

## Improved Health Record Mining using Supervised Machine Learning with Recommendation

K.Devipriya<sup>1</sup>, V.priya<sup>2</sup>, K.Dhineshkumar<sup>3</sup>

<sup>1</sup>Scholar, PG and Research Department of Computer Science, Nehru Memorial College (Autonomous), Puthanampatti, INDIA

<sup>2</sup>Assistant Professor PG and Research Department of Computer Science, Nehru Memorial College (Autonomous), Puthanampatti, INDIA

<sup>3</sup>Assistnt Professor, Agni College of Technology, Chennai, INDIA

### ABSTRACT

In medical dataset retrieving of particular health record and classification of patient at risk in earlier stages by data mining technique is widely utilized. Finding out hidden patterns from unlabeled data is a difficult and crucial process. In this thesis a classification algorithm is implemented to classify the records of patient from a huge database and by CADs exact retrieval is attained. Insertion and query are the two methods available in CADs. Patient details are added through insertion and by means of query extraction of specific data is done. It accurately extracts features. Based on this, an attempt to implement an efficient health record classification and utilized it to improve the performance of health classification.. CADs algorithm is used to extract details of a patient and identifies whether the patient is in critical, normal or medium condition. If the patient is going to suffer from a disease it will be attained as result from this algorithm. In addition to this our approach suggests the food to a particular user is done by collaborative filtering is taken further in order to safeguard themselves from upcoming disease problems.

**Keywords--** Data Mining Techniques, Supervised Machine Learning, CADs algorithm and Risk Factor

### I. INTRODUCTION

Leveraging temporal observations to predict a patient's health state at a future period is a very challenging task. The problem is formulated as an optimization based binary classification task addressed in three steps.[1] First, the time series data is transformed into a binary matrix representation suitable for application of classification methods. Second, an oval convex-concave optimization problem is defined to extract multivariate patterns from the constructed binary matrix. Then, a mixed integer discrete optimization formulation is provided to reduce the dimensionality and extract interpretable multivariate patterns. Finally, those interpretable multivariate patterns are used for early

classification in challenging clinical applications. [2-3]Predicting or prioritizing the human genes that cause disease, or "disease genes", is one of the emerging tasks in biomedicine informatics. Research on network-based approach to this problem is carried out upon the key assumption of "the network-neighbor of a disease gene is likely to cause the same or a similar disease", and mostly employs data regarding well-known disease genes, using supervised learning methods. This work aims to find an effective method to exploit the disease gene neighborhood and the integration of several useful data sources, which potentially enhance disease gene predictions.[3-4]Active learning has been extensively studied and successfully applied to solve real problems. The typical setting of active learning methods is to query absolute questions. In a medical application where the goal is to predict the risk of patients on certain disease using Electronic Health Records (EHR),the absolute questions take the form of "Will this patient suffer from Alzheimer's later in his/her life?", or "Are these two patients similar or not?".[5-6]. It alternatively focuses on designing *relative* questions that can be easily answered by domain experts. The proposed relative queries take the form of "Is patient A or patient B more similar to patient C?", which can be answered by medical experts with more confidence. These questions poll relative information as opposed to absolute information, and even can be answered by non-experts in some cases. [6-7]The surprising empirical results on real-world medical problems demonstrate the usefulness of our ARP method, as querying for the relative similarities of patients can achieve comparable and in some cases even better prediction performance than querying absolute questions on patients, while the latter type of questions is significantly more difficult to answer. Cancer classification is the critical basis for patient-tailored therapy.[7-8]. Conventional histological analysis tends to be unreliable because different tumors may have similar appearance. In the training stage, the method performs cross validation to learn optimal amount of shrinkage,

and hence, to sift genes for classification [9-10]. The authors introduced a family of simple methods, TSP and k-TSP, for inducing decision rules based on the top or k-top scoring pairs of genes. The approaches are based on the concept of relative expression reversals of gene pairs. Some pairs of genes are chosen for classification if the difference of their expression levels can correctly classify the samples of different types. K-TSP is, in fact, an enhanced ensemble version of TSP [10-11]. A new approach is introduced for classifying different types of cancers based on the rough sets theory. By dynamically constructing implicit hyper cuboids, the approach selects potential functional genes for inducing classifiers. Experimental results show that the induced classifiers are capable of classifying cancers with high accuracy and ROC, while only a small number of genes are involved. The results suggest that the proposed method is a feasible way of classifying different cancer types in applications. [11-12] Future efforts can be devoted to the enhancement of the method for Objective, To determine whether cognitive impairment assessed at annual geriatric health examinations is associated with increased mortality in the elderly. Method, This cohort study was based on data obtained from the government-sponsored Annual Geriatric Health Examination Program for the elderly in Taipei City between 2006 and 2010. The study sample consisted of 77,541 community-dwelling Taipei citizens aged 65 years older. [13-14]

The Short Portable Mental Status Questionnaire (SPMSQ) was selected to measure cognitive impairment. Mortality was ascertained by matching cohort IDs with national death files. Results: There was a dose-response relationship between cognitive impairment and mortality (increased one score of SPMSQ, Hazard ratio [HR]: 1.12, 95% confidence interval [CI]: 1.10–1.14). [15-16] Relative to no cognitive impairment, the HRs were 1.67 (95% CI: 1.43–1.94), 2.26 (95% CI: 1.90–2.70), and 2.68 (95% CI: 2.25–3.19) for mild, moderate, and severe cognitive impairments, respectively. The causes of death associated with cognitive impairment were circulatory, respiratory, and other causes, but not death from cancer [17].

## II. NEED AND NECESSITY OF DATA MINING AND RECORDING

Selective information delegacy is any fact, numbers, house painting or text variant that can be processed and manipulated by an electronic computer.

Data such as, gross revenue, fact, single value or set of values, cost, inventory, payroll department, health data and accounting, industry sales, forecast information, and macroeconomic data, logical database design. The patterns, associations, or relationship among all this data can provide information.

For example, analysis of retail point of sale operation data can yield information on which products are selling.

### Knowledge

Information can be converted into cognition about historical patterns and time to come trends. For example, sum-up information on retail supermarket sales can be analyzed in light of promotional efforts to provide knowledge of consumer buying demeanor. Thus, a maker or retailer could determine which tokens are most susceptible to promotional efforts. Dramatic forward motion in information capture, cognitive operation powerfulness, information transmission, and storage capabilities are enabling establishment to integrate their various database into data warehouse. Data reposition is defined as a process of centralized data management and retrieval.

### Data Warehouse

Data warehousing, like data mining, is a relatively new term although the construct itself has been around for years. Data warehousing represents an ideal visual sense of maintaining a central secretary of all organizational data. Centralization of data is needed to maximize user access and depth psychology. Dramatic technological advances are making this vision a world for many companies. And, equally dramatic advances in data analysis software package are allowing users to access this data freely. The data analysis software is what supports data mining.

### What Is Data Mining

Usually, data mining (every now and then referred to as records or understanding discovery) is the method of analyzing statistics from exclusive views and summarizing it into beneficial records - records that may be used to increase revenue, cut charges, or both. Data mining software program is one in all some of analytical equipment for reading facts.

It permits users to analyze data from many distinctive magnitude or angles, catalog it, and summarize the relationships identified. Technically, data mining is the technique of finding correlations or patterns among dozens of fields in massive relational databases.

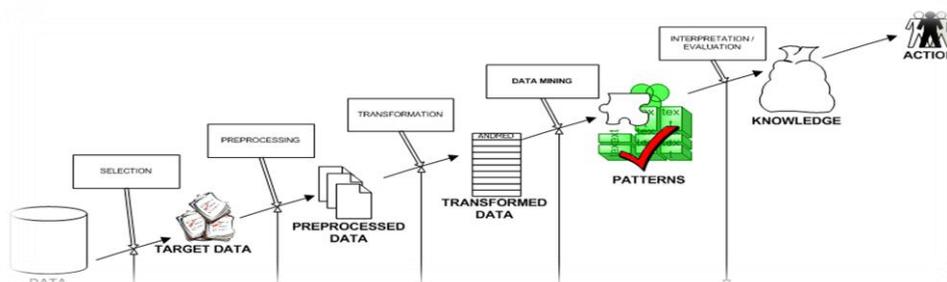


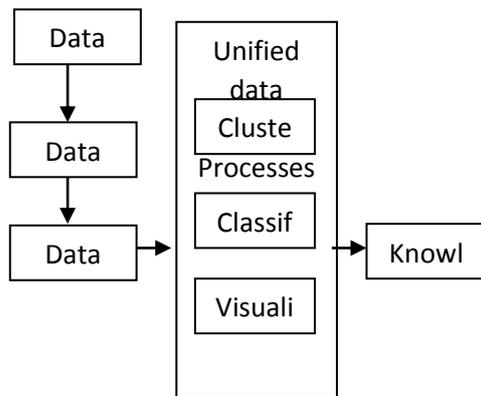
Fig 1, Structure of Data Mining

### How Data Mining Works?

Whilst large-scale data technology has been evolving separate transaction and analytical systems, information mining affords the hyperlink among the two. Data mining software program analyzes relationships and patterns in saved transaction information primarily based on open-ended person queries. Numerous kinds of analytical software program are to be had: statistical, gadget studying, and neural networks.

#### WORKING OF DATA MINING

Data collected from various informants are gathered together and then preprocessing is done. The result of preprocessing is data curing . This data set is consider as exam set and data mining process is performed from this knowledge will be extracted.



Several types of analytical software are available: statistical, machine learning, and neural networks. Generally, any of four types of relationships are sought.

**Classes:** Stored data is used to find data in predestined groups. For example, an hotel chain could mine customer purchase data to determine when customers visit and what they typically decree. This information could be used to augment traffic by having daily specials

**Clusters:** Information items are grouped according to logical relationships or consumer preferences. For example, data can be mined to identify market segment or consumer kinship.

**Associations:** Information can be mined to identify tie. The beer -diaper example is an example of associative mining.

**Sequential normal:** Data is mined to predict demeanor normal and manner. For example, an outdoor equipment retailer could predict the likelihood of a backpack being leverage d based on a consumer's purchase of sleeping old grip and hiking shoes.

#### Elements of Data Mining

Extract, transform, and load transaction information onto the data warehouse organization Depot and manage the data in a multidimensional database system.

Provide data access to stage business analysts and in data formatting ion technology professionals. Analyze the data by application software.

Present tense the data in a useful format, such as a graphical record or mesa

Data Mining Consists Of Five Major Elements

Extract, transform, and load transaction information onto the data warehouse organization

Depot and manage the data in a multidimensional database system.

Provide data access to stage business analysts and in data formatting ion technology professionals.

Analyze the data by application software.

Present tense the data in a useful format, such as a graphical record or mesa

#### Data Mining Functionalities

There are two methods on Data mining one was Predictive Method and Descriptive Method.

**Classification:** Predicts an item class

**Clustering:** Finds Clusters in data

**Association:** Frequently occurring events

**Summarization:** Describes a Group

Assortment is one of the most researched inquiries in motor car encyclopedism and data mining. In machine eruditeness, assortment refers to an algorithmic process for designating a given stimulation data into one among the different categories given. A wide compass of real problems has been stated as Classification Problems, for example mention scoring, stock market prediction, medical diagnosis, convention credit, text categorization and many more.

An algorithm that implements classification is known as a classifier. The input data can be termed as an instance and the categories are known as stratum. The feature of the instance can be described by a vector of features. These features can be nominal phrase, ordinal, integer -valued or real-valued

The Classification is a supervised procedure that learns to classify new illustration based on the knowledge learnt from a previously classified ad training set of instances. This piece of work has been carried out to make a execution evaluation of Machine Learning Algorithmic rule: CAD. The classification of information mining is probably its most developed percentage it has the greatest voltage pay-off and the most precise description.

In data mining, the choice of algorithmic program to use in analyzing a data set depends on the understanding of the psychoanalyst. In most cases, a lot of time is wasted in trying every single prediction proficiency (bagging, boosting, stacking, and meta-learning) in a dictation to breakthrough the best solution those paroxysm the analyst's needs. Hence, with the arrival of improved and modified prediction techniques, there is a need for an analyst to know which prick performs best for a particular type of data set.

In this thesis java coding and Weak 3.9 is used on three uniquely different data sets to compare the classification abilities of each of the classification algorithm on these different data samples. The thesis also deals with some of the data preprocessing techniques that will help to reveal the nature of the data

sets, with the aim of appropriately using the right classification algorithm in making classification.

The advantages and disadvantages of this classification are discussed also. Hence, this study will be helpful to learners and experts alike as they choose the best approach to solving basic data-mining problems. This will help in reducing the lead time for getting the best classification possible.

#### **DATA ANALYSIS:**

Data analysis is a method wherein uncooked statistics is prepared and structured so that treasured information can be extracted from it. The procedure of organizing and considering data is manner to accepting what the records does and does not incorporate. There is an expansion of methods wherein public can method information analysis and it is notoriously smooth to direct information during the analysis section to push positive conclusions or agendas. Evaluation of data is a method of examining, cleansing, remodeling, and modeling facts with the objective of highlighting beneficial statistics, suggesting conclusions, and helping decision making.

Data evaluation has multiple aspects and approaches, encompassing various techniques below an array of names, in distinct enterprise, science, and social technology domain names. Statistics Mining is the discovery of unknown information determined in databases. Information mining capabilities include clustering, type, prediction, and institutions. One of the most essential data mining applications is that of mining association regulations. Association policies, first introduced in 1993, are used to identify relationships amongst a set of gadgets in databases. these relationships aren't primarily based on inherent properties of the records themselves, however as an alternative based on co-prevalence of the facts gadgets. Emphasis on this paper is on the basket marketplace analysis facts.

Numerous algorithms had been proposed to discover common item units in transaction databases. Data mining provides new views for data analysis. The purpose of data mining is to extract and discover new information from statistics. Over the past few a long time, new methods were advanced about the abilities of records collection and statistics generation. Statistics collection tools have supplied us with a big quantity of information. Data mining processes have integrated techniques from multiple disciplines such as, statistics, machine learning, database technology, pattern recognition, neural networks, information retrieval and spatial data analysis. Data mining techniques have been used in many different fields such as, business management, science, engineering, banking, data management, administration, and many other applications. Data mining is a repetitive process consisting of several steps. Starting with the understanding and definition of a problem and ending with the analysis of results and determine a strategy with using the result.

The health care enterprise requires the use of DM due to it generates large and complex volumes of

statistics. Therefore, un-computerized evaluation has emerge as each high priced and impractical. The life of coverage fraud and abuse impels insurers to apply DM. DM can generate information that may be beneficial to all stakeholders in health care, along with patients through identifying effective treatments and excellent practices.

DM came into prominence in mid 90s due to the fact computers made possible the fast production of massive statistics warehouses, containing probably massive amounts of statistics. The present day statistical strategies and the advances in possibility idea offered the important analytical tools. The records of statistics and its contents are a great deal older. Huge collections of information had been constructed over loads and thousands of years by means of numerous sorts of government and scientists.

A well-known case is the large series of very correct planetary observations of the Danish astronomer Tycho Brahe. The know-how embedded in this information the legal guidelines of the moves of the planets have been located by his successor Johannes Kepler and have been confirmed by way of the paintings of Newton. The main DM sports encompass description and visualization, in search of institutions between records elements, grouping facts into units of comparable records (a procedure known as clustering), statistics classification, prediction based on developments that can be extracted from statistics, and so forth.

DM packages in fitness care are numerous and already well set up: comparing treatment effectiveness, health care management, the analysis of relationships between patients and vendors of care, pharma covigilance, fraud and abuse detection. Despite the apparent benefits, there exist many barriers and problems in adapting DM evaluation strategies. DM can be restrained via the accessibility to facts that frequently is shipped in distinct settings (clinical, administrative, insurers, labs, etc.). Records can be incomplete, corrupted, noisy, or inconsistent. There exist ethical, legal and social issues (statistics ownership, privateness worries). Many patterns find in DM may be the end result of random fluctuations; such a lot of such styles can be vain, which requires a severe statistical analysis. DM of clinical facts calls for specific scientific know-how in addition to expertise of DM technology and, final but no longer least, DM requires institutional dedication and funding.

Wide deployment of mobile devices, such as smart phones equipped with low cost sensors, has already shown great potential in improving the quality of healthcare services. Remote mobile health monitoring has already been recognized as not only a potential, but also a successful example of mobile health (mHealth) applications especially for developing countries. The Microsoft launched project "MediNet" is designed to realize remote monitoring on the health status of diabetes and cardiovascular diseases in remote areas in Caribbean countries.

In such a remote mHealth monitoring system, a client could deploy portable sensors in wireless body sensor networks to collect various physiological data, such as blood pressure (BP), breathing rate (BR), Electrocardiogram (ECG/EKG), peripheral oxygen saturation (SpO<sub>2</sub>) and blood glucose. Such physiological data could then be sent to a central server, which could then run various web medical applications on these data to return timely advice to the client. These applications may have various functionalities ranging from sleep pattern analyzers, exercises, physical activity assistants, to cardiac analysis systems, providing various medical consultation.

Moreover, as the emerging cloud computing technologies evolve, a viable solution can be sought by incorporating the software as a service (SaaS) model and pay-as-you-go business model in cloud computing, which would allow small companies (healthcare service providers) to excel in this healthcare market. It has been observed that the adoption of automated decision support algorithms in the cloud-assisted mHealth monitoring has been considered as a future trend. Unfortunately, although cloud-assisted mHealth monitoring could offer a great opportunity to improve the quality of healthcare services and potentially reduce healthcare costs, there is a stumbling block in making this technology a reality.

The health care industry calls for the usage of DM due to it generates large and complicated volumes of data. As a result, un-computerized analysis has come to be each highly-priced and impractical. The lifestyles of insurance fraud and abuse impels insurers to use DM. DM can be useful to all stakeholders in health care, together with sufferers through identifying great treatment and satisfying practice. DM got prominence in mid 90s because computers made feasible the short introduction of massive information warehouses, containing probably huge quantities of records. The cutting-edge today statistical strategies and the development in opportunity precept presented the important logical gear. The records of information and its contents is lots older. Big collections of data were constructed over masses and hundreds of years thru numerous types of government and scientists. The huge series of very correct planetary observations of the Danish astronomer Tycho Brahe. The statistics embedded on this information the criminal guidelines of the actions of the planets have been found through the usage of his successor Johannes Kepler and were showed with the useful resource of the paintings of Newton.

The primary DM sports consist of explanation and visualization, attempting to find institutions amongst records elements, grouping facts into sets of comparable facts (a technique known as clustering), statistics class, prediction primarily based totally on developments that can be extracted from records, and so forth. DM applications in fitness care are numerous and already nicely set up: comparing remedy effectiveness, health care management, the assessment of relationships among patients and carriers of care, pharma covigilance, fraud and misuse detection despite the obvious benefits, there

exist many limitations and problems in adapting DM analysis strategies.

DM can be limited with the help of the person-friendliness to facts that often is deliver in special settings like clinical, administrative, insurers, labs, and so on. Data can be incomplete, corrupted, noisy, or inconsistent. There exist ethical, legal and social problems (information possession, privateness worries). Many styles locate in DM may be the end result of random fluctuations; so many such patterns may be useless, which requires a severe statistical analysis. DM of medical records calls for particular medical information in addition to information of DM generation and, last but no longer least, DM calls for institutional dedication and funding

### III. EXISTING METHOD

In this study, formulated the task of risk prediction as a multi-class classification problem using the Cause of Death (COD) information as labels, regarding the health-related death as the "highest risk". The goal of risk prediction is to effectively classify 1) whether a health examination participant is at risk, and if yes, 2) predict what the key associated disease category is. In other words, a good risk prediction model should be able to exclude low-risk situations and clearly identify the high-risk situations that are related to some specific diseases.

A fundamental challenge is the large quantity of unlabeled data. For instance, 92.6% of the 102,258 participants in our aged physical condition assessment dataset do not have a COD tag. The semantics of such "alive" cases can vary from generally healthy to seriously ill or anywhere in between. In other words, there is no ground truth available for the "healthy" cases. It simply treat this set of "alive" cases as the negative class, it would be a highly noisy majority class. On the other hand, if take this large alive set as genuinely unlabeled, as opposed to cases with known labels removed, it would become a multi-class learning problem with large unlabeled data.

Iris records of patient m1 in UCI machine learning repository with test items in different categories (e.g., sepal length, sepal width, Petal width, petal length, age-at-heart-attack, Survival etc.) and irregular results marked black. this example suggests that 1) a participant may have a series of irregularly time stamped longitudinal data, each of that's probable to be sparse in terms of unusual outcomes, and 2) check items are obviously in classes, each conveying distinctive semantics and in all likelihood contributing in another way in risk identification. Therefore this heterogeneity should be respected in the modeling.

It does no longer attention at the data fusion for the health examination records to be included with different kinds of datasets including the hospital-primarily based electronic fitness records and the individuals' living situations (e.g., diets and general sporting events).by means of integrating facts from

multiple available information resources, extra powerful prediction can be accomplished. In the preceding present gadget class techniques on healthcare data do no longer consider the problem of unlabeled records. It both has professional-defined low-hazard or organizer steering and simply treat non-high quality instances as terrible.

Techniques that keep in mind unlabeled data are usually based totally on Semi-Supervised learning (SSL) that learns from each categorized and unlabeled information. Mining health exam information and gaining knowledge of strategies that cope with unlabeled fitness records. The trouble in modern-day kingdom of art work unlabeled facts offers an in depth account of those taking element in being in healthful exam whose being healthy conditions can special substantially from healthy to very-ill.

There is no get onto land fact for differentiating their states of being wholesome. Identifying the members at threat is important for early caution and preventive intervention. The essential assignment of gaining knowledge of a classification model for threat prediction lies within the unlabeled records that constitute most of the people of the gathered dataset. Especially, the unlabeled information describes the contributors in fitness examinations whose fitness situations can vary substantially from healthful to very-ill. There is no ground reality for differentiating their states of health.

#### IV. PROPOSED DATA MINING ALGORITHMS AND TECHNIQUES

Various algorithms and techniques like Classification, Clustering, Regression, Artificial Intelligence, Neural Networks, Association Rules, Decision Trees, Genetic Algorithm, Nearest Neighbor method etc., are used for knowledge discovery from databases.

Classification is the most typically carried out facts mining method, which employs a set of pre-categorized examples to develop a version that could classify the population of statistics at huge. Scheme discovery and recognition make risk applications are especially well suitable to this type of analysis. This method frequently employs decision tree based type algorithms. The data classification manner involves getting to know and type. In learning the training records are analyzed with the aid of class algorithm.

In classification test data are used to estimate the accuracy of the class guidelines. If the accuracy is suitable the guidelines may be applied to the new information tuples. For a fraud detection software program, this may embody entire information of each fraudulent and valid activities decided on a report-by-report. The classifier training set of rules makes use of those preclassified examples to decide the set of parameters required for correct discrimination. The set of rules then encodes those parameters into a version called a classifier.

*Types of classification models:*

Classification by decision tree induction

Bayesian Classification

Neural Networks

Support Vector Machines (SVM)

Classification Based on Associations

**CADS (COLLABORATIVE ADAPTIVE DATA SHARING):**

To suggest CADS, a Collaborative Adaptive facts sharing platform, which allows records annotation at insertion-time and leverages these annotations at question-time? CADS study with time the facts demand (query workload), which is then used to create adaptive insertion and question bureaucracy. A number of the collaborative information sharing packages in order to gain from a success CADS platform are catastrophe management, company context management, news portals, social networking, and scientific collaboration.

The CADS gadget has sorts of actors: producers and purchasers. Producers add information within the CADS gadget the usage of interactive insertion paperwork and clients search for applicable facts using adaptive question bureaucracy.

**Insertion phase:** The insertion phase starts off evolved with the submission of a brand new document to be blanketed within the repository. After the user uploads the report, CADS analyzes the text and creates an adaptive insertion shape with the set of the most likely {characteristic call, attribute value} pairs to annotate the new file. The consumer fills this form with the desired facts and submits it. The final level includes the storage of the associated file and metadata inside the CADS repository.

**Query phase:** within the query part, the consumer is supplied with an adaptive query form, which supports {characteristic call, characteristic price} situations. Initially, earlier than CADS has started mastering the information demands via processing the query workload, the query form handiest specifies the default attributes (e.g., "report kind", "Date", "region"). The person can specify extra {attribute call, attribute cost} conditions. There's also a typical "Description" characteristic in which the consumer sorts keywords while she does no longer recognize a way to position them in {characteristic name, attribute cost} situations. The machine discourages the consumer from just using the "Description" characteristic, due to the fact this does not allow the gadget to research the user information demand in a established manner, which in turn helps evolving the schema and performing schema mappings.

In a few instances the conditions may additionally cause additional attributes recommendation, which CADS believes might be beneficial for the person to in addition refine the question. As an instance, if the person specifies the characteristic "hurricane category" and former customers who special "hurricane class" also exact "Wind velocity", then the adaptive question form will recommend to the user the characteristic "Wind speed". In addition, if the attribute designated by way of a consumer is much like some other existing characteristic, CADS will recommend a mapping

between the 2 attributes, in the spirit of pay-as-you-pass integration. Additionally, the machine may propose replacing the textual content inside the conventional "Description" attribute fee with a few (attribute calls, attribute cost) conditions.

Algorithm for CADs:

Input: Set of web pages from search engine  $W_s$

Output: Set of key terms  $T_s$

Step1: Read all web pages given for training  $W_s$

Step2: Read stop word list  $S_w$

Step3: For each web page  $W_i$  from  $W_s$

$C$  = Read content of the web page  $W_i$

$C$  = Apply html parser to remove html tags from  $C$

$T_s$  = Split  $C$  with pattern single space

For each term  $T_i$  in the term set  $T_s$

If  $T_i$  present in stop word list  $S_w$  then

Remove  $T_i$  from  $T_s$

End

If ( $T_i$  contains "ing")

$T_i$  = Remove "ing" from  $T_i$   
End

If ( $T_i$  contains "ed")

$T_i$  = Remove "ed" from  $T_i$

End

End

End

Step4: For each  $T_i$  from  $T_s$

Identify presence of bigram  $B_i$

If  $B_i$  Presents

Update  $T_s$

Else

Continue

End

Step5: Return set of textual term set

Step6: Stop

## V. SIMULATION

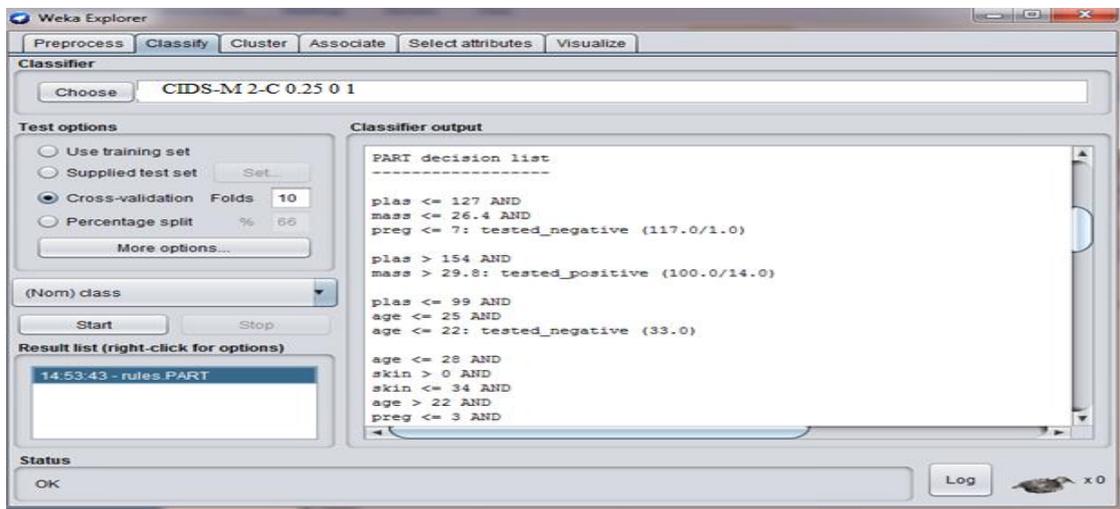


Figure.3

## VI. SIMULATION

OUTPUTS:

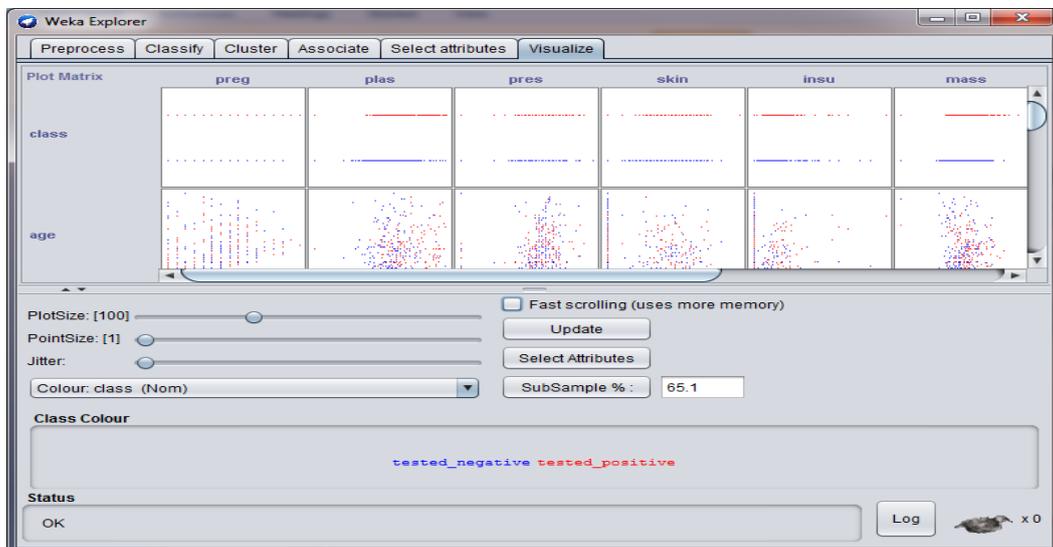


Figure.4

### VII. RESULT AND DISCUSSION

The accomplishment of the novel technique is assessed in terms of evaluation metrics such as the accuracy, running time, sensitivity and the specificity. At the outset, the entire elements are chosen from the two datasets and a column consisting of an attribute is extracted from them. Subsequently, the attribute is effectively transformed into uncertain data attributes by

processing in accordance with the step mentioned . In the performance evaluation procedure, samples of 100, 200 and 300 groups are shortlisted and the three set of uncertain data samples are processed with the novel classification technique. The outcome of the investigation is plotted in the following section. Figure 5.5.1 depicts the comparison graph in terms of Accuracy between SHG and CADs.

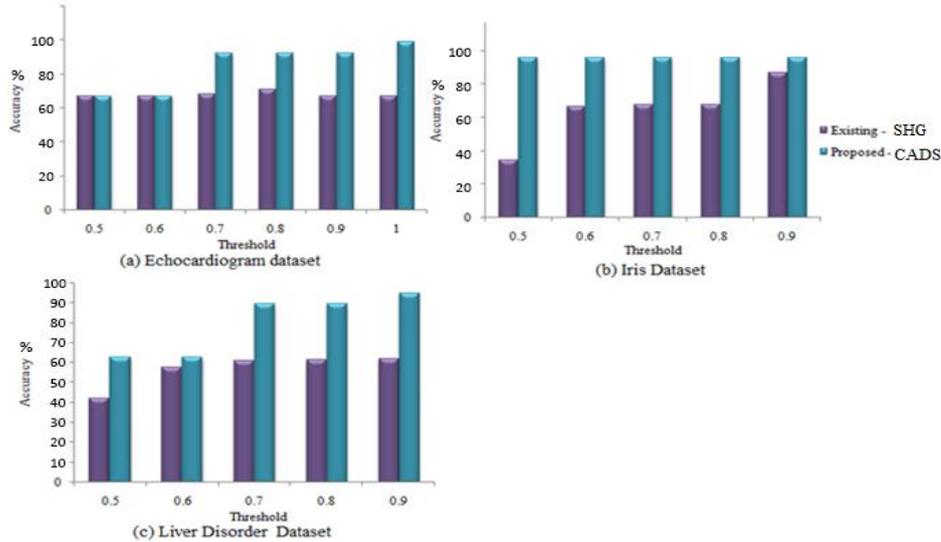


Fig. 5. Accuracy Graph for Different Datasets

Figure 5.5.1 illustrates the Accuracy graph for proposed collaborative adaptive data sharing and modern technique viz Semi-supervised heterogeneous graph in Iris, liver disorder and echocardiogram dataset. The above Figure 5.5.1(b) exhibits the comparison of accuracy of the innovative technique with the modern method for the iris dataset. The accuracy of the novel

algorithm CADs is found to be superior to that of the modern algorithm SHG.

From Figure 5.5.1 (a) the novel CADs algorithm accomplishes superior performance than the modern SHG algorithm with regard to the accuracy in echocardiogram dataset.

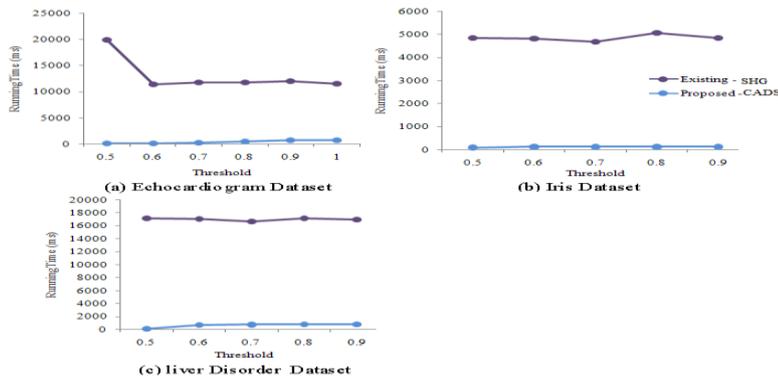


Figure6. Execution time Graph for Different Datasets

The Figure 5.5.2 illustrates the comparison of running time of the novel technique with that of the modern method in respect of the iris dataset. The execution time for the innovative CADs technique is established to be lesser than that for the modern SHG method, thereby leading to the conclusion that the CADs technique is the most endearing candidate for the iris dataset with regard to the running time.

Default Node Size	Proposed –CADs		Existing –SHG	
	Time	Accur acy	Time	Accur acy
1	704	0.9924	11521	0.6742
2	669	0.9924	11506	0.6364

3	623	0.9848	11550	0.5455
4	639	0.9848	11550	0.4697

**Table 1. Performance of CADS in Echocardiogram Dataset**

Default Node Size	Proposed –CADS		Existing –SHG	
	Time	Accuracy	Time	Accuracy
1	170	0.9933	4733	0.8667
2	140	0.9933	4673	0.8667
3	140	0.9933	4710	0.8667
4	130	0.9933	4710	0.8667

**Table 2. Performance of CADS in Iris dataset**

Tables 1 and 2 effectively exhibit the Execution time and accuracy for the Default node in collaborative adaptative procedure. In respect of all the datasets the running time is found to be of minimum value in relation to the existing approach. In case of node 1 the accuracy values are 95.26, 99.3, and 97.25 in echocardiogram, iris disorder databases respectively.

The modern SHG technique needs lesser running time than that required for the innovative CADS method in respect of the width values 0.5, 0.6 and 0.7 respectively. The relative computation takes up a lot of time duration than that of the modern technique, as the novel takes only reasonable time to classify the uncertain data.

## VIII. CONCLUSION

The proposed approach to identify the disease of particular patient by extracting the particular patient's record from a huge database. Our proposed system focuses on medical database hence confidentiality is an important factor. Extracting particular record and processing it to identify disease is done through classification and CADS algorithm. In past CADS platform is available to make exact insertion of patient details and query processing to attain better results. With this idea CADS algorithm was developed where in insertion phase details of patient will be uploaded and in query processing, a particular record is retrieved and prescription will be delivered.

The experimental analysis is conducted for evaluating the performance of the proposed approach. The iris dataset, ECG dataset and liver dataset from the UCI data repository is selected for the performance analysis. Here, evaluation of proposed CADS algorithm is done against existing SHG algorithm with three different real datasets to prove the efficiency of the proposed CADS algorithm is better than SHG algorithm in terms of running time and accuracy. The results from

the experimental analysis show that the collaborative method has achieved a maximum average accuracy of 97.56% while the existing approach has achieved only 92.3%. The analysis indicates that the proposed research used for selecting the feature value enhances the accuracy of classification of uncertain data. Through collaborative filtering recommendation in food with respect to identified disease is provided with respect to patient's current body condition. Hence our approach attains better results when compared to existing approaches where suggestion for identified diseases alone is provided. The last said concept is bound to be implemented as future research work.

## REFERENCES

- [1] Parvez Ahmad, SaqibQamar and Syed QasimAfserRizvi, "Techniques of Data Mining In Healthcare: A Review" International Journal of Computer Applications (0975 – 8887) Volume 120 – No.15, June 2015.
- [2] Ionuț ȚĂRANU, "Data mining in healthcare: decision making and precision" Database Systems Journal vol. VI, no. 4/2015.
- [3] Truyen Tran, DinhPhung, Wei Luo, SvethaVenkatesh, "Stabilized Sparse Ordinal Regression for Medical RiskStratification".
- [4] Ling Chen, Xue Li, Sen Wang, Hsiao-Yun Hu, Nicole Huang, Quan Z. Sheng and Mohamed Sharaf, "Mining Personal Health Indexfrom Annual Geriatric Medical Examinations" 2014 IEEE International Conference on Data Mining.
- [5] Marco Eichelberg, Thomas Aden, And Jo' RgRiesmeier, "A Survey and Analysis of Electronic Healthcare Record Standards" ACM Computing Surveys, Vol. 37, No. 4, December 2005, pp. 277–315.
- [6] Surabhi Thorat, Seema Kute, "Medical Data Mining Life Cycle and its Role inMedical Domain" International Journal of Computer Science and Information Technologies, Vol. 5 (4) , 2014, 5751-5755.
- [7] Mohamed F. Ghalwash, VladanRadosavljevic and ZoranObradovic, "Extraction of Interpretable Multivariate Patterns forEarly Diagnostics" 2013 IEEE 13th International Conference on Data MiningDOI 10.1109/ICDM.2013.19.
- [8] Thanh-Phuong Nguyena, Tu-Bao Hob, "Detecting disease genes based on semi-supervised learning and protein–proteininteraction networks" Artificial Intelligence in Medicine 54 (2012) 63– 71.
- [9] BuyueQian · Xiang Wang · Nan Cao, Hongfei Li · Yu-Gang Jiang, "A relative similarity based method for interactivepatient risk prediction" Data Min Knowl DiscDOI 10.1007/s10618-014-0379-5 Published online: 09 September 2014.
- [10] Jin-Mao Wei, Shu-Qin Wang, and Xiao-Jie Yuan, "Ensemble Rough Hypercuboid Approachfor Classifying Cancers" IEEE transactions on knowledge and data engineering, vol. 22, no. 3, march 2010.
- [11] Chen-Yi Wu, Yi-Chang Choua, Nicole Huangc, Yiing-Jenq Chou, Hsiao-Yun Hua, Chung-Pin Li,

“Cognitive impairment assessed at annual geriatric health examinations predicts mortality among the elderly” *Preventive Medicine* 67 (2014) 28–34.

[12] Vagelis Hristidis Eduardo Ruiz, “CADS: A Collaborative Adaptive Data Sharing Platform”.

[13] Michael D. Ekstrand, John T. Riedland Joseph A. Konstan, “Collaborative Filtering Recommender Systems” *Foundations and Trends in Human-Computer Interaction* Vol. 4, No. 2 (2010) 81–173 @2011.

[14] Serge Abiteboul, Neoklis Polyzotis, The Data Ring: Community Content Sharing, In *CIDR*, pages 154-163, 2007.

[15] G. Jeh, and J. Widom. SimRank: a measure of structural-context similarity. *ACM SIGKDD international Conference on Knowledge Discovery and Data Mining*. KDD 2002.

[16] J. Banerjee, W. Kim, H. Kim, and H. F. Korth. 1987. Semantics and implementation of schema evolution in object-oriented databases. *SIGMOD Rec.* 16, 3 (Dec. 1987), 311-322.

[17] Heymann, P., Koutrika, G., and Garcia-Molina, H. Can social bookmarking improve web search?. *International Conference on Web Search and Web Data Mining*. WSDM '08.