# Implementing Association Rule Summarization for Predicting Relative Risk for Diabetes Mellitus

**Taslim N. Kureshi, Prof. Hemlata Dakhore**

Department of Computer Science and Engineering, G. H. Raisoni institute of Technology and Engineering, Nagpur, Maharashtra, India

## ABSTRACT

Diabetes is a developing pandemic of non-transmittable malady which influences the greater part of the general population on the planet. Keeping in mind the end goal to stifle the development of diabetes mellitus we utilize affiliation control rundown to electronic medicinal records to find set of hazard variables and the comparing sub-populace which speaks to patients at especially high danger of creating diabetes. Typically affiliation control mining creates huge volume of informational collections which we have to outline for any therapeutic record or any clinical utilize. We join four strategies to locate the basic components which prompt high danger of diabetes all these four techniques created synopses that depicted sub populaces at high danger of diabetes with every strategy having its unmistakable quality. As per our motivation we utilize bottom up summarization (BUS) calculation which delivers more appropriate rundown.

**Keywords :** Data Mining, Association Rule Mining, Survival Analysis, Association Rule Summarization

## I. INTRODUCTION

Diabetes mellitus is developing pestilence malady which influences more than 25.8 million individuals and around 7 million of them don't know they have this sickness. Normally diabetes is a gathering of ailments portrayed by high glucose (blood glucose). At the point when a man has diabetes the body either creates enough insulin or not able to utilize its own particular insulin successfully. At the point when glucose gets develop in our blood, that glucose ought to be controlled or should be adequately utilized else it might to lead passing. The danger of death of a man who has diabetes is twice as the individual who does not have diabetes of same age.

The real entanglements of diabetes are coronary illness and stroke. Grown-ups with diabetes have coronary illness demise rates around 2 to 4 times higher than grown-ups without diabetes the danger of stroke is 2 to 4 times higher among individuals with diabetes. It likewise prompts hypertension and 67% of diabetic patients have circulatory strain more noteworthy than or equivalent to 140/90 millimetres of mercury or

utilized professionally prescribed prescription for hypertension. Diabetes is a main source of visual deficiency among grown-ups matured 20-74 years. Around 60% 70% of individuals with diabetes have mellowed to extreme types of sensory system harm. The aftereffect of such harm incorporate disabled sensation or agony in the feet or hands demonstrated processing of sustenance in the stomach carpal passage disorder or other nerve issue.

Right around 30% of individuals with diabetes matured 40 years or more seasoned have hindered sensation in the feet. Diabetes may likewise prompt complexity amid pregnancy, ineffectively controlled diabetes before origination and amid the principal trimesters of pregnancy among ladies with sort 1.

Diabetes can bring about significant birth surrenders in 5% to 10% of pregnancy and unconstrained premature births in 15% to 20% of pregnancies. On other hand for a ladies prior diabetes enhancing blood glucose levels before and amid early pregnancy can be diminish the danger of birth deformities in their newborn children. Ineffectively controlled diabetes amid the second and

third trimesters of pregnancy can bring about unreasonably vast infants representing a hazard to both mother and kid.

Affiliation decides are suggestion that partner set possibly cooperating conditions (e.g.: high BMI and the nearness of hypertension conclusion). The utilization of affiliation tenets is especially useful on the grounds that notwithstanding measuring the diabetes hazard, the likewise promptly furnish the doctor with a "legitimization" in particular the related arrangement of conditions. These conditions can be utilized to guide treatment towards a more customized and focused on preventive care or diabetes administration.

Diabetes is a piece of the metabolic disorder, which is a heavenly body of illnesses including hyperlipidaemia (lifted triglyceride and low HDL levels), (hypertension) and focal weight (with body mass file surpassing 30 kg/m2). These ailments communicate with each other, with cardiovascular and vascular illnesses and in this manner comprehension and demonstrating these associations is critical. Affiliation standards are suggestions that partner an arrangement of conceivably collaborating conditions (e.g. high BMI and the nearness of hypertension analysis) with lifted hazard. The utilization of affiliation standards is especially advantageous, on the grounds that notwithstanding evaluating the diabetes chance, they additionally promptly give the doctor a "defense", specifically the related arrangement of conditions. This arrangement of conditions can be utilized to guide treatment towards a more customized and focused on preventive care or diabetes administration. While affiliation rules themselves can be effectively deciphered, the subsequent manage sets can some of the time be substantial, disintegrating the translate capacity of the lead set in general.

Particularly, in this work, we consider a rich arrangement of hazard components, in particular co-sullen ailments, lab results, drugs and statistic data that are usually accessible in electronic therapeutic record (EMR) frameworks. With such a broad arrangement of hazard variables, the arrangement of found standards becomes combinatorial expansive, to a size that extremely prevents translation. To beat this test, we connected manage set summarization methods to pack the first lead set into a more conservative set that can be translated effortlessly. Various effective affiliation

govern set synopsis strategies have been proposed [10] yet no unmistakable direction exists with respect to the materialness, qualities and shortcomings of these procedures. The concentration of this original copy is to audit and describe four existing affiliation lead synopsis systems and give direction to specialists in picking the most reasonable one. A typical deficiency of these systems is their failure to take diabetes risk–a nonstop outcome–into account. Keeping in mind the end goal to make these methods more suitable, we needed to negligibly alter them: we stretch out them to consolidate data about ceaseless result factors.

In particular, our key commitments are as per the following.

- We introduce a clinical utilization of affiliation lead mining to distinguish sets of co-bleak conditions (and the patient sub populaces who experience the ill effects of these conditions) that suggest altogether expanded danger of diabetes.
- Association control mining on this broad arrangement of factors brought about an exponentially expansive arrangement of affiliation standards. We developed four prevalent affiliation run set synopsis systems (basically from the audit [10]) by fusing the danger of diabetes into the way toward finding an ideal rundown.
- Our fundamental commitment is a relative assessment of these amplified outline systems that gives direction to experts in choosing a fitting calculation for a comparable issue.

## II. RELATED WORKS

A diabetes list is fundamentally a prescient model that allocates a score to a patient in view of his evaluated hazard of diabetes. Collins et al. [7] directed a broad review of diabetes records depicting the hazard variables and the displaying system that these files used. They found that most files were added substance in nature and none of the overviewed files have considered cooperation among the hazard elements.

While we don't know about any new diabetes file distributed after the overview, a current review [12] concentrating on the metabolic disorder (of which diabetes is a part) speaks to a critical advancement. Kim et al. utilized affiliation run mining to deliberately investigate co-occurrences of conclusion codes. The

subsequent affiliation rules don't constitute a diabetes file on the grounds that the review does not assign a specific result of intrigue and they don't evaluate or anticipate the danger of diabetes in patients, however they found some huge relationship between finding codes.

We have as of late attempted a diabetes think about [4] where we expected to find the connections among sicknesses in the metabolic disorder. We utilized an indistinguishable companion from this present review; in any case, we included just eight finding codes and age as predictors. We found affiliation rules including some of these eight finding codes, surveyed the danger of diabetes that these principles give on patients and introduced the tenets as a movement diagram delineating how patients advance from a solid state towards diabetes. We exhibited that the approach discovered clinically important affiliation decides that are steady with our restorative desire.

With just eight indicator factors, the extent of the found govern set was modest–13 noteworthy rules– and thusly, elucidation was direct. Actually, no govern set synopsis was fundamental.

## III. LITERATURE SURVEY

Chaudhari et al [13] Disease assurance is a champion among the most basic employments of such system as it is one of the principle wellsprings of passing wherever all through the world. Predict the human use the commitments from complex tests coordinated in labs moreover expect the disease considering risk components, for instance, tobacco smoking, alcohol confirmation, age, family history, diabetes, hypertension, raised cholesterol, physical torpidity, weight. Investigators have been using a couple data mining strategies to help restorative administrations specialists in the examination of coronary disease. K-Nearest-Neighbor (KNN) is one of the viable data mining techniques used as a piece of request issues. Starting late, investigators are showing that uniting particular classifiers through voting is defeating other single classifiers. This paper inquires about applying KNN to help human administrations specialists in the finish of disease exceptionally coronary ailment. It moreover inquires about if planning voting with KNN can redesign its precision in the assurance of coronary disease patients. The results exhibit that applying KNN could achieve higher accuracy than neural framework assembling in the finding of coronary disease patients. The results furthermore show that applying voting couldn't enhance the KNN precision in the assurance of coronary sickness.

Prof.Mythili et al [12] Diabetes mellitus, in essential terms called as diabetes, is a metabolic sickness, where a man is affected with high blood glucose level. Diabetes is a metabolic issue brought on due to the mistake of body to make insulin or to properly utilize insulin. This condition rises when the body does not make enough insulin, or in light of the way that the cells don't respond to the insulin that is conveyed. Blood glucose test is the imperative system for diagnosing diabetes. Also, there have been different robotized methodologies proposed for finish of diabetes.

Each one of these procedures has some data values which would be the result of different tests that should be finished in recuperating focuses. This paper proposes a framework those arrangements to encourage the patients encountering diverse remedial tests, which by far most of them consider as a dull undertaking and repetitive.

The parameters recognized for diagnosing diabetes have been created in a way that, the customer can expect if he is impacted with diabetes himself. Back Propagation count is used for conclusion.

Ahmed et al [15] Heart illness is an imperative purpose behind grimness and mortality in cutting edge society. Therapeutic conclusion is basic yet ensnared undertaking that should be performed decisively and adequately. The competent data examination instruments are used to remove accommodating gaining from the immense measure of restorative data. There is huge data open inside the therapeutic administrations structures. In any case, there is a task of convincing examination gadgets to discover covered associations and examples in data. Learning exposure and data mining have found various applications in business and trial space.

One of the applications is sickness finding where data mining instruments are exhibiting productive outcomes. This investigation paper proposed to find the heart disorders through data mining, Support Vector

Machine (SVM), Genetic Algorithm, unsavoury set speculation, connection rules and Neural Networks. In this review, we immediately assessed that out of the above systems Decision tree and SVM is best for the coronary sickness. So it is watched that, the data mining could help in the recognizing verification or the desire of high or by and large safe heart disorders.

Thangaraju et al [16] Data mining is the demonstration of taking a gander at enormous earlier databases with a particular true objective to make new data. There are different sorts of data mining strategies are available. Course of action, Clustering, Association Rule and Neural Network are likely the hugest frameworks in data mining. In Health mind organizations, Data mining accept an imperative part. Most a great part of the time the data mining is used as a piece of human administrations undertakings for the route toward foreseeing diseases. Diabetes is an unending condition. This suggests is continues for a long time, consistently for some individual's whole life [11]. This paper considers the examination of diabetes gaging approaches using gathering frameworks. Here we are using three different sorts of batching frameworks named as Hierarchical gathering; Density based gathering, and Simple K-Means grouping. Weka is used as a gadget.

Durairaj et al [17] Neural Networks are one of the sensitive enlisting strategies that can be used to make desires on restorative data. Neural Networks are known as the Universal pointers. Diabetes mellitus or essentially diabetes is a contamination achieved due to the extension level of blood glucose.

Distinctive traditional procedures, considering physical and invention tests, are available for diagnosing diabetes. The Artificial Neural Networks (ANNs) based system can sufficiently associate for hypertension danger desire. This improved model segregates the dataset into both of the two social events. The earlier revelation using sensitive enrolling methods help the specialists to decrease the probability of quitting any and all funny business of the disease. The data set chose for request and exploratory entertainment relies on upon Pima Indian Diabetic Set from (UCI) Repository of Machine Learning databases. In this paper, a separated audit is driven on the utilization of different fragile preparing frameworks for the desire of diabetes. This survey is relied upon to perceive and propose a fruitful method for earlier estimate of the illness.

## IV. IMPLEMENTATION

We attempt to utilize association rule mining to the electronic therapeutic record (EMR); All the hazard calculate about a patient in particular co-sullen malady and research facility results and solutions are being accessible in the EMR, there are less opportunities to miss insights about a patient with the broad arrangement of hazard variables the arrangement of found hazard turns out to be amazingly vast to conquer this we utilize rule set outline method which is utilized to pack the first rule set into a conservative set. We utilize the accompanying systems

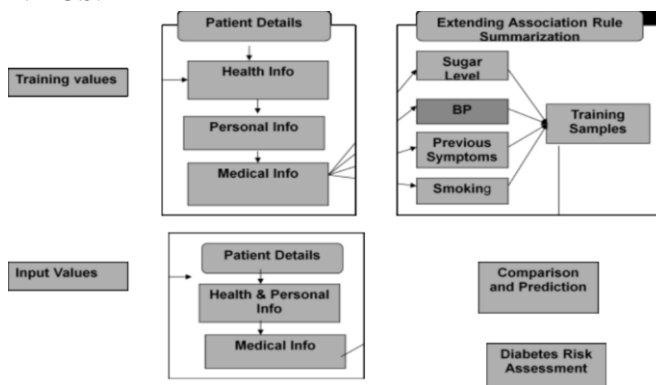1. APRX-collection
2. RPG-global
3. TOPK
4. BUS.



**Figure 1:** System Architecture

We first check the support of individual things and figure out which of them extensive (ie.) we have least support are. In each pass we begin with a seed set of things observed to be expansive in the past pass. We utilize this seed set for creating new possibly substantial thing sets called competitor key, thing set and number the genuine support for these applicant thing set amid the disregard the information.

Toward the finish of the pass we figure out which hopeful thing set are in reality expansive and they progress toward becoming seed for not pass. This procedure contains until no new substantial thing set are found. Testing measurable criticalness: for each found thing set we need to test whether the result

dispersion in the influenced and unaffected subpopulation is without a doubt distinctive.

Step-2 the arrangement of thing set is separated so that lone the factually critical ones are returned as distributional association rule, this rule is portrayed by the accompanying insights from the quantity of thing set gathered. Let OR be the watched number of diabetes episode in the subpopulation DR secured by R. give ER a chance to signify the normal number of diabetes occurrences in the subpopulation secured by R. ER = OR-i$\varepsilon$DRyi where yi is the martingale for patient. The relative risk factor is defined by R

RR = OR/ER.

| Parameter | Weightage | Values |
|---|---|---|
| Male & Female | Age<30 >30to<50 | 0.1 0.3 0.7 0.8 |
| Smoking | Never Past Current | 0.1 0.3 0.6 |
| Overweight | Yes No | 0.8 0.1 |
| Alcohol intake | Never Past Current | 0.1 0.3 0.6 |
| Heart rate | Low(<60 bpm) Normal(60 to 100bpm) High(>100bpm) | 0.9 0.1 0.9 |
| Blood sugar | High(>120&<400) Normal(>90&<120) Low(<90) | 0.5 0.1 0.4 |
| Bad cholesterol | Very high>200 High(160 to 200) Normal<160 | 0.9 0.8 0.1 |

Table 1: Description of the risk factors that appeared in any of the summarized rules

When we attempt to apply distributional rule mining with our electronic medicinal records it created an expansive number of (measurably critical) rules. Rules that were produced marginally vary from each other prompting muddling of clinical patters. With a specific end goal to conquer the issue of this extensive number of rules which were produced we go for abridging the rule set into littler set for our less demanding review. We first survey the current rule set and database outline techniques then we attempt to fuse a nonexclusive structure with a specific end goal to get a persistent result of variable into record.

Presently we show the rule set produced by the expanded outline algorithm, for every algorithm we utilized the parameter setting that gave the best outcomes to APRXcollection we utilized $\alpha = 0.1$, $\lambda = 1$ for RPG worldwide we utilized $\delta = 0.5$, $\sigma = 0.2$, $\lambda = 0.98$ for top K we utilized $\lambda = 0.2$ and for BUS we utilized $\lambda = 0.1$

A. APRX-collection

The APRX collection algorithm is utilized to discover the supersets of the condition (hazard consider) in the rule with the end goal that most subsets of rundown rule will be legitimate rules in the first (un condensed) set and these subset rules suggest comparable hazard for diabetes.

| R | RR | ER | OR | RULE |
|---|---|---|---|---|
| 1 | 1.96 | 36.24 | 71 | Fibra |
| 20 | 1.34 | 271.71 | 363 | Bmi trigal acerab Statin aspirin htn |
| 16 | 1.19 | 426.78 | 506 | Hdl trigl acearb Aspirin htn |
| 15 | 1.31 | 348.92 | 457 | Bmi trigal statin aspirin ihd |
| 10 | 1.23 | 534.58 | 660 | Bmi sbp ccb htn |

Table 2: Rule set summarized by APRX- collection

The APRX collection focuses just on articulation of the rule subsequently it needs data about which patients are as of now secured accordingly patients can get secured by different rules prompting rules with fundamentally the same as condition this strategy additionally needs in accuracy and data about high hazard subgroups.

B. RPG-global

The principle downsides of APRX collection were the excess in the rule set and the weakening of the hazard. The RPGlobal synopsis is like APRX collection n in that it is mostly worried with the outflow of the rule and thus it plays out an extremely forceful pressure. RPG Global has two downsides by considering Patient scope and by building the outline from rules in the first rule set.

| RR | ER | OR | RULE |
|---|---|---|---|
| 1.69 | 32 | 55 | Bmi trigal acearb diuret htn |
| 1.23 | 52 | 65 | Acearb bb diuret aspirin htn |
| 1.29 | 42 | 55 | Sbp tchol acearb diuret htn |
| 2.10 | 25 | 54 | Hdl trigal diuret aspirin htn |
| 1.28 | 42 | 54 | Bmi tchol hdl trigl tobacco |

Table 3: Rule set created by RPGlobal.

## C. TOPK

Top-K algorithm decreases the repetition in the rule set which was conceivable through working on patients instead of the declaration of the rules. This approach relinquished the extraordinary pressure rates of past two algorithm TOP-K still accomplishes high pressure rate and it effectively distinguished rules with high hazard and low excess.

| RR | ER | OR | RULE |
|---|---|---|---|
| 2.40 | 21.70 | 52 | Fibra htn |
| 1.58 | 37.97 | 60 | Bmi hdl ihd |
| 1.47 | 45.52 | 67 | Sbp htn tobacoo |
| 1.46 | 317.03 | 464 | Bmi htn |
| 1.62 | 32.16 | 52 | Sbp tchol trigal statin htn |

Table 4: Rule created by the top-k algorithm

## D. BUS

The outlines made by BUS (appeared in Table 5) and Top-K are comparable in quality. The BUS outline shows less changeability (it tends to utilize similar conditions: bmi and trigl co-happen in 40% of the rules), yet this diminished fluctuation does not convert into expanded excess in the patient space.

BUS (rather than TopK) works on the patients and not on the rules. Accordingly, excess as far as rule expression can happen. Nonetheless, BUS unequivocally controls the excess in the patient space through the parameter commanding the base number of new (already revealed) cases (patients with diabetes

occurrence) that should be secured by each rule. In this manner the diminished changeability in the rule expression does not convert into expanded repetition.

| RR | ER | OR | RULE |
|---|---|---|---|
| 2.34 | 24 | 57 | Bmi trigal acearb statin htn |
| 2.10 | 25 | 54 | Hdl trigal diuret aspirin htn |
| 1.91 | 56 | 107 | Bmi trigal statin htn |
| 1.54 | 78 | 121 | Bmi trigal tobacco |
| 1.37 | 39 | 54 | Dbp diuret htn |

Table 5: Top 10 summarized rule created by BUS

The BUS rule set figured out how to incorporate coronary illness prior (rule #3) and with higher hazard (2.15) than the Top-K rule set. Likewise, BUS involves tobacco use in a blend of hazard variables with higher relative hazard than Top-K. Generally speaking, regardless of the distinctions in the algorithms, BUS and Top-K produce comparable fantastic rundowns.

## V. EXPERIMENTAL RESULT

The quantity of rules should have been decreased to a level where clinical elucidation is doable. To this end, we concentrated four techniques to compress these rules into sets of 10-20 rules that clinical examiners can assess.

While every one of the four strategies made sensible synopses, every strategy had its unmistakable quality. In any case, not these qualities are essentially gainful to our application. We found that the most imperative differentiator between the algorithms is whether they utilize a determination model to incorporate a rule in the outline in view of the outflow of the rule or in light of the patient sub populace that the rule covers.

APRX-COLLECTION and RPGlobal basically work on the outflow of the rules with an essential target of expanding pressure. They utilize delegate rules, each of which speaks to various unique rules. Such illustrative rules accomplish high pressure, however weaken the danger of diabetes over the regularly huge subpopulation they cover. TopK and BUS work basically on the patients and their objective–especially

in the event of TopK–can be considered as limiting excess. They delivered great outlines in light of the fact that a helpful reaction of diminishing repetition is to accomplish great pressure. The opposite is not valid: high pressure rate does not bring about low repetition. Amongst TopK and BUS, we found that BUS held marginally more repetition than TopK, which permitted it to have better patient scope and better capacity to remake the first information base. This favorable position made BUS the most appropriate algorithm for our motivation.
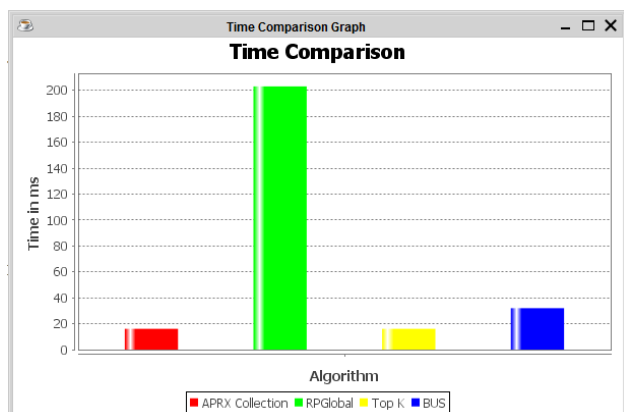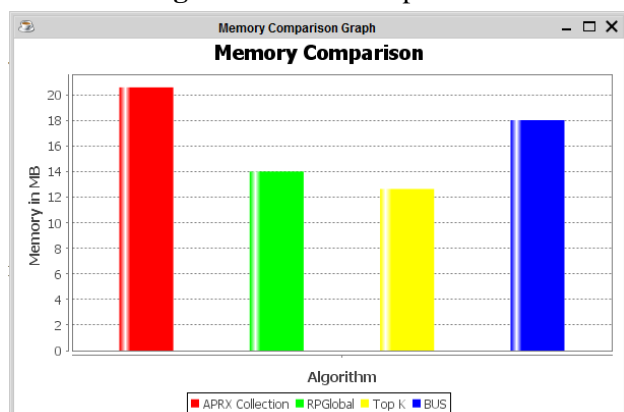


**Figure 2:** Time Comparison



**Figure 3:** Memory Comparison

Our Result also shows the Time and space utilization of the above four algorithms. Figure show the results generated for the same.

## VI. CONCLUSIONS

The electronic information produced by the utilization of EMRs in routine clinical practice can possibly encourage the revelation of new learning. Affiliation control mining coupled to a synopsis system gives a basic device to clinical research. It can reveal concealed clinical connections and can propose new examples of conditions to divert counteractive action,

administration, and treatment approaches. In our particular illustration, we utilized distributional affiliation control mining to distinguish sets of hazard elements and the relating understanding subpopulations that are at altogether expanded danger of advancing to diabetes. An over the top numbers of association rules were found hindering the clinical elucidation of the outcomes.

## VII. REFERENCES

[1]. F. Afrati, A. Gionis, and H. Mannila, "Approximating a collection of frequent sets," in Proc. ACM Int. Conf. KDD, Washington, DC, USA, 2004.

[2]. R. Agrawal and R. Srikant, "Fast algorithms for mining association rules," in Proc. 20th VLDB, Santiago, Chile, 1994.

[3]. Y. Aumann and Y. Lindell, "A statistical theory for quantitative association rules," in Proc. 5th KDD, New York, NY, USA, 1999.

[4]. P. J. Caraballo, M. R. Castro, S. S. Cha, P. W. Li, and G. J. Simon, "Use of association rule mining to assess diabetes risk in patients with impared fasting glucose," in Proc. AMIA Annu. Symp., 2011.

[5]. Centers for Disease Control and Prevention. "National diabetes fact sheet: National estimates and general information on diabetes and prediabetes in the United States," U.S. Department of Health and Human Services, Centers for Disease Control and Prevention, 2011 [Online].

[6]. V. Chandola and V. Kumar, "Summarization – Compressing data into an informative representation," Knowl. Inform. Syst., vol. 12, no. 3, pp. 355–378, 2006.

[7]. G. S Collins, S. Mallett, O. Omar, and L.-M. Yu, "Developing risk prediction models for type 2 diabetes: A systematic review of methodology and reporting," BMC Med., 9:103, Sept. 2011.

[8]. Diabetes Prevention Program Research Group, "Reduction in the incidence of type 2 diabetes with lifestyle intervention or metformin," N. Engl. J. Med., vol. 346, no. 6, pp. 393–403, Feb. 2002.

[9]. G. Fang et al., "High-order SNP combinations associated with complex diseases: Efficient discovery, statistical power and functional interactions," PLoS ONE, vol. 7, no. 4, Article e33531, 2012.

[10]. M. A. Hasan, "Summarization in pattern mining," in Encyclopedia of Data Warehousing and Mining, 2nd ed. Hershey, PA, USA: Information Science Reference, 2008.

[11]. Rakesh Motka, Viral Parmar, Balbindra Kumar, A. R. Verma, " Diabetes Mellitus Forecast Using Different Data Mining Techniques", International conference on computer and Communication Technology

[12]. Prof.Sumathy, Prof.Mythili, Dr.Praveen Kumar, Jishnujit T M, K Ranjith Kumar, "Diagnosis of Diabetes Mellitus based on Risk Factors", International Journal of Computer Applications, Vol.10, Issue No.4, November.2010

[13]. Anand A. Chaudhari, Prof.S.P.Akarte, " Fuzzy and Data Mining based Disease Predection using K-NN Algorithm", International Journal of Innovations in Engineering and Technology, Vol. 3, Issue No. 4, April 2014

[14]. Prof.Sumathy, Prof.Mythili, Dr.Praveen Kumar, Jishnujit T M, K Ranjith Kumar, " Diagnosis of Diabetes Mellitus based on Risk Factors", International Journal of Computer Applications, Vol. 10, Issue No. 4, November 2010

[15]. Aqueel Ahmed, Shaikh Abdul Hannan, " Data Mining Techniques to Find Out Heart Diseases: An Overview", International Journal of Innovative Technology and Exploring Engineering, Vol. 1, Issue No. 4, September 2012

[16]. P. Thangaraju, B.Deepa, T.Karthikeyan, "Comparison of Data mining Techniques for Forecasting Diabetes Mellitus", International Journal of Advanced Research in Computer and Communication Engineering, Vol. 3, Issue No. 8, August 2014

[17]. M. Durairaj, G. Kalaiselvi, " Prediction Of Diabetes Using Soft Computing Techniques- A Survey", International Journal of Scientific & Technology Research, Vol. 4, Issue No.3, March 2015

[18]. S.F.B, Jaafar and Darmawaty Mohd Ali. "Diabetes Mellitus Forecast using Artificial Neural Network (ANN), Asian conference on sensors and the international conference on new techniques in pharmaceutical and medical research proceedings (IEEE), Kuala Lumpur, Malaysia, 5-7 September 2005, pp 135-139.

[19]. S. Alby, B. L. Shivakumar," A survey on data-mining technologies for prediction and diagnosis of diabetes", International conference of IEEE 2014.

[20]. Gyorgy J. Simon, Terry M. Therneau, Steven S. Cha, " Extending association rule summarization techniques to assess risk of diabetes mellitus", IEEE VOL 27, no. 1, January 2015