# GIHAT: An Efficient Prediction Technique for Measure for Diabetes Mellitus

R. Vanathi
Research Scholar
P.G. and Research Department of CS
Khadir Mohideen College, Adirampattinam
Thanjavur (Dist), Tamil Nadu, India

Dr. A. Shaik Abdul Khadir
Associate Professor
P.G. and Research Department of CS
Khadir Mohideen College, Adirampattinam
Thanjavur (Dist), Tamil Nadu, India.

**Abstract---** The medical service industry is a consistently developing field, producing trillions of information consistently. The modernization of the area has an immediate association with this incremental extent. These acquired informational collections are somewhat organized however for the most part unstructured in nature. These acquired information must be prepared with most extreme care to determine finish usable examples for subjective and prescient investigations. These gigantic records of information, in the wake of handling, when utilized, will turn out to be very unpredictable. Diabetes is a lifetime disease marked by elevated levels of sugar in the blood. It is the second leading cause of sightlessness and renal disease worldwide. Sort 2 diabetes mellitus (S2DM) is genuine and expensive metabolic illness that is a developing worries among peoples .S2DM is related with various comorbid conditions that can prompt negative patient results. Comorbid endless torment is extremely basic in S2DM because of the nearness of diabetic neuropathy and musculoskeletal conditions that are related with delayed hyperglycemia. This Paper using **General Integrated High Availability Transaction (GIHAT)** algorithm concentrates on the causes, sorts, and factors influencing DM (diabetes mellitus), preventive measures, and treatment of diabetes other than those directly associated with Diabetic Patients structured and unstructured data-sets .This algorithm executed in "R" Programming used for statistical analysis which provides the accurate results comparing existing algorithms.

Keywords---*Sort2diabetes Mellitus, Predictive Diabetes Analysis, R programming, and Statistical Analysis*

_____*****_____

## I.    INDRODUCTION

Datamining is an iterative procedure in which advance is characterized by revelation, through either manual or programmed strategies. Information mining is most helpful in an exploratory investigation situation in which there are no foreordained thoughts regarding what will constitute an "intriguing" result. Information mining is the look for new, profitable, and nontrivial data in extensive volumes of information.

Practically speaking, the two essential objectives of information mining have a tendency to be forecast and depiction [1]. Expectation includes utilizing a few traits or fields in the informational index to anticipate obscure or future estimations of different characteristics of intrigue. Then again, portrayal concentrates on discovering designs for portraying the information so people can translate it. To accomplish the objectives of forecast and depiction one must take after an information mining process. There are a wide range of variants of information mining forms and numerous feelings on the most proficient method to approach them. This paper concentrates on the RapidMiner programming bundle to pre-process and dissect diabetes information and mine a diabetes expectation show. "Kidney disappointment is a destructive complexity of diabetes, and Pimas, so far as should be obvious, have the world's heightest rate of sort 2 diabetes." [3] The goal of this investigation is to see any broad connections between various patients attribute furthermore, the inclination to create diabetes.

RapidMiner programming bundle bolsters all ventures of information mining process [2]. It is a Java-based open-source programming and can be utilized as a Java API. It additionally gives a basic and well disposed of GUI. RapidMiner utilizes interior XML portrayals to guarantee institutionalized trade configuration of information mining tests. Utilizing RapidMiner, we will effortlessly convey an examination report and a forecast display. The examination report will outline the information and their relationship to treating diabetes. The forecast model will be a choice tree that should help in foreseeing whether a patient will create diabetes utilizing the information accumulated. The informational index utilized as a part of this task is excerpted from the UCI Machine Learning Repository [4]. The Pima Indians Diabetes Data Set contains 8 classifications and 768 examples assembled from a bigger database having a place with the National Institute of Diabetes and Stomach related and Kidney Diseases. The determinations of these examples are: All patients are females at any rate 21 years of age of Pima Indian legacy.

## II.    DIABETES DATASETS DESCRIPTIONS

By and large diabetes is classified in to two sorts Type 1 and Type 2 illness. Sort 1 may account for 5% to 10%, Type-2 90% to 95%, Gestational diabetes amid Pregnancy 5% to10% different sorts Diabetes Miletus 1% to 5% .

It is use to calculations of the Data Mining Classification Methods

| Types | Levels |
|-------|--------|
| **Beginning Stage** | **5-10%** |
| **Final Stage** | **90-95%** |
| **During Pregnancy** | **5- 10%** |
| **Diabetes Miletus** | **1-5%** |

**Table 1. Classification of diabetes**

## 2.1 PREPROCESSING OF DATASETS

A large portion of the informational collections utilized as a part of Data mining was definitely not essentially assembled in light of a particular objective. Some of them may contain mistakes, anomalies or missing values. Keeping in mind the end goal to utilize those informational collections in the information mining process, the information needs to experience pre-processing, utilizing information cleaning, discretization and information change [5]. It has been assessed that information planning alone records for 60% of constantly and exertion extended in the whole information mining process [6]

The Department of Industrial Policy and Promotion (DIPP) released a press note with some highlights pertaining to the Foreign Direct Investments in the field of Healthcare. It has few noteworthy points which are stated below

- OrbiMed, a healthcare-dedicated investment firm, plans to invest around 4 crores in INR in Kolkata-based pathology and radiology services chain Suraksha Diagnostics. The main deal for this investment is to expand the diagnostics chain's laboratory network across India and to bring state of the art technologies at its disposal to serve people better.

- Attune Technologies Private Limited, a Chennai-based healthcare technology firm, has raised 1 crore INR in a Series B funding from Qualcomm Ventures and Norwest Venture Partners. This investment paves way to expand its digital healthcare solutions from the current 200 hospitals and laboratories to 30,000 such facilities globally.

- Sanofi-Synthelabo (India) Limited invested Rs 90 crore in Apollo Sugar Clinics Limited (ASCL), a unit of its subsidiary Apollo Health and Lifestyle Limited and the investment also covers Diabetes Mellitus related treatment research.

Some of the major initiatives taken by the Government of India to promote Indian healthcare industry are as follows:

## III. RELATED WORK

Information examination that should be possible with the RapidMiner programming incorporates charts and tables, and in addition different outlines and plots. The RapidMiner Histogram Color Network was utilized to outwardly look at the estimations of the traits and see the associations with the Class trait esteems (Yes, No), which finds that the patients with higher Plasma-Glucose esteems are probably going to create diabetes and most with low Plasma-Glucose values don't create diabetes inside five years. Additionally breaking down this connection between Plasma- Glucose and Class by utilizing a container plot certifies the above perception. To help clear up whether the perception may be of esteem, a Naïve Bayes learning device is connected, where the ascribes are considered to be irregular factors, and the information are thought to be known. The parameters are viewed as originating from a dissemination of conceivable esteems, and Bayesians look to the watched information to give data on likely parameter esteems [5]. This checks the underlying perception gives off an impression of being right. After the information preprocessing, our next objective is to discover for the most part relationship in the information keeping in mind the end goal to comprehend the connections between the traits and regardless of whether the patients go ahead to create diabetes. With the discretization of numerical traits, we will center on the sub- gatherings (canisters) made rather than the singular estimations of the credits to limit the unpredictability of the examination without losing precision. RapidMiner gives an exceptionally helpful apparatus, BasicRuleLearner, for helping thin perception down, which filters through the information and discovers general relationship rules, for example,

In the event that Plasma-Glucose = high then Yes (124/60)
On the off chance that Pregnant = medium the No (28/65)
On the off chance that DPF = low the No (26/50)

One may understand that a few guidelines have low exactness, in this manner might delude. It would be off base to just take a gander at one trait and after that look at the outcome. It must be smarter to see the outcomes when at least two characteristics are joined, also to consolidate all characteristics[18] . This can be accomplished utilizing CHAID (Chi-squared Automatic Interaction Detector) choice tree [11]. CHAID recognizes communication between factors in the informational index by recognizing discrete gatherings of respondents, and looks to anticipate what the effect will be on the reliant factors by taking their reactions to logical factors. Since CHAID requires insights information, it isn't important to discretize numerical factors. Information examination and concealed relationship uncover that have been made so far can be utilized to either alter the

88

trait or help show signs of improvement comprehension of the characteristic esteems. Presently we have to develop a prescient model to appraise whether a patient will create diabetes inside a satisfactory level of conviction, rather than essentially revealing insight about the information itself. RapidMiner gives intends to this reason. pick two principle choices: the ID3 Algorithm and the Decision Tree.

A decision tree can be learnt by part the source informational collection into subsets in view of a property estimation test [12]. This procedure is rehashed on each determined subset in a recursive way. The recursion is finished when part is either non-achievable or a particular grouping can connected to every component of the inferred subset. An arbitrary timberland classifier utilizes a number of decision trees, so as to enhance the Order rate. The decision isn't just useful in speaking to the present information connections, yet additionally ready to apply other information to the calculation and test how well it works at anticipating the result. RapidMiner backings to create a decision tree. A some portion of the decision tree naturally created by RapidMiner [12] where the Plasma-Glucose trait is picked as the root hub. This additionally strengthens our unique perception in the earlier segment. The decision tree discloses to us that Plasma-Glucose is the principle characteristic that will lead us to knowing regardless of whether a patient will create diabetes [14]. The

genuine information set expresses that 248 patients create diabetes and 475 do not. This decision tree predicts that there are 200 patients that create diabetes and 523 that don't. Of those 200 patients, 19 aren't right, which brings the remedy expectations down to 181. This implies the decision tree has 72% of exactness.

## IV. PROPOSED SYSTEM

### R Programming

R is an open source programming dialect and programming condition for measurable processing and designs that is bolstered by the R Foundation for Statistical Computing. [6] The R dialect is generally utilized among analysts and information diggers for creating factual programming and information analysis. Polls, studies of information mineworkers, and investigations of academic writing databases demonstrate that R's ubiquity has expanded considerably in late years.

R is a GNU package. The source code for the R programming condition is composed fundamentally in C, FORTRAN, and R. R is uninhibitedly accessible under the GNU General Public License, and pre-accumulated twofold forms are accommodated different working frameworks. While R has a charge line interface, there are a few closures available.
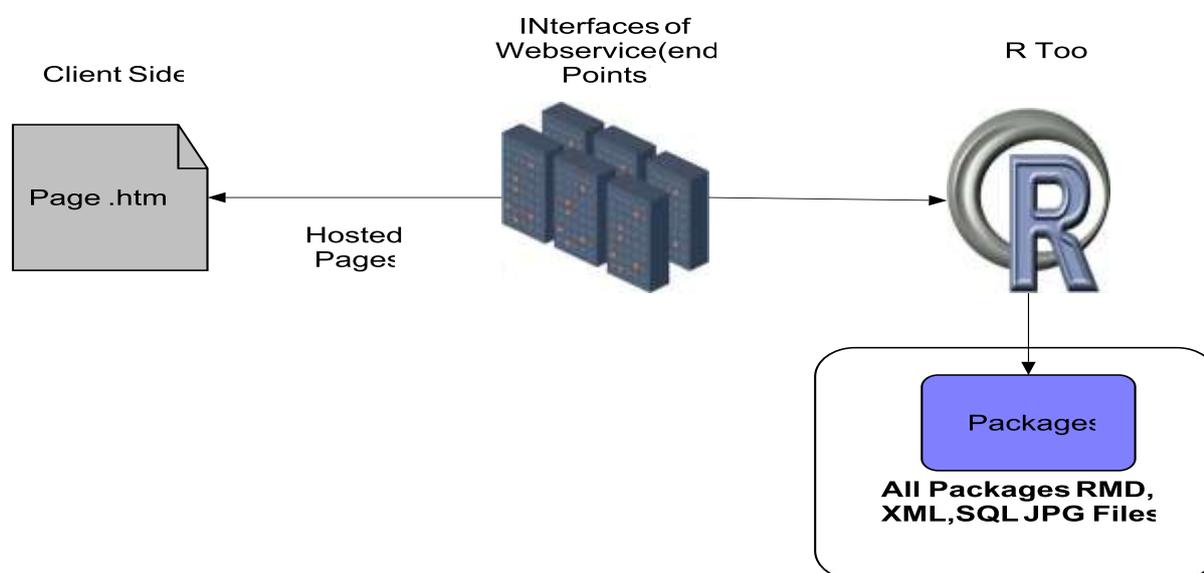


**Fig 1. R Tool Structure**

### Data-set Incrimination

The data obtained for research comprises of both structured and unstructured data-sets. These data-sets comprises of data's from both EHR's and also patient's paper records. In this stage, cleanse the data using GIHAT [General Integrated High Availability Transaction] algorithm. In this algorithm, prepare the data and make it
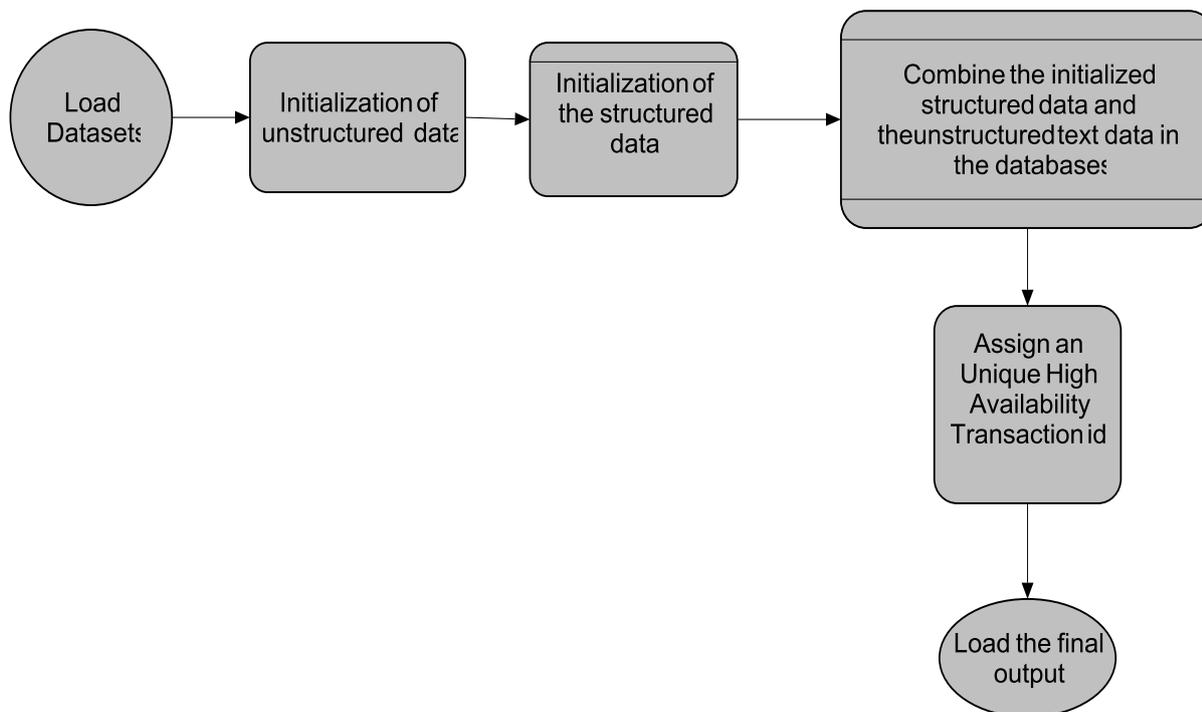
ready for the next phase of research where the data is fed for predictive analysis. Algorithm concentrates on the causes, sorts, and factors influencing DM (diabetes mellitus), preventive measures, and treatment of diabetes other than those directly associated with hypoglycaemia and serious metabolic unsettling influences.

**Proposed Architecture**

The proposed architecture employs R tool ecosystem. The proposed General Integrated High Availability Transaction [GIHAT] algorithm works at the third phase of the architecture and is presented in a R tool ecosystem. The detailed architecture is presented in the Figure 2.

This includes the process of collecting Structured and unstructured data about the diabetes Patients. Then the unstructured data is initialized by array initialization and then put in to the structured attributes. Then the database is loaded with all the structured data initialization and both the structured and unstructured datas are combined. Assign an Unique High Availability Transaction id for each patients which makes them easy to get their prescriptions.



**Fig 2. Proposed System Model**

**Algorithm Name: General Integrated High Availability Transaction [GIHAT] algorithm**

---

**Input :** Diabetic Patients structured and unstructured data-sets

Structured data [$STR_{data}$]

Unstructured Text data [$USTXT_{data}$]

**Output:** Processed data

---

// Load the database with all the unstructured text data

1. Database = [$USTXT_{data}$]

// Initialization of the unstructured text data [$USTXT_{data}$]

2. [ $USTXT_{data}$ ] + array initialization procedure [$^{init}$] =$USTXT^{init}_{data}$

// Assign a structure attribute to the initialized unstructured text data

3. for each $USTXT^{init}_{data}$

4. assign pre-defined array elements[$AR^{el}$] where each element is a structural property

5. $USTXT^{init}_{data}$ + $Ar^{el}$[1,2,3,4,5] = $Ar^{el}$ [$USTXT_{data}$]$^{(1:5)}$

// $Ar^{el}$ [$USTXT_{data}$]$^{(1:5)}$ is the new initialized unstructured text data

7. $USTXT^{init}_{data}$ = $Ar^{el}$ [$USTXT_{data}$]$^{(1:5)}$

**90**

// Load the database with all the structured data  [STR$_{data}$]

8. Database = [STR$_{data}$]

// Initialization of the structured data

9.  [STR$_{data}$]+ array initialization procedure [$^{init}$] =STR$^{init}_{data}$

// Combine the initialized structured data and theunstructured text data in the databases

10. Database =STR$^{init}_{data}$ + USTXT$^{init}_{data}$

// Database with structured elements are obtained

11. STR$^{init}_{data}$ + USTXT$^{init}_{data}$ = STRUC$^{[STRinitdata + USTXTinitdata]}$

11. Structured elements [STRUC$^{[STRinitdata + USTXTinitdata]}$] are arranged

12. **if** (STRUC$^{[STRinitdata + USTXTinitdata]}$ == initialized state) **do**

// Assign an Unique High Availability Transaction id

13. STRUC$^{[STRinitdata + USTXTinitdata]}$ = [STRUC$^{[STRinitdata + USTXTinitdata]}$ +UID$^{HAT}$]

// Load the final output to res_data$^{[complete]}$

14. STRUC$^{[STRinitdata + USTXTinitdata]}$= res_data$^{[complete]}$

15. else, Go to step 1

16. Restart the algorithm

Process of collecting Structured and unstructured data about the diabetes Patients. Then the unstructured data is initialized by array initialization and then put in to the structured attributes. Then the database is loaded with all the structured data initialization and both the structured and unstructured data are combined. Assign an Unique High Availability Transaction id for each patients which makes them easy to get their prescriptions.

## V.    RESULT AND DISCUSSION

A total of 1500 patients were looked over Out Patient Department of Periodontics, Government Dental College and Hospital, the Diabetic Clinic, Government Medical College and Hospital, and Diabetes Care and Research Center at Aurangabad. These patients were examined as having diabetes mellitus and were under treatment. The patients were picked by the going with thought criteria:

•       Under treatment or had diabetes mellitus broke down for at any rate latest one year or more.
•       Not having some other foundational ailments.
•       Not having any history of diabetic perplexities like neuropathy, nephropathy, and     retinopathy et cetera.
•       Not using medicines, for instance, phenytoin, nephidipine et cetera.
•       Not encountered any periodontal treatment since latest one year.

•       Readiness to share in the examination.

The vital history was recorded for each one of the patients. A wary oral examination was finished with the help of mouth reflect and graduated periodontal test.
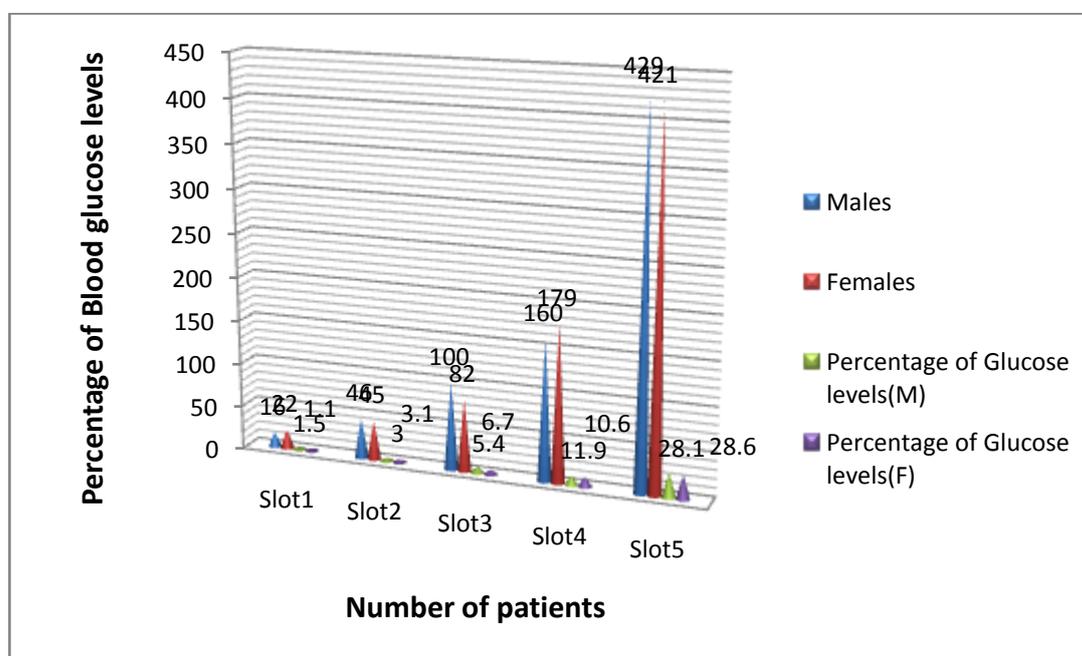
**Confirmation of blood glucose levels**
In each one of the patients, different blood was assembled under strict aseptic conditions, after an overnight snappy and one and half hour after dinner. The fasting and postprandial blood glucose levels were controlled by means of autoanalyzer.

As a result Of the 1500 patients, 3.4% of patients had insulin-subordinate diabetes mellitus (ISDM) and 96.6% had non-insulin-subordinate diabetes mellitus (NISDM). The assembled data was analyzed quantifiably General Integrated High Availability Transaction association coefficient examination was used to investigate the association among inescapability and reality of periodontal contamination and diverse elements, for instance, age, sex, glycemic status, and length of diabetes mellitus. Out of 1500 patients, 751 (50.1%) were male and 749 (49.9%) were female. The age scope of the patients was 15 years to 76 years with a mean age of 53.24±11.91 years. The patients were characterized into five gatherings as appeared in Table 1

**Table 2. Distribution of patients according to age and Gender**

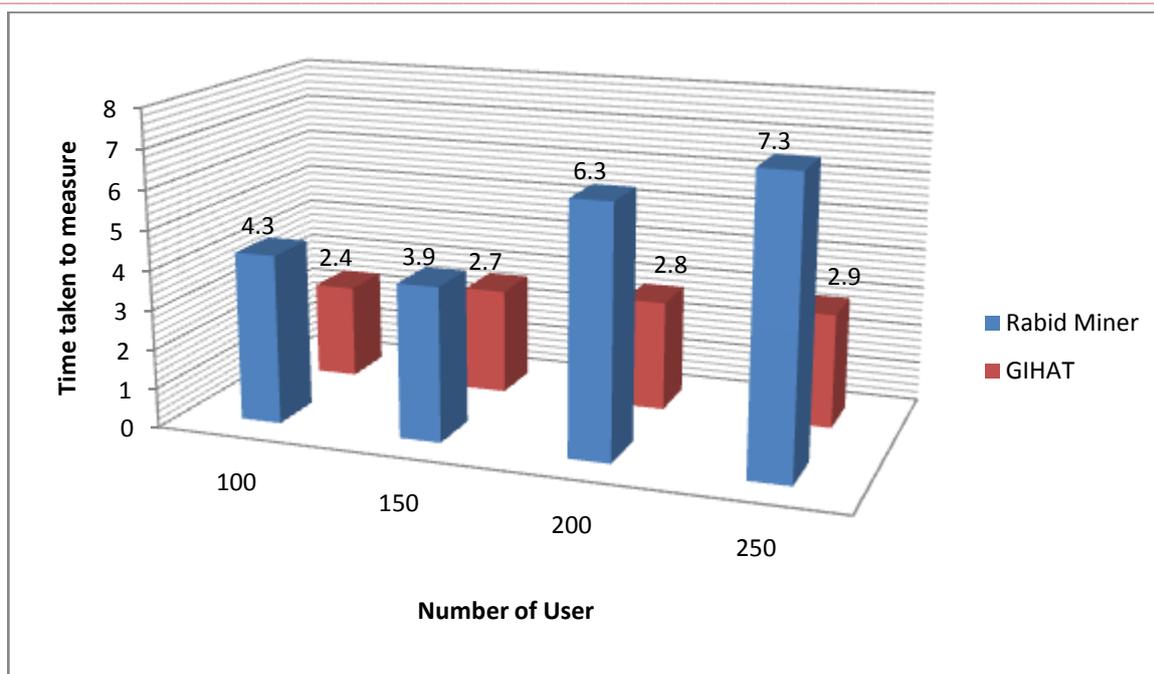| GROUP | AGE | NO.OF. PATIENTS | PERCENTAGE |
|---|---|---|---|
| **Males** | | | |
| **Slot I** | 15-24 | 16 | 1.1 |
| **Slot II** | 25-34 | 46 | 3.1 |
| **Slot III** | 35-44 | 100 | 6.7 |
| **Slot IV** | 45-54 | 160 | 10.6 |
| **Slot V** | 55& Above | 429 | 28.6 |
| **TOTAL** | | 751 | 50.1 |
| **Females** | | | |
| **Slot I** | 15-24 | 22 | 1.5 |
| **Slot II** | 25-34 | 45 | 3.0 |
| **Slot III** | 35-44 | 82 | 5.4 |
| **Slot IV** | 45-54 | 179 | 11.9 |
| **Slot V** | 55& Above | 421 | 28.1 |
| **TOTAL** | | 749 | 49.9 |



**Fig 3.No. of patients Vs Blood Glugose Levels**

**Efficiency**

The time taken to measure blood Glugose levels more efficient in General Integrated High Availability Transaction [GIHAT] algorithm, when comparing with rabid miner algorithm. Thus, the patient's feels better on identifying their issues or illness and it helps to prevent their diseases (Diabetes) as soon as possible

## VI. CONCLUSION

This paper presents a General Integrated High Availability Transaction [GIHAT] algorithm, The challenges in big data analysis on personal healthcare is Analysing of data sets which are in various form of structured and unstructured datasets. The data sets are then loaded in the database. Assign an Unique High Availability Transaction id for each patients which makes them easy to get their prescriptions. Thus this is more useful to the patients on comparing with other methods. Predict the onset of Diabetes mellitus at an early stage

### REFERENCE

[1] Jiawei Han and kamber "Data Mining Concepts and techniques".

[2] Tina R. Patil and Mrs. S. S. Sherekar"Performance Analysis of Naive Bayes and J48 Classification Algorithm for Data Classification" International Journal of Computer Science and Applications Vol. 6, No.2, Apr 2013 ISSN: 0974-1011 (Open Access).

[3] Raj Kumar, Rajesh Verma "Classification Rule Discovery for Diabetes Patients by Using Genetic Programming" "International Journal of Soft Computing and Engineering (IJSCE),ISSN: 2231- 2307, Volume-2, Issue-4, September 2012.

[4] Jahanvi Joshi and RinalDoshi Dr. Jigar Patel "Diagnosis and Prognosis Breast Cancer Using Classification Rules" International Journal of Engineering Research and General Science Volume 2, Issue 6, October-November, 2014,ISSN 2091-2730.

[5] Pyle, D. (1999) Data Preparation for Data Mining, San Francisco: Morgan Kaufmann.

[6] Cios, K. J., Pedrycz, W., Swiniarski, R.W., Kurgan, L. A. (2007) Data Mining: A Knowledge Discovery Approach, New York: Springer.

[7] Seibel, J. A. (2007) Diabetes Guide, WebMD, http://diabetes.webmd.com/guide/oral-glucose-tolerance-test.

[8] Stein, D. W. (2006) Hypertension / High Blood Pressure Guide, WebMD, http://www.webmd.com/ hypetensiondiagnosing- high-blood-pressure

[9] Zelman, K. M. (2008), How Accurate is Body Mass Index, or BMI? WebMD, http://www.webmd.com/diet/features/how-accurate-body-mass-index-bmi.

[10] Kass, G. V. (1980) An Exploratory Technique for Investigating Large Quantities of Categorical Data. Journal of Applied Statistics 29(2): 119-127.

[11] Quinlan, J. R. (1992) C4.5: Programs for Machine Learning, San Francisco: Morgan Kaufmann.

[12] Han, J., Kamber, M. (2006) Data Mining: Concepts and Techniques, 2nd ed. San Francisco: Morgan Kaufman

[13] Gaganjot Kaur and Amit Chhabra " Improved J48 Classification Algorithm for the Prediction of Diabetes" International Journal of Computer Applications (0975 – 8887) Volume 98 – No.22, July 2014.https://archive.ics.uci.edu/ml/datasets/Pima+Indians +Diabetes.

[14] Panigrahi Srikanth,D.Dharmaiah and Ch.Anusha "a computational intelligence technique for effective medical diagnosis using decision tree algorithm" imanager's Journal on Computer Science, Vol. 3 l No. 1 l March - May 2015.

[15] Dharmaiah Devarapalli, Allam Apparao, Amit Kumar, G R Sridhar "A Novel Analysis of Diabetes Mellitus by Using Expert System Based on Brain Derived Neurotrophic Factor" , Helix Vol. 1:251-256 (2013)

[16] Panigrahi srikanth,Dharmaiah Devarapalli "Identification of AIDS Disease using Genetic Algorithm "Springer, Briefs in Forensic and Medical Bioinformatics .ISBN:978-981-287-337- 82449

[17] "NO                SQL"                [Online]
Avilable:https://en.wikipedia.org/wiki/NoSQL.[Accessed:
16-Sept-20I6].

[18] A. B. M. Moniruzzaman and S. A. Hossain, "Nosql
database: New era of databases for big data analytics-
classification, characteristics and comparison," Nosql
database New era databases big data Analytics-
classification, Characteristics Comparison, vol. 6, no. 4,
Jun. 2013, pp. 1-13,.

[19] R. Hecht and S. Jablonski, "NoSQL evaluation: A use case
oriented survey," Proc. International Conference on Cloud
Service Computing CSC 2011,pp. 336-341, 2011.

[20] "Apache Cassandra Project. " [Online]. Available:
http://cassandra. apache.org. [Accessed: 16-Sept-2016]

[21] F. Chang,et aI. ,"Bigtable: A distributed storage system for
structured data," ACM Transactions on Computer Systems
(TOCS), vol. 26, no. 4,Jun. 2008, doi:I0.
114511365815.1365816.

[22] M. N. Vora,"Hadoop-HBase for large-scale data," Proc.
International Conference on Computer Science Network
Technology ICCSNT, vol. I,Dec. 2011,pp. 601-605.

[23] R. C. Taylor, "An overview of the
Hadoop/MapReduce/HBase framework and its current
applications in bioinformatics. ," BMC Bioinformatics,vol.
11, no. 12,2016.

[24] Sun and K. Chandan, "Big data analytics for healthcare,"
Proc. ACM. The 19th SIGKDD International Conference
on Knowledge Discovery and Data Mining,Aug. 2013, pp.
1525-1525.