN	546ms	42	12	39	57	66	59%	78%	0.34	0.482	0.1739	0.3406
	BLR	515ms	32	22	19	77	72.6667	80%	59%	0.2733	0.373	0.222	0.2754
	MLR	530ms	32	22	19	77	72.6667	80%	59%	0.2733	0.373	0.222	0.2754
	CRT	515ms	8	46	8	88	64	92%	15%	0.36	0.5	0.3433	0.3153
	CS-CRT	531ms	37	19	5	94	84.5161	95%	66%	0.36	0.119	0.1681	0.3153
	PLS-DA	452ms	25	21	16	83	74.4828	84%	54%	0.2667	0.314	0.2019	0.2726
	PLS-LDA	593ms	36	20	16	83	76.7742	84%	64%	0.2667	0.308	0.1942	0.2726

Legend: CTime- Computing Time, TP-True Positive, FN-False Negative, FP-FalsePositive, TN-True Negative, Acc- Accuracy, Spec- Specificity, Sen- Sensitivity, CV- Coefficient of Variation, Erate- Error Rate, P- Precision, N- Recall, BVErate- Balanced Error Rate

Table 1: The performance of various algorithms is discussed in this paper which I have shown in the table below.

### 3. PRELIMINARY

Diabetes is a long lasting constant condition that influences the human body by lessening the insulin which conveys glucose into the platelets. This expands the sugar level in the body prompting diverse entanglements like stroke, coronary illness, visual deficiency, kidney disappointment and demise. Diabetic patients by and large have the accompanying manifestations.

- Increased thirst
- Frequent pee
- Weight misfortune
- Increased yearning
- Slow-recuperating diseases
- Blurred vision
- Nausea and Vomiting

The accompanying therapeutic tests are utilized to analyze the diabetic mellitus [3]

- Urine test
- Fasting blood glucose level
- Random blood glucose level
- Oral glucose resistance test
- Glycosylated haemoglobin (HbA1c)

#### 3.1 DATASET

The information set decided for grouping and test reproduction depends on Pima Indian Diabetic Set from University of California, Irvine (UCI) Repository of Machine Learning databases. The patients under thought are the Pima Indian populace living in Arizona, USA. More than half Pima Indian Population is experiencing diabetes and 95% of them are because of the overweight. Number of exploration has been done on these populaces demonstrated that weight is the primary driver for the diabetes. The information set essentially contain 9 properties and 768 number of occasions [4]. The 8 such traits alongside images are recorded in Table 1.

- 1) Total no of times pregnant
- 2) Glucose resilience test to discover the plasma glucose level focus in salivation.
- 3) Diastolic circulatory strain measured in mmHg for weight level(BP)
- 4) Body Mass Index (BMI)  
BMI= Patients weight in kg/(Patients stature in meter)<sup>2</sup>
- 5) Skin rashes and thickness fold in mm (Triceps)
- 6) 2-hour Serum Insulin in mu U/ml (INSULIN)
- 7) Age in years
- 8) Diabetes family work
- 9) Diabetes Class Variable (1 shows diabetic test is certain (nearness) and 0 demonstrates test is negative (nonappearance))

### 4. PERFORMANCE COMPARISONS OF ALGORITHMS

Order precision (ACU) is the most widely recognized strategy utilized for assessment

of execution. Computation of exactness is performed by taking proportion of really ordered examples (genuine negative, genuine positive) to the aggregate number of tests.

$$\text{Exactness} = \text{Truly ordered examples} / \text{all out specimens} \quad (1)$$

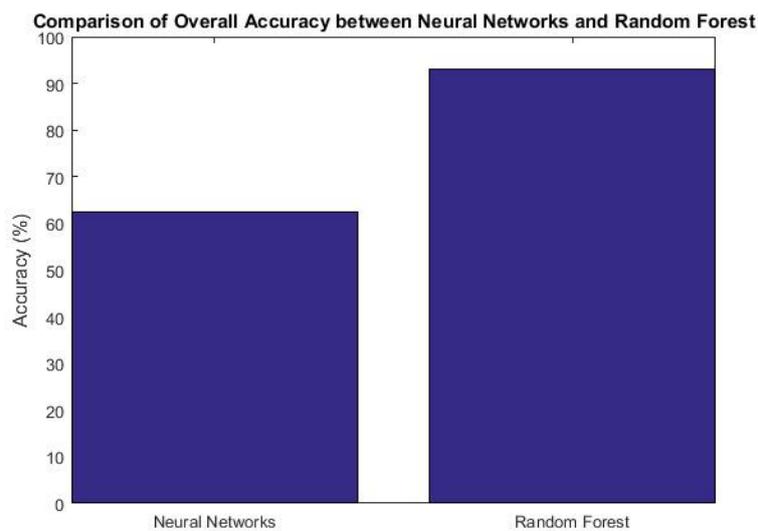
Another assessment strategies utilized for measuring execution are Sensitivity and Specificity. Affectability is computed by partitioning the genuine positive (TP) tests to the entirety of genuine positive (TP) and false negative (FN) tests.

$$\text{Affectability} = \text{TP} / (\text{TP} + \text{FN}) \quad (2)$$

Specificity is computed by partitioning the genuine negative (TN) tests to the total of genuine negative and false positive (FP) tests.

$$\text{Specificity} = \text{TN} / (\text{TN} + \text{FP}) \quad (3)$$

SCG have the accuracy esteem 56%. MSE calculation is having the minimum precision rate and it show errors when high dimensional information sets are given as info. By examining the table obviously random forest tree demonstrates the most extreme precision contrasted with different calculations. To accomplish the better exactness, results were examined for various k values. Random forest tree has the accuracy of 92.96% accuracy.



**Fig 1:** The comparison between the accuracy of scaled conjugate gradient and Random Forest Tree

## 5. SUGGESTION

Random forest tree calculations play out the arrangement with a higher effectiveness and decreased unpredictability. The prediction can be increased by increasing the k fold cross validation but due to it the error will increase and accuracy will be varied. So it can be tries but not implemented with knowing consequences.

## 6. CONCLUSION

Information mining and machine learning calculations in the therapeutic field extricates diverse concealed examples from the medicinal information. They can be utilized for the investigation of vital clinical parameters, expectation of different maladies, estimating errands in solution, extraction of medicinal learning, treatment arranging backing and patient administration. Various calculations were proposed for the expectation and determination of diabetes. These calculations give more precision than the accessible customary frameworks. This paper incorporates calculations of random forest tree and scaled conjugate gradient. From the perception SCG have the slightest characterization exactness and Random forest tree give the better grouping precision results.

## 7. REFERENCES

- [1] D Reby, S Lek, I Dimopoulos, J Joachim, Ja Lauga, S Aulagnier “Artificial neural networks as a classification method in the behavioural sciences” *Elsevier Behavioural Processes* 40 (1997) 35-43
- [2] Dr S kumar N M, Eswari T, Sampath P & Lavanya S “Predictive Methodology for Diabetic Data Analysis in Big Data” *ScienceDirect ISBCC’15 Procedia Computer Science* 50 ( 2015 ) 203-208
- [3] Dr.V.Karthikeyani, I.Parvin Begum “Comparison a Performance of Data Mining Algorithms (CPDMA) in Prediction Of Diabetes Disease” *International Journal on Computer Science and Engineering (IJCSE)* Vol. 5 No. 03 Mar 2013 205-210
- [4] D.Senthil Kumar, G.Sathyadevi and S.Sivanesh “Decision Support System for Medical Diagnosis Using Data Mining” *IJCSI International Journal of Computer Science Issues*, Vol. 8, Issue 3, No. 1, May 2011
- [5] F Marir, H Saida, and F Al-Obeidata “Mining the Web and Literature to Discover New Knowledge about Diabetes” *ScienceDirect Procedia Computer Science* 83 (2016) 1256-1261
- [6] Indiramma M, Raghavendra S “Classification and Prediction Model using Hybrid Technique for Medical Datasets” *International Journal of Computer Applications (0975-8887)* Volume 127-No.5, October 2015
- [7] K Polat, S Günes “ An expert system approach based on principal component analysis and adaptive neuro-fuzzy inference system to diagnosis of diabetes disease” *ScienceDirect Digital Signal Processing* 17 (2007) 702-710
- [8] M. Durairaj, V. Ranjani “Data Mining Applications In Healthcare Sector: A Study” *International Journal of Scientific & Technology Research* Volume 2, Issue 10, October 2013
- [9] N Khana, D Gaurava, T Kandl “Performance Evaluation of Levenberg-Marquardt Technique in Error Reduction for Diabetes Condition Classification” *SciVerse ScienceDirectProcedia Computer Science* 18 ( 2013 ), *ICCS 2013*,2629-2637

- [10] N Chandgude, S Pawar “A survey on diagnosis of diabetes using various classification algorithm” *International Journal on Recent and Innovation Trends in Computing and Communication* ISSN: 2321-8169 Volume: 3 Issue: 12 6706-6710
- [11] S Perveen, M Shahbaz, A Guergachi, K Keshavjee “Performance analysis of data mining classification technique to predict diabetes” *ScienceDirect, Procedia Computer Science* 82(2016) SDMA 2016, 115-121
- [12] Kaur S and Dr. R.K.Bawa “Future Trends of Data Mining in Predicting the Various Diseases in Medical Healthcare System” *International Journal of Energy, Information and Communications* Vol.6, Issue 4 (2015), pp.17-34
- [13] Vijayan V and Aswathy R “Study of Data Mining Algorithms for Prediction and Diagnosis of Diabetes Mellitus” *International Journal of Computer Applications* (0975-8887) Volume 95-No.17, June 2014
- [14] Hayashi Y. and Yukita S. “Rule extraction using Recursive-Rule extraction algorithm with J48 graft combined with sampling selection techniques for the diagnosis of type2 diabetes mellitus in the Pima Indian dataset” *Elsevier Informatics in Medicine Unlocked* 2(2016)92-104

