

Mini review Article

A model for **Coronary Heart Disease** Prediction using Data Mining Classification Techniques (Decision Trees, Naive Bayes, and KNN)

8 **ABSTRACT**

9 Nowadays the guts malady is one amongst the foremost causes of death within the world. Thus it's early prediction and diagnosing is vital in medical field, which might facilitate in on time treatment, decreasing health prices and decreasing death caused by it. The treatment value the disease isn't cheap by most of the patients and Clinical choices are usually raised supported by doctors" intuition and skill instead of on the knowledge-rich information hidden within the stored data. The model for prediction of heart disease using a classification techniques in data mining reduce medical errors, decreases unwanted exercise variation, enhance patient well-being and improves patient results. The model has been developed to support decision making in heart disease prediction based on data mining techniques. The experiments were performed using the model, based on the three techniques, and their accuracy in prediction noted. The decision tree, naïve Bayes, KNN (K-Nearest Neighbors) and WEKA API (Waikato Environment for Knowledge Analysis-application programming interface) were the various data mining methods that were used. The model predicts the likelihood of getting a heart disease using more input medical attributes. 13 attributes that is: blood pressure, sex, age, cholesterol, blood sugar among other factors such as genetic factors, sedentary behavior, socio-economic status and race has been use to predict the likelihood of patient getting a Heart disease until now. This study research added two more attributes that is: Obesity and Smoking. 740 Record sets with medical attributes was obtained from a publicly available database for heart disease from machine learning repository with the help of the datasets, and the patterns significant to the heart attack prediction was extracted and divided into two data sets, one was used for training which consisted of 296 records & another for testing consisted of 444 records, and the fraction of accuracy of every data mining classification that was applied was used as standard for performance measure. The performance was compared by calculating the confusion matrix that assists to find the precision recall and accuracy. High performance and accuracy was provided by the complete system model. Comparison between the proposed techniques and the existing one in the prediction capability was presented. The model system assists clinicians in survival rate prediction of an individual patient and future medication is planned for. Consequently, the families, relatives, and their patients can plan for treatment preferences and plan for their budget consequently.

35 Keywords: WEKA API; Decision Tree; Naïve Bayes; KNN, Cardiovascular disease, KDD.

37 **1. INTRODUCTION**

38 The Heart is a strong organ, situated close to the middle of the chest; it is duty is pumping blood to different parts of the body and together with system of vessels and blood from the human body's cardiovascular framework; Interferences to this dissemination of blood can result in serious medical issue including death [5]. People have been affected by dangerous sicknesses all through the past. The system for prediction can assist to lower the dangers of the disease. Prediction is done dependent on the present data fed to the framework model Using WEKA API which is open source information mining programming

44 in Java. The model is being created dependent on three distinct information mining strategy that is Nave
45 Bayes, KNN, decision tree with WEKA API. The input dataset is analyzed using different classification
46 algorithms and comparison is done for accuracy.

47
48 Nowadays an immense measure of information is gathered and kept in a daily basis. There is a
49 significant need to break down this information yet with no scientific device, this appears to be
50 unimaginable. This has prompted the improvement of Knowledge Discovery in Databases (KDD) which
51 changes the low dimension information to a top state learning. KDD comprises of different procedures at
52 various advances and Data mining is one of those procedures. Information mining is the way toward
53 finding fascinating learning from huge measure of information kept in databases, information stockrooms
54 or other data vaults. The fundamental point of information mining procedure is to separate data from a
55 dataset and change it into a reasonable structure so as to help basic conclusions [45]. A tremendous
56 measure of information is accessible in healthcare industry however the mining of this information is poor.
57 In this way, the investigation of the medicinal services information is a must. Information Discovery in
58 databases is getting to be famous research instrument for open human services information. In this study,
59 we will do the exhibition investigation of various information mining grouping strategies on medicinal
60 services information from the Cleveland, Hungary, Switzerland and the VA Long Beach Clinics
61 Foundation, medical records department. This work will help discovering the best information mining
62 arrangement method as far as precision on the specific dataset. The examined characterization systems
63 are K-closest neighbor (KNN), Naive Bayes, Decision tree. The exhibition of these procedures is
64 estimated dependent on their exactness. This investigation will assist the future scientists with getting
65 proficient outcomes in the wake of realizing best information mining grouping method for specific dataset.

66 Information Mining is the nontrivial procedure of recognizing substantial, novel, conceivably valuable and
67 at last reasonable example in information with the wide utilization of databases and the touchy
68 development in their sizes. Information mining refers to removing or "mining" learning from a lot of
69 information. Information digging is the quest for the connections and worldwide examples that exist in
70 enormous databases however are tucked away among a lot of information [17]. The fundamental
71 procedure of Knowledge Discovery is the change of information into learning so as to help in making
72 judgments is known as information mining. Information Discovery procedure comprises of an iterative
73 grouping of information cleaning, information coordination, information determination, information mining
74 design acknowledgment and learning introduction. Information digging is the quest for the connections
75 and worldwide examples that exist in enormous databases bramble are tucked away among a lot of
76 information.

77 Many hospitals have put in databases systems to manage their clinical data or patient data. These data
78 systems generally generate giant amounts of information which may be in any format like numbers, text,
79 charts and pictures however sadly, this info that contains made information isn't used for clinical deciding.
80 There's abundant data keep in repositories that may be used effectively to support deciding in attention.
81 Data processing techniques is wide utilized in medical field for extracting information from info. In data
82 processing call tree may be a technique that is employed extensively. Call trees are non-parametric
83 supervised learning technique used for classification.

84 The most aim is to form a model that predicts the worth of a target variable by learning straightforward call
85 rules inferred from the info options. The structure of the choice tree is within the type of tree and leaf
86 nodes. Decision trees are most typically utilized in research, principally in call analysis. Blessings are that
87 they're straightforward to know and interpret. They're strong, performed well with giant datasets, able to
88 handle each of the numerical and categorical information.

89 By providing economical treatments, it will facilitate to scale back prices of treatment. Mistreatment data
90 processing techniques it takes less time for the prediction of the un-wellness with a lot of accuracy.
91 The most necessary step a company will absorb terms of information mining is to require advantage of
92 the opportunities afforded by it. Collect information and place it to smart use with data processing, and
93 you'll before long begin reaping the advantages that's ; more cash by Learning that varieties of
94 merchandise customers have purchased and maximize that insight to individualize expertise, increase

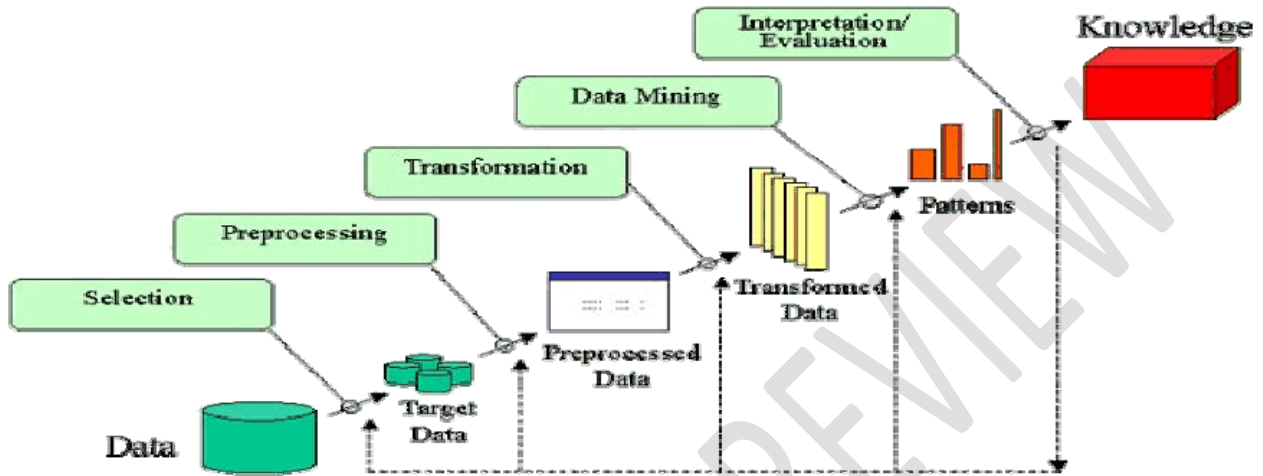
95 client loyalty, and boost client time period price. Improve stigmatization and promoting through Get
96 feedback and use data processing to spot what's operating and what isn't with branding and marketing.
97 contour reach by creating all of your outreach a lot of timely and relevant with data processing, faucet into
98 new markets by Use different databases to spot potential customers and conduct relevant reach, Learn
99 from the past by comparison current information to past data to search out trends to stay in mind once
100 creating business choices.

101
102 Data mining has become more and more necessary, particularly in recent years, once nearly all industries
103 and sectors everywhere the planet face issues on information explosion. All of unforeseen, there's just
104 too abundant information, and this fast rise within the quantity of information demands a corresponding
105 increase in the amount of knowledge and knowledge. Thus, there's a requirement to quickly,
106 expeditiously and effectively method all that information into usable data and data processing offers the
107 answer. In fact, you'll say that data processing is that the resolution. You'll realize data processing to be
108 most frequently used or applied in organizations or businesses that maintain fairly giant to large
109 databases. The sheer size of their databases and also the quantity of knowledge contained among them
110 need over a little live of organization and analysis that is wherever data processing comes in. Through
111 data processing, users are able to investigate information from multiple views in their analysis. It'll
112 additionally build it easier to categorize the knowledge processed and establish relevant patterns,
113 relationships or correlations among the assorted fields of the data. Therefore, we are able to deduce that
114 data processing involves tasks of a descriptive and prognosticative nature. Descriptive, as a result of it
115 involves the identification of patterns, relationships and correlations among giant amounts of information,
116 and prognosticative, as a result of its application utilizes variables that are accustomed predict their future
117 or unknown values. The use of information mining (DM) model allows machine intelligence in nosology
118 processes.

119 DM is that the machine intelligence-based process of extracting important data from the set of huge
120 quantity of information. DM may be a speedily growing field in a very big selection of health science
121 applications. Applicable DM-based classification techniques and sensible cardiovascular disease
122 prediction systems will lead toward quality health care in terms of accuracy and low economical health
123 care services. The most motivation behind digitization of health information and utilization of sentimental
124 computing tools is to lower the value of health care and cut back the quantity of preventable errors.
125 Among numerous DM techniques, like agglomeration, association rule classification and regression, the
126 classification is one among the foremost necessary techniques used for categorization of information
127 patterns. In DM, essentially the classification-based machine learning algorithms are accustomed predict
128 membership perform for labeling CVD information instances. Classification will be information analysis
129 technique that extracts labels describing necessary data categories. The classifier's model is portrayed as
130 classification rules, call trees or mathematical formulae, and it's termed as supervised learning. The
131 model is employed for classifying future or unknown objects. The classification algorithmic program
132 predicts un-wellness categorical class (eg, negative and positive) and build classifier model supported the
133 coaching set. If the accuracy of the model is suitable, the model may be applied to categorize information
134 tuples whose class labels are unknown. The classification contains 2 basic steps of learning and
135 classification. In learning, coaching information is analyzed by classification algorithmic program and
136 classifier's model is made. Within the classification section, check information are utilized to estimate the
137 accuracy of the classification model. A healthy range of researchers are applying numerous algorithms
138 and techniques like classification, clustering, multivariate analysis, artificial neural networks (ANNs), call
139 trees, genetic algorithmic program (GA), KNN strategies, single DM model and hybrid and ensemble
140 approaches to help health care professionals with improved accuracy within the identification of
141 cardiovascular disease. During this study, the analysis quest of however the burden of artery un-wellness
142 may be considerably reduced through soft machine strategies is explored. The final drawback statement
143 of this study is to develop approach-based classifier's model that may be applied to CVD information sets
144 to boost model prediction's outcomes for higher prediction accuracy and responsibility. Additionally to the
145 current, the study presents example of intelligent cardiovascular disease prediction system supported
146 associate degree approach with totally different classifiers, namely, Naive theorem and KNN. The

147 planned prediction system is computer program primarily based, having the power of scaling and
148 enlargement as per user's additional demand.

149 The figure beneath illustrates Steps of the Knowledge Discovery in Databases process on the most
150 proficient method to separate learning from information with regards to enormous databases Fayyad et.al
151 [14].



152

153 **Figure1.0 Steps of Knowledge Discovery in Databases process by Fayyad et.al [14]**

154

155 Various health industry information systems are structured to help patient charging, stock organization
156 and making some simple calculation. A couple of health sectors utilize decision model systems yet are,
157 as it were, limited. They can address simple inquiries like "What is the ordinary time of patients who have
158 coronary disease?" "What number of therapeutic techniques had achieved crisis facility stays longer than
159 10 days?", "Recognize the female patients who are single, more than 30 years old, and who have been
160 treated for coronary sickness." However they can't respond to complex inquiries like "Given patient
161 records, foresee the probability of patients getting a coronary disease." Clinical decisions are as often as
162 possible made subject to experts' impulse and experience rather than on the learning rich data concealed
163 in the database.

164 This preparation prompts bothersome tendencies, botches and super helpful costs which impacts the
165 idea of care provided for patients. The proposed structure that coordinates the clinical decision help with
166 PC based patient records could reduce therapeutic errors, overhaul tolerant security, decrease
167 bothersome practice assortment, and improve getting result. This suggestion is promising as data
168 modeling and analysis tool like data mining can make a learning rich condition which can help to in a
169 general sense improve the idea of clinical decisions.

170 In this fast moving world people need to continue with an extravagant life so they work like a machine to
171 win some portion of money and continue with a pleasant life appropriately in this race they disregard to
172 manage themselves, because of this there sustenance affinities change in their entire lifestyle change, in
173 this sort of lifestyle they are logically stressed they have heartbeat, sugar at a young age and they don't
174 give enough rest for themselves and eat what they get and they even don't overemphasize the idea of the
175 sustenance whenever cleared out the go for their own special prescription in light of all these little
176 indiscretion it prompts a significant threat that is the coronary disease [7]. On account of this people go to
177 therapeutic administrations experts but the prediction made by them isn't 100% definite [25].

178 Quality facility proposes diagnosing patients precisely and controlling medications that are convincing.
179 Poor clinical decisions can incite tragic outcomes which are along these lines unsatisfactory. Medicinal
180 centers ought to in like manner limit the cost of clinical tests. They can achieve these results by using
181 fitting PC based information or decision support system.

182 The treatment cost of heart disease is not affordable by most of the patients, and the Clinical decisions
183 are often made based on doctors' intuition and experience rather than on the knowledge-rich data hidden
184 in the database. This practice leads to unwanted biases, errors and excessive medical costs which
185 affects the quality of service provided to patients. The proposed model for Heart Disease Prediction using
186 Data Mining Classification Techniques reduces medical errors, enhances patient safety, decrease
187 unwanted practice variation, reduce treatment cost and improves patient outcome. This suggestion is
188 promising as data modeling and analysis tools have the potential to generate a knowledge-rich
189 environment which can help to significantly improve the quality of clinical decisions [32].

190 2. LITERATURE REVIEW

191 This part goes for investigating the different information mining methods presented as of late for coronary
192 illness expectation. The man-made brainpower methods centering K-closest neighbor (KNN), Naive
193 Bayes and Decision tree will be presented. Recently distributed papers in displaying survival will be talked
194 about and the recommendations for another strategy are introduced

195 2.1 Theoretical and Empirical Review

196 Various information mining systems have been utilized in the analysis of cardiovascular disease (CVD)
197 over various Heart illness datasets. A few papers utilize just a single method for conclusion of coronary
198 illness and different scientists utilize more than one information mining technique for the finding of
199 coronary illness.

200 In [23,27] Jyoti et.al presented three classifiers Decision Tree, Naïve Bayes and Classification by
201 methods for gathering to break down the proximity of coronary sickness in patients. Request by methods
202 for bundling: Clustering is the route toward social occasion relative segments. This framework may be
203 used as a preprocessing adventure before urging the data to the portraying model. Preliminaries were
204 driven with WEKA 3.6.0 gadget Enlightening list of 909 records with 13 particular properties. All properties
205 were made supreme and anomalies were made due with straightforwardness. To update the desire for
206 classifiers, innate request was joined. Observations show that the Decision Tree data mining technique
207 beats other two data mining methods in the wake of intertwining feature subset assurance yet with high
208 model improvement time.

209
210 [27] Nidhi et.al discernments revealed that the Neural Networks with 15 characteristics improved in
211 examination with other data mining frameworks [27]. The investigation concentrate assumed that
212 Decision Tree technique showed better execution with the help of innate figurings using included subset
213 assurance. This examination work furthermore proposed a model of Intelligent Heart Disease Prediction
214 structure using data mining frameworks explicitly Decision Tree, Naïve Bayes and Neural Network. An
215 aggregate of 909 records were obtained from the Cleveland Heart Disease database. The results
216 declared in the investigation work guarded the better execution of Decision Tree methodology with 99.6%
217 accuracy using 15 qualities. In any case, Decision tree technique in mix with inherited estimation the
218 introduction declared was 99.2% using 06 qualities.

219
220 In [8,9] Chaitrali et.al exhibited that Artificial Neural Network outmaneuvers other data mining
221 methodology, for instance, Decision Tree and Naïve Bayes. In this investigation work, Heart disorder
222 desire system was made using 15 characteristics [8,9]. The investigation work included two extra
223 properties weight and smoking for capable finish of coronary sickness in making convincing coronary
224 disease desire system.

225
226
227 [31] Researchers in year 2013 showed Hybrid Intelligent Techniques for the figure of coronary ailment.
228 Some Heart Disease gathering system was researched in this examination and shut with legitimization
229 noteworthiness of data mining in coronary sickness end and course of action. Neural Network with
230 separated getting ready is helpful for sickness conjecture in starting time and the extraordinary execution

231 of the structure can be gotten by preprocessed and institutionalized dataset. The game plan precision can
232 be improved by decline in features.

233
234 [47] Vikas et.al, in their investigation work used three standard data mining figuring's CART (Classification
235 and Regression Tree), ID3 (Iterative Dichotomized 3) and decision table (DT) removed from a decision
236 tree or rule based classifier to develop the conjecture models using a greater dataset. Discernment
237 showed that presentation of CART computation was better when differentiated and other two course of
238 action procedures.

239
240 V. Manikandan et.al in [46] recommended that association standard mining is used to remove the thing
241 set relations. The data game plan relied upon MAFIA counts which achieved better precision. The data
242 was surveyed using entropy based cross endorsement and bundle strategies and the results were
243 considered. MAFIA (Maximal Frequent Item set Algorithm) used a dataset with 19 characteristics and the
244 goal of the examination work was to have exceedingly definite audit estimations with bigger measures of
245 precision.

246
247 Beant et.al in [6] circulated an investigation paper in IJRITCC "Review on Heart Disease using Data
248 Mining Techniques". The maker referenced created by gigantic number of experts and investigated
249 diverse data mining strategies reliant on execution and accuracy.

250
251 Methaila et.al [3] in their examination work focused on using different counts and mixes of a couple of
252 target qualities for amazing heart ambush figure using data mining. Decision Tree has beaten with
253 99.62% precision by using 15 characteristics. Moreover the exactness of the Decision Tree and Bayesian
254 Classification further improves in the wake of applying inherited computation to diminish the genuine data
255 size to get the perfect subset of value satisfactory for coronary disease estimate.

256
257 The experts [19] proposed a model for desire for coronary ailment using J48, Bayes Net, and Naïve
258 Bayes, Simple CART and REPTREE Algorithms using understanding educational accumulation from
259 Medical Practitioners.

260
261 Appraisal of the disorder matrix showed that J48, REPTREE and SIMPLE CART exhibit a figure model of
262 89 cases with a peril factor positive for heart attacks. The strategies immovably prescribed that data
263 mining counts can foresee a class for judgments.

264
265 B.Venkatalakshmi et.al [5] played out an examination on coronary disease finding using data mining
266 methodology Naïve bayes and Decision Tree techniques. Different sessions of examinations were
267 coordinated with the proportional datasets in WEKA 3.6.0 contraption. Instructive gathering of 294
268 records with 13 attributes was used and the results revealed that the Naïve Bayes beat the Decision tree
269 frameworks.

270
271 The synopsis of looked into writing alongside the quantity of properties utilized for the forecast of
272 Cardiovascular Disease (CVD) is given in table beneath

273
274 **Table 1.0: Table shows different data mining techniques used in the diagnosis of Heart**
275 **disease.**

Author/Researcher	Data Mining Technique used	Year	Number of Attributes Selected
Jyoti Sonia, et.al.	Naïve Bayes, Decision Tree, KNN	2011	13

K.Srinivas et.al.	Naïve Bayes, knn and D.L.	2011	14
Nidhi Bhatla et.al.	Naïve Bayes, Decision Tree, Neural Network	2012	15 and 13
Chaitrali S.Dangare & Sulabha S.Apte	Naïve Bayes, Decision Tree, Neural Network	2012	13 and 15
Abhishek Taneja	Naïve Bayes, J48 unpruned tree, Neural Network	2013	15 and 8
R. Chitra et. al.	Hybrid Intelligent Techniques	2013	15
Vikas Chaurasia, et.al.	CART, ID3, Decision Table	2013	Not mentioned
V. Manikandan et al.	K-Mean based on MAFIA, K-Mean based on MAFIA with ID3, K-Mean based on MAFIA with ID3 and C4.5	2013	19
Beant Kaur & Williamjeet	Papers Reviewed	2014	Nil
Aditya Methaila et. al.	Decision Tree, Naive Bayes, Neural Network, Genetic Algorithm	2014	15 and 16
Hlaudi Daniel Masethe, Mosima Anna Masethe	J48, REPTREE, Naïve Bayes, Bayesnet, Simple CART	2014	15
B.Venkatalakshmi and M.V Shivsankar	Decision Tree and Naïve Bayes	2014	13

277 **2.2 Artificial Intelligence Techniques in Heart Disease Prediction**

278 Information mining has been generally connected in the therapeutic field as this give enormous measure
279 of information. Different scientists had connected the various information mining procedures on social
280 insurance information [11]. connected 5 arrangement calculations for example choice tree, fake neural
281 system, strategic relapse, Bayesian systems and credulous Bayes and stacking-sacking technique for
282 structure arrangement models and thought about the precision of the plain and outfit model to foresee
283 whether a patient will return to a medicinal services Center or not. From results, the best order model
284 relies upon informational collection for example ANN (**Artificial neural networks**) in 3M informational
285 index, choice tree in 6M and strategic relapse in 12M informational collection [23, 26] contrasted the
286 information mining and conventional insights and expresses a few focal points of mechanized information
287 framework. This paper gives an outline of how information mining is utilized in social insurance and
288 medication. Patil Dipti [29] decides if an individual is fit or unfit dependent on authentic and constant
289 information utilizing grouping calculations that is K-means and D-stream are connected. The presentation
290 and precision of D-stream calculation is more than K-implies [4] utilized choice tree to construct an
291 arrangement model for anticipating representative's exhibition. To manufacture a characterization model
292 CRISP-DM was received.

293 PC reproduction demonstrates that the strategic relapse, neural system model and troupe model
294 delivered best generally speaking grouping precision. Koç et al [24] connected ANN and strategic relapse
295 to anticipate if the customer will buy in a term store or not subsequent to promoting effort. ANN orders
296 84.4% information accurately while calculated relapse characterizes 83.63% information effectively
297 however LR takes 54 seconds and ANN takes 11 seconds to run. Along these lines, with more
298 information and higher dimensional element space, utilizing ANN will be progressively productive. Fartash
299 et.al [13] contrasted the different order calculations with anticipate the transmission capacity use design in
300 various time interims among various gatherings of clients in the system correlation of various
301 characterization calculations including. Choice Tree and Naïve Bayesian utilizing Orange is finished. The
302 Decision Tree calculation accomplished 97% exactness and effectiveness in foreseeing the required data
303 transfer capacity inside the system. Sakshi and Prof.Sunil Khare [35] gave a total examination of various
304 information mining characterization procedures that incorporates choice tree, Bayesian systems, k-closest
305 neighbor classifier and fake neural system.

306 Clinical databases have gathered enormous amounts of data about patients and their ailments. The term
307 Heart illness includes the assorted sicknesses that influence the heart. Coronary illness is the real reason
308 for setbacks on the planet. The term Heart illness includes the assorted ailments that influence the heart.
309 Coronary illness kills one individual at regular intervals in the United States [48]

310 **2.3 Data Mining Review**

311 Notwithstanding the way that data burrowing has been around for more than two decades, its potential is
312 simply being recognized now. Data mining solidifies quantifiable examination, AI and database
313 advancement to think hid models and associations from gigantic databases Fayyad portrays data mining
314 as "a method of nontrivial extraction of saw, in advance darken and possibly profitable information from
315 the data set away in a database" [44] describes it as "a method of assurance, examination and showing
316 of colossal measures of data to discover regularities or relations that are at first cloud with the purpose of
317 getting clear and accommodating results for the owner of database" [17]

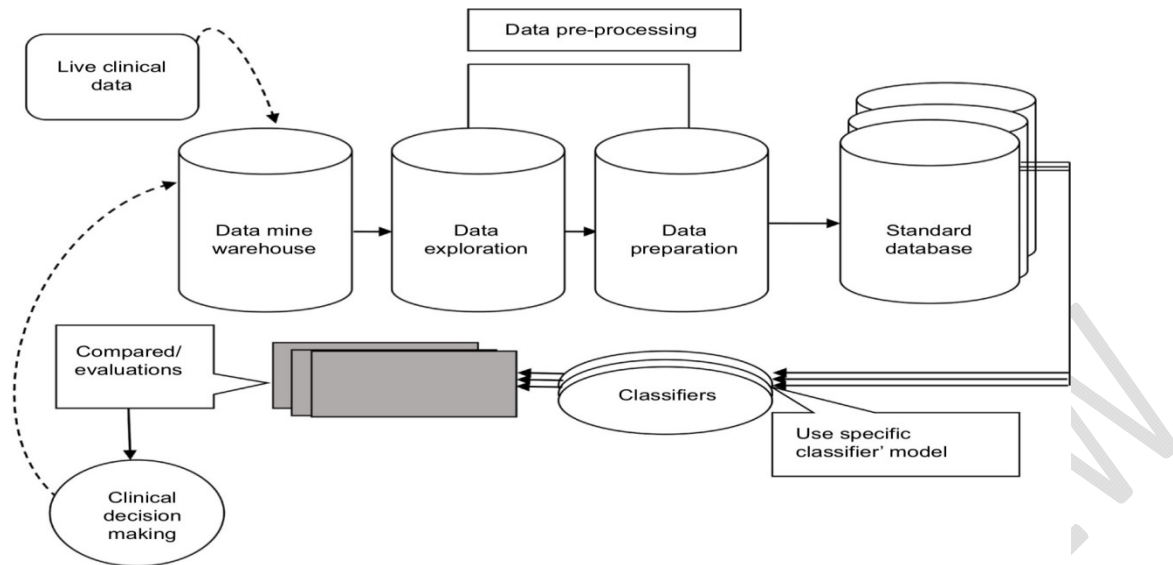
318 Data mining uses two systems: oversight and unsupervised learning. In oversight learning, a planning set
319 is used to learn model parameters however in unsupervised adjusting no arrangement set is used (e.g., k
320 means grouping is unsupervised) [28]. Each data mining methodology fills another need dependent
321 upon the exhibiting objective. The two most ordinary showing goals are gathering and figure. Game plan
322 models predict full scale names (discrete, unordered) while estimate models envision steady regarded
323 limits Decision Trees and Neural Networks use portrayal counts while Regression, Association Rules and
324 Clustering use desire figurings [10]. Decision Tree figurings consolidate CART (Classification and
325 Regression Tree), ID3 (Iterative Dichotomized [10] and C4.5. These computations shift in selection of
326 parts, when to keep a center point from part, and undertaking of class to a non-split center [11] CART
327 uses Gini rundown to check the dirtying impact of a package or set of getting ready tuples [17]. It can
328 manage high dimensional unmitigated data.

329 Decision Trees can moreover manage constant data (as in backslide) yet they ought to be changed over
330 to straight out data. Gullible Bayes or Bayes' Rule is the explanation behind a few, AI and data mining
331 methods [42]. The standard (estimation) is used to make models with insightful capacities. It gives better
332 methodologies for researching and getting data. It gains from the "evidence" by figuring the association
333 between the goal (i.e., subordinate) and other (i.e., independent components). Neural Networks includes
334 three layers: input, concealed and yield units (factors). Relationship between data units and concealed
335 and yield units rely upon centrality of the doled out worth (weight) of that particular data unit. The higher
336 the weight the more huge it is. Neural Network computations use Linear and Sigmoid trade limits. Neural
337 Networks are sensible for setting up a ton of data with few wellsprings of information. It is used when
338 various systems are unacceptable.

339 3. RESEARCH DESIGN

340 Methodology provides a framework for endeavor the projected DM modeling. The methodology may be a
341 system comprising steps that remodel information into recognized data patterns to extract information for
342 users. The DM methodology framework breaks down the mining method of vast knowledge into phases. It
343 shows associate degree unvaried DM method for implementing machine learning strategies on the vast
344 knowledge set taken for application. The projected methodology includes steps, stated because the
345 preprocessing stage wherever the thoroughgoing exploration of the information is disbursed. It'll account
346 for handling missing values, equalization knowledge and normalizing attributes counting on algorithms
347 used. Once pre-processing of information is performed, prognostic modeling of the information is
348 disbursed victimization classification models and ensemble approach. Finally, prescriptive modeling is
349 undertaken, wherever the prognostic model is evaluated in terms of performance and accuracy
350 victimization varied performance metrics. The figure below shows a framework break down of the
351 unvaried data mining process of vast knowledge into phases

352



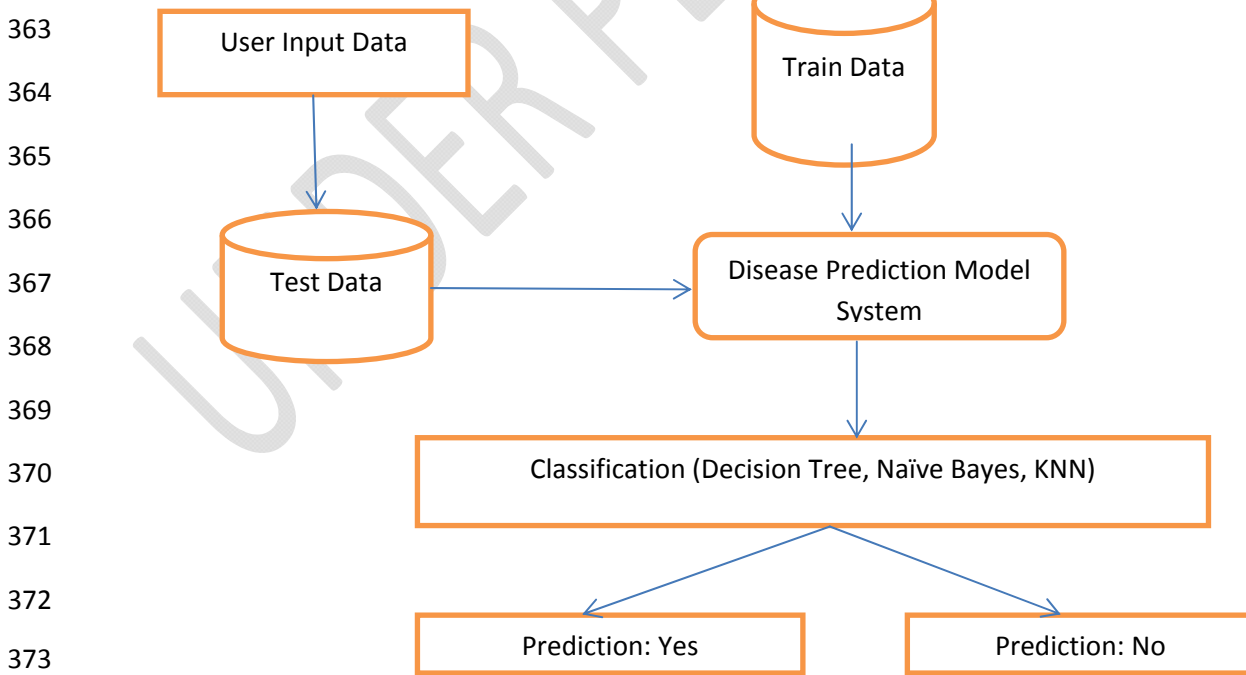
353

354 **Figure 2.0: Methodology for mining heart disease data.**

355 In this examination, three information digging procedures for prescient information mining assignment
 356 were utilized, that incorporates Decision tree, K-NN and Naïve Bayes. These strategies were utilized for
 357 producing learning to settle on it valuable for basic leadership. Every strategy delivered various outcomes
 358 to arrange the locale into centered or non-centered states involving the accessible factors in dataset .The
 359 experimentation was performed utilizing WEKA programming interface.

360 **3.1. Proposed Model**

361 The proposed engineering of coronary illness forecast framework is given below. The figure illustrates the
 362 frame work of the coronary heart disease prediction model steps activities.



374 **Figure 3.0: The System Model**

375 It comprises of preparing dataset and client contribution as the test dataset. Weka information mining
376 apparatus with programming interface was utilized to actualize the coronary illness forecast framework.
377 The source code of Weka is in java. The framework is planned with java swing and use Weka
378 programming interface to call the various techniques for Weka. The segments utilized are cases, various
379 classifiers and strategies for assessment. Administered learning strategy is utilized here. A directed
380 learning calculation examinations the preparation information and derives a capacity from the named
381 preparing set. It tends to be utilized for mapping new models. The preparation information got from ucl
382 repository coronary illness database is the preparation model. This preparation information comprise of
383 the class name and its comparing esteem. Credulous Bayes, KNN and choice tree classifiers are
384 administered learning calculations. They gain from the given preparing models. At the point when another
385 case with same characteristics as in preparing information with various qualities other than those in the
386 preparation model comes, these calculations effectively characterize the new case dependent on the
387 speculation made from the preparation set. Gullible Bayes, KNN and choice tree classifiers are order the
388 new perception into two classifications based on preparing dataset. The preparation dataset is in the
389 ARFF group. The preparation set comprises of 296 traits including the class characteristic. Coronary
390 illness forecast framework acknowledges contribution from the client through a graphical UI. Every one of
391 the traits required for grouping is gotten from a content field. The graphical UI is fabricated utilizing swing.
392 The following procedure is to move the client information acquired from graphical UI into a record of CSV
393 (Comma isolated Value) augmentation. At that point the CSV record is changed over into ARFF
394 document. Weka programming interface give local strategies to changing over from CSV to ARFF. The
395 changed over client information is treated as test information. The test informational index will contain
396 every one of the characteristics of preparing dataset. In the event that the client did not enter a property
397 estimation a '?' will be relegated at the estimation of that comparing trait. Weka will deal with this missing
398 worth. This test information is kept running on Naive Bayes, KNN and choice tree calculation. These
399 calculations order the occasions got from the client and foresee the opportunity to have coronary illness.
400 Netbeans IDE is utilized to code in Java.

401 3.1.1 Decision Tree

402 A call tree could be a decision support tool that uses a tree-like model of selections and their
403 doable consequences, as well as happening outcomes, resource prices, and utility. It's a way to
404 show AN algorithmic program that solely contains conditional management statements.
405 Decisions trees are ordinarily utilized in research, specifically in call analysis, to assist determine
406 a technique possibly to succeed in a goal, however also are preferred tools in machine learning.
407 Classification is that the method of building a model of categories from a collection of records
408 that contains category labels. Decision Tree algorithmic program is to seek out the method the
409 attributes-vector behaves for variety of instances. Additionally on the bases of the coaching
410 instances the categories for the freshly generated instances are being found. This algorithmic
411 program generates the principles for the prediction of the target variable. With the assistance of
412 tree classification algorithmic program the vital distribution of the information is well
413 comprehensible [50]. J48 is AN extension of ID3. The extra options of J48 are accounting for
414 missing values, call trees pruning, continuous attribute worth ranges, derivation of rules, etc.
415 within the wood hen data processing tool, J48 is AN open supply Java implementation of the
416 C4.5 algorithmic program. The wood hen tool provides variety of choices related to tree pruning.
417 Just in case of potential over fitting pruning is used as a tool for précising. In different algorithms
418 the classification is performed recursively until each single leaf is pure, that's the classification of
419 the information ought to be as excellent as doable. This algorithmic program it generates the
420 principles from that specific identity of that knowledge is generated. The target is more and more
421 generalization of a call tree till it gains equilibrium of flexibility and accuracy.

422 423 3.1.2 Naïve Bayes

424 This technique depends on probabilistic information. The gullible Bayes principle yields probabilities for
 425 the anticipated class of every individual from the arrangement of test example. Gullible Bayes depends on
 426 administered learning. The objective is to foresee the class of the experiments with class data that is
 427 given in the preparation information.

428 The quality "Analysis" is distinguished as the anticipated characteristic with worth "1" for patients with
 429 coronary illness and worth "0" for patients with no coronary illness. "Quiet Id" is utilized as the key; the
 430 rest are info traits. It is expected that issues, for example, missing information, conflicting information, and
 431 copy information have all been settled.

432 It is a probabilistic classifier supported Bayes' theorem such by the previous possibilities of its root nodes.
 433 The mathematician theorem is given in Equation one and social control constant is given in Equation a
 434 pair of. It proves to be associate best formula in terms of diminution of generalized error. It will handle
 435 statistical-based machine learning for feature vectors $Y = [Y_1, Y_2, \dots, Y_n]^T$ and assign the label for feature
 436 vector supported supreme probable among on the market categories. It means feature "y" belongs to X_i
 437 category, once posterior likelihood $P(X_i|Y)$ is most i.e $Y=X_i; P(X_i|Y)_{Max}$. The Bayesian classification
 438 downside is also developed by a-posterior possibilities that assign the category label ω_i to sample X
 439 specified $P(X_i|Y)$ is supreme. The Bayesian classification downside is also developed by a-posterior
 440 possibilities that assign the category label ω_i to sample X specified $P(X_i|Y)$ is supreme.

$$P(X_i | \underline{y}) = \frac{p(\underline{y} | X_i) P(X_i)}{p(\underline{y})} \quad (1)$$

$$p(\underline{y}) = \sum_{i=1}^2 p(\underline{y} | X_i) P(X_i) \quad (2)$$

441
 442 Application of Bayes' rule with the mutual exclusivity in diseases and also the conditional independence in
 443 findings is understood because of Naïve theorem Approach. It's a probabilistic classifier supported Bayes'
 444 theorem with robust independence assumptions between the options. Naïve theorem classifier despite its
 445 simplicity, it astonishingly performs well and infrequently outperforms in advanced classification.
 446 Straightforward Naïve theorem will be enforced by plugging within the following main Bayes formula

447 $P(X_1, X_2, \dots, X_n | Y) = P(X_1 | Y) P(X_2 | Y) \dots P(X_n | Y) \quad (3)$
 448 The abovementioned Naïve theorem network produces a mathematical model, that is employed for
 449 modeling the sophisticated relations of random variables of un-wellness attributes and call outcome. The
 450 formula uses the formula to calculate chance with regard to un-wellness condition attributes worth and
 451 call attribute value supported by previous information, the formula classifies the choice attribute into
 452 labels allotted, and thus the conditional support is computed for every variable attribute [51].

453 Implementation of Bayesian Classification

454 The Naïve Bayes Classifier strategy is especially fit when the dimensionality of the sources of info is high.
 455 In spite of its effortlessness, Naive Bayes can frequently outflank increasingly advanced grouping
 456 techniques. Gullible Bayes model recognizes the attributes of patients with coronary illness. It
 457 demonstrates the likelihood of each information trait for the anticipated state.

458 Why favored Naive Bayes calculation

459 Credulous Bayes or Bayes' Rule is the reason for some, AI and information mining techniques. The
 460 standard (calculation) is utilized to make models with prescient abilities. It gives better approaches for
 461 investigating and getting information.

462 **Why preferred naive Bayes implementation:**

- 463 a. At the point when the information is high.
- 464 b. At the point when the properties are free of one another.
- 465 c. When we need increasingly proficient yield, when contrasted with different strategies yield

466 **Bayes Rule**

467 A restrictive likelihood is the probability of some end, C, given some proof/perception, E, where a reliance
 468 relationship exists among C and E.

469 This likelihood is meant as $P(C | E)$ where

470 $P(C/E) = P(E/C) P(C)/p(E)$

471 **3.1.3 K-NN – K-Nearest Neighbors**

472 K-NN is a kind of occasion based learning or apathetic realizing, where the capacity is just approximated
 473 locally and all calculation is conceded until characterization. K-NN arrangement, the yield is class
 474 participation. An article is ordered by a dominant part vote of its neighbors, with the item being doled out
 475 to the class most basic among its k closest neighbors (k is a positive whole number, normally little). In the
 476 event that $k = 1$, at that point the item is just appointed to the class of that solitary closest neighbor. K-
 477 Nearest Neighbors have been used in statistical estimation and pattern recognition i.e

478 If $K=1$, select the nearest neighbor

479 •If $K>1$, for classification select the most frequent neighbor, for regression calculate the average of K
 480 neighbors

X	Y	Distance
Attribute 1	Attribute 1	0
Attribute 1	Attribute2	1

481

482 $X=Y \Rightarrow D=0$

483 $X \neq Y \Rightarrow D=1$

484 **3.2 Experiments Data Set**

485 The information set for this analysis was taken from UCI data repository [49]. Information accessed from
 486 the UCI Machine Learning Repository is freely obtainable. This info contains seventy six attributes, and
 487 when neglecting redundant and unsuitable attributes, fifteen attributes were hand-picked. Below is that
 488 the list of fifteen attributes and their temporary description. Specially, the Cleveland, Hungarian,
 489 Switzerland and therefore the VA urban center databases are employed by several researchers and
 490 located to be appropriate for developing a mining model, attributable to lesser missing values and
 491 outliers. The information were cleansed and preprocessed before they were submitted to the planned
 492 algorithmic rule for coaching and testing. The 740 record sets were the valid instances for supervised
 493 machine-learning model building. The below shows the chosen vital risk factors from databases and their
 494 corresponding values Predictable attribute

495 1. Diagnosis (value 0: <50% diameter narrowing (no heart disease); value 1: >50% diameter narrowing
496 (has heart disease))

497 **Key attribute**

498 Patient Id – Patient's identification number

499 **Input attributes (Description of attributes)**

500 1. Age in Year

501 2. Sex (value 1: Male; value 0: Female)

502 3. Chest Pain Type (value 1: typical type 1 angina, value 2: typical type angina, value 3: non angina
503 pain; value 4: asymptomatic)

504 4. Fasting Blood Sugar (value 1: >120 mg/dl; value 0: <120 mg/dl)

505 5. Restecg – resting electrographic results (value 0: normal; value 1: having ST-T wave abnormality;
506 value 2: showing probable or definite left ventricular hypertrophy)

507 6. Exang - exercise induced angina (value 1: yes; value 0: no)

508 7. Slope – the slope of the peak exercise ST segment (value 1: unsloping; value 2: flat; value 3:
509 downsloping)

510 8. CA – number of major vessels colored by floursopy (value 0-3)

511 9. Thal (value 3: normal; value 6: fixed defect; value 7: reversible defect)

512 10. Trest Blood Pressure (mm Hg on admission to the hospital)

513 11. Serum Cholestrol (mg/dl)

514 12. Thalach – maximum heart rate achieved

515 13. Oldpeak – ST depression induced by exercise

516 14. Smoking – (value 1: past; value 2: current; value 3: never)

517 15. Obesity – (value 1: yes; value 0: no)Execution of Bayesian Classification

518 Attribute choice or feature sub-setting technique was applied for any reduction of information to form
519 patterns easier and comprehensible, however found negligible effects on performance measures of the
520 model engaged during this study. Visible of the above, all the thirteen attributes were taken into the
521 thought for developing a classifier's model and getting CVD prognostic outcome. The info mining
522 approach was used for evaluating the classification algorithms engaged and the DM tool was accustomed
523 to build the model. In these experiments, 10-fold cross-validations were utilized to partition the info set
524 into coaching and check sets; this fulfills the necessity of model training and testing purpose

525 **3.3 Data Source**

526 The publicly available heart disease database from Cleveland, Hungary, Switzerland and the VA Long
527 Beach Clinical databases [49] have aggregated enormous amounts of data about patients and their
528 ailments. The term Heart infection includes the assorted illnesses that influence the heart. Coronary
529 illness is the real reason for setbacks on the planet. Coronary illness kills one individual at regular
530 intervals in the United States. Coronary illness, Cardiomyopathy and Cardiovascular infection are a few
531 classifications of heart ailments. The expression "cardiovascular malady" incorporates a wide scope of
532 conditions that influence the heart and the veins and the way where blood is siphoned and coursed
533 through the body. Cardiovascular ailment (CVD) results in extreme disease, incapacity, and passing.
534

535 740 Record sets with therapeutic qualities will be gotten from a freely accessible database for coronary
536 illness from AI archive will be utilized, that is Cleveland, Hungary, Switzerland and the VA Long Beach
537 Heart Disease databases [49] with the assistance of the datasets, and the examples noteworthy to the
538 heart assault forecast are separated.

539 3.4. Processing and Analysis

540 The record sets were split into 2 datasets: coaching dataset and testing dataset. A complete 740 record
541 sets with fifteen medical attributes were obtained from the guts illness info. The records were split into 2
542 datasets like coaching dataset (296 record sets) and testing dataset (444 record sets). To avoid bias, the
543 records for every set were hand-picked willy-nilly in a very quantitative relation of 1 to 1.5.
544 In machine learning, a coaching set consists of associate degree input vector and a solution vector, and
545 is employed along with a supervised learning methodology to coach the information (e.g. decision tree,
546 KNN or a Naive Thomas Bayes classifier) employed by associate degree in AI machine. In a very dataset
547 a coaching set is enforced to make up a model, whereas a check (or validation) set is to validate the
548 model designed. Knowledge points within the coaching set are excluded from the check (validation) set.
549 When a model has been processed by victimization the coaching set checks the model by creating
550 predictions against the test set as a result of the information within the testing set already contained in the
551 celebrated values for the attribute to predict.

552
553 The table below shows the description of dataset selected for this work. The total record sets divided into
554 two with 13 and 15 attributes respectively.

555
556
557
558

Dataset	No. Of Attributes		Instances	Classes
Health Services Data	A	B	740	2
	13	15		

559 **Table 2.0 Dataset Description**
560

561 The model was developed and the first 13 input attributes were used then two more other attributes which
562 are **obesity and smoking** were added, as these attributes are considered as important attributes for
563 heart disease.

564 Also the deaths due to heart disease in many countries occur due to: work overload, mental stress and
565 many other problems, these are the other factor attributes we had considered in observing the prediction
566 change.

567 Most of the research papers referred upon have used 13 input attributes for prediction of Heart disease,
568 to get more appropriate results two more important attributes were added that is obesity and smoking.

569 Healthcare industry is generally "information rich", but unfortunately not all the data are mined which is
570 required for discovering hidden patterns & effective decision making- that's why we looked for more other
571 attributes which contribute to the heart disease

572 4. EXPERIMENTS AND RESULTS

573 The exhibition survey of a model for Heart Disease Prediction, utilizing Decision Trees, Naive Bayes, and
574 KNN displaying strategies were assessed concerning AI calculations. The targets of the trials were: To
575 break down the exhibition for the coronary illness expectation procedures, and portray how to improve
576 their forecast power, Efficient and precise in coronary illness forecast; To examine the centrality of
577 symptomatic highlights that best depict coronary illness information utilizing information mining strategies.

578 The Experiments demonstrated that the proposed technique gives the exact conclusion of coronary
579 illness than the current strategies

580 4.1 Experimental Setup

581 This exploration utilized classifiers given by Weka. The informational indexes were utilized as contribution
582 to three AI calculations; Naive Bayes (NB), K-Nearest Neighbor (KNN) and Decision Trees (DT). The
583 investigations began with 13 info properties and then 15 information traits esteems. Investigation results
584 were then exhibited in tables, broke down and deciphered as definite

585 4.2 Experimental Results and Analysis

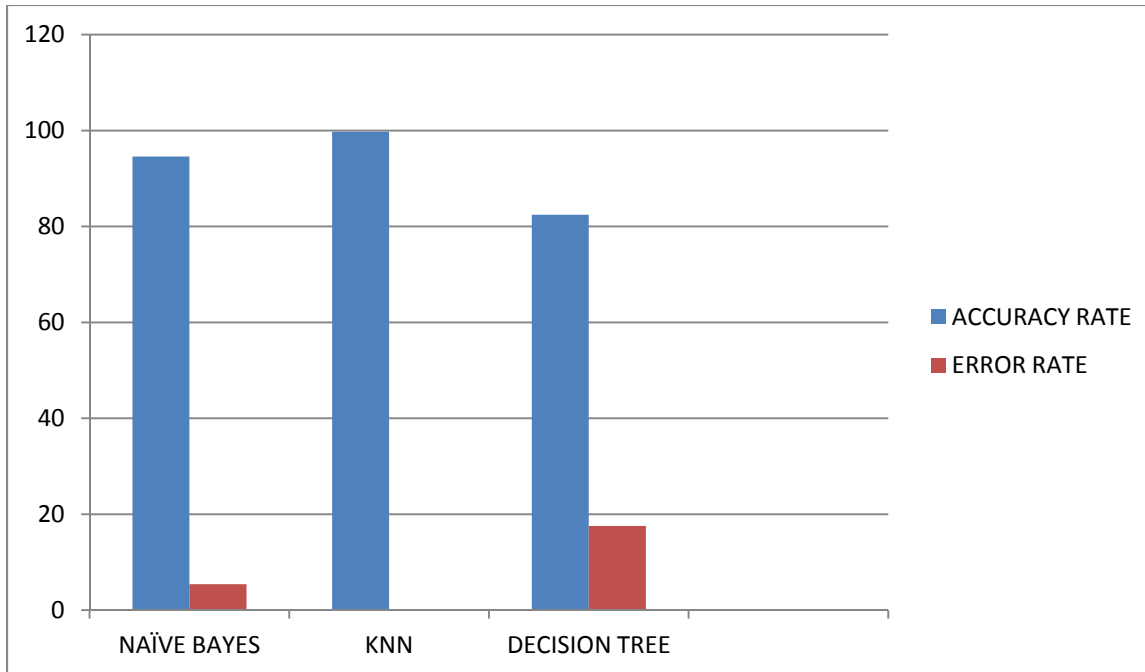
586 The test results and investigation accomplished for this examination was spoken to as in the tables
587 beneath. The exploration system has been clarified in the past area. For the tests, different information
588 mining grouping strategies were connected on the dataset. In this investigation, WEKA AI apparatus for
589 information mining was utilized to achieve the goals. The level of precision rate and mistake rate of
590 information mining Classification methods was utilized as the estimation parameters for investigation.
591 These parameters recommend that the classifier having a higher exactness rate and lower estimation of
592 blunder rate arrange the dataset in very amended way and the other way around. In this examination, the
593 information was right off the bat isolated into preparing information and testing information. The
594 preparation set was utilized to build the classifier and test set utilized for approval. In this examination, the
595 level of dataset utilized for preparing and testing information were 40% and 60% individually. At that point,
596 the 10 overlay cross approval technique was connected to create the classifiers utilizing recently
597 referenced AI apparatuses. At last the outcomes were archived as far as precision rate and mistake rates.
598

599 The table beneath Displays the results for classification techniques applied on health facility services data
600 in WEKA. Considering accuracy and error rates as performance measure the classification techniques
601 with highest accuracy are obtained for health facility Services data in given different techniques used.
602

603 **Table 3.0 Results Using WEKA API**
604

Technique Used	Accuracy Rate		Error Rate	
	13 Attributes	15 Attributes	13 Attributes	15 Attributes
Naive Bayes	90.76	94.59	9.24	5.41
Decision Tree	97.07	99.77	2.93	0.23
KNN	79.28	82.43	20.72	17.57

605
606 The graph below displays the performance analysis of classification techniques for 15 attributes using
607 WEKA. The best classifier for this particular data set will then be chosen.
608



609
610 **Fig 4.0 Performance analysis of classification techniques using WEKA API**

611 **4.3. Results**

612 The dataset comprised of all **740 Record sets** in Heart illness database. The records were then divide into
 613 two, one utilized for preparing comprises of 296 records and another for testing comprises of 444 records.
 614 The information mining apparatus Weka 3.6.6 was utilized for trial. At first dataset contained a few fields,
 615 in which some incentive in the records was absent. These were recognized and supplanted with most
 616 fitting qualities utilizing Replace Missing Values channel from Weka 3.6.6. The Replace Missing Values
 617 channel checks all records and replaces missing qualities with mean mode technique. This procedure is
 618 known as Data Pre-Processing. After pre-handling the information, information mining order procedures,
 619 for example, KNN, Decision Trees, and Naive Bayes were connected. A disarray lattice is acquired to
 620 figure the exactness of arrangement. A perplexity grid demonstrates what number of occurrences has
 621 been doled out to each class. In our analysis we have two classes, and in this manner we have a 2x2
 622 perplexity network

623 Class A= YES (Has coronary illness)

624 Class B = (No coronary illness)

625 **Table 4.0 a Disarray Network**

	A(Has heart disease)	B(Has no heart disease)
A(has heart disease)	TP	FN
B(has no heart disease)	FP	TN

626

627 TP (True Positive): It indicates the quantity of records named genuine while they were in reality evident.
 628 FN (False Negative): It signifies the quantity of records delegated false while they were in reality evident.
 629 FP (False Positive): It indicates the quantity of records named genuine while they were in reality false. TN
 630 (True Negative): It means the quantity of records named false while they were in reality false. Results got
 631 with 13 properties are determined beneath

632 **Table 3.3 Confusion Networks Got For Three Arrangement Techniques with 13 Qualities**

633 **Confusion matrix for Naive Bayes:**

	A	B
A	182	13
B	28	221

634

635 **Confusion matrix for Decision Trees:**

	A	B
A	205	6
B	7	226

636

637 **Confusion matrix for KNN:**

	A	B
A	160	30
B	62	192

638

639 Results obtained by adding two more attributes i.e. total 15 attributes are specified below.

640 Table 3.4 Confusion matrixes obtained for three classification methods with 15 attributes

641 **Confusion matrix for Naive Bayes:**

	A	B
A	187	11
B	13	233

642

643 **Confusion matrix for Decision Trees:**

	A	B
A	168	0
B	1	275

644

645 **Confusion matrix for KNN**

	A	B
A	153	36
B	42	213

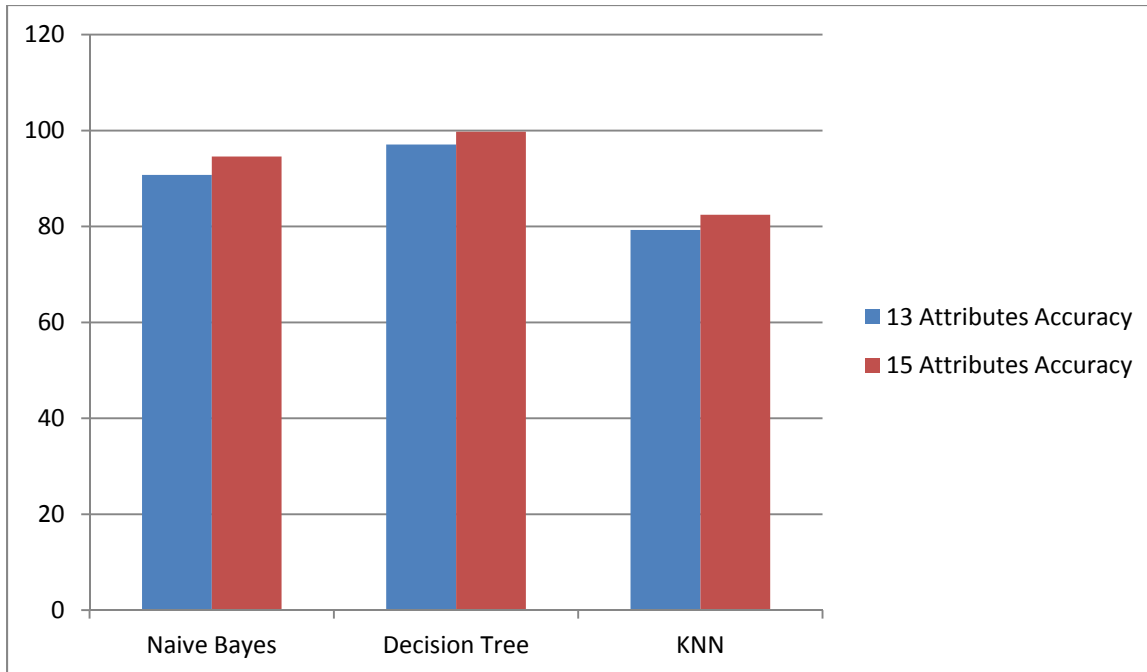
646

647 **Table 5.0 shows accuracy for different classification methods with 13 input attributes and 15 input**
648 **attributes values.**

Classification Techniques	Accuracy with	
	13 Attributes	15 Attributes
Naive Bayes	90.76	94.59
Decision Tree	97.07	99.77
KNN	79.28	82.43

649

650 The accuracy of each of the method is plotted on a graph as below:



651

652 **Figure 5.0: Graphical representation of accuracy for each method.**

653 The performance and accuracy of every experiment are evaluated through performance measures like
 654 true positive rate, precision, F-measure, receiver in operation characteristic (ROC) space, letter statistics
 655 and root mean sq. (RMS) error. Identical measures are used for comparative analysis of enforced
 656 algorithms. Once the experiments, subsequent step is to match algorithms employed in these
 657 experiments for lightness the most effective one in terms of un-wellness prediction chance and classifier's
 658 accuracy. Having a glance at the results, it becomes apparent that the goal to supply AN ensemble
 659 classifier for early diagnostic screening with needed level of accuracy is triple-crown. A correlation
 660 between accuracy and therefore the quantity of attributes employed in the creation of the classifier was
 661 found. In general, additional attributes offer larger accuracy as visualized by results. With relation to
 662 mythical creature space as performance live, AN optimal/perfect classifier can score one on this take a
 663 look at, therefore this will build our results trying less dimmed with results quite 0.9 mythical creature for
 664 all classifiers. The comparative performance outline of enforced algorithms is given in table above.
 665 In general, the results of all the enforced rules are far better by all algorithms with specially the choice
 666 tree algorithm leading in accuracy and prediction chance. The accuracy of enforced algorithms on the
 667 given heart condition knowledge set is given within the table given above, and therefore the lowest
 668 accuracy is 84.43% for KNN analysis and therefore the highest accuracy is 98.17% for the choice rule
 669 supported on the top mentioned results and comparisons with relation to the chosen performance
 670 measures, the naïve Bayes and decision tree performed well and every rule with quite 94 prediction
 671 chance increased responsibility of the prediction system. Additional stress is given to pick out the
 672 algorithms having high true positive rate, as being the core live for early designation of heart condition

673

674

675

676 5. CONCLUSION AND FUTURE WORK

677 5.1 Knowledge Contributions

678 This research that proposed the use of a model for Heart Disease Prediction using Data Mining
679 Classification Techniques provided a set of contributions that can be summarized while considering
680 different points of view. On the more theoretical and modeling side, heart disease model for prediction
681 analysis was proposed.

682 On the implementation side, this research improved results on accuracy with increase in number of
683 attributes. This is supported by the high levels of classification accuracy exhibited when data sets that
684 were used showed that there is increase in classification accuracy as the number of the attributes used
685 for testing increased.

686 5.2 Conclusion

687 This approach-based paradigm for cardiopathy prediction model has been projected as a system
688 whereas utilizing Naïve Bayesian, decision tree and KNN classifiers. The projected system is GUI-based,
689 easy, scalable, reliable and expandable system, that has been enforced on the maori hen platform. The
690 projected operating model can even facilitate in reducing treatment prices by providing Initial medical
691 specialty in time. The model can even serve the aim of coaching tool for medical students and can be a
692 soft diagnostic tool obtainable for MD and heart specialist. General physicians will utilize this tool for initial
693 diagnosing of cardio patients. Various information mining characterization procedures were connected on
694 the particular dataset, the order procedure inside the framework model is performed with traits like age,
695 sex, heart beat rate, cholesterol level and so on. The expectation is then made dependent on this
696 arrangement results. Here the AI ability of the PC framework can be stretched out into the medicinal field.
697 The proposed framework model is best for lessening the blunder event during the illness expectation. In
698 this examination the exactness and precision of three unique classifiers are estimated, the outcome
699 demonstrates choice tree arrangement has high precision and less mistake rate, Naïve Bayer
700 characterization strategy creates preferred outcome over KNN grouping. This investigation can assist
701 scientists with getting productive outcomes in the wake of knowing the best order strategy for this specific
702 dataset. The general target of the examination was to foresee all the more precisely the nearness of
703 coronary illness. In this exploration, more information characteristics weight and smoking were utilized to
704 get progressively precise outcomes.

705 5.3 Future Work

706 Heart Disease Prediction using Data Mining Classification Techniques can be used largely in hospital
707 based sectors for disease prediction, However, there is need for more research to be done on contextual
708 knowledge being incorporated as part of feature selection and model creation for specific domains where
709 precise context, which does not depend on attributes needs to be used in learning and prediction is
710 required also. There is need to experiment the prediction models with real live testing of heart disease.
711 This research can also be enhanced by experiment with more attributes in training and testing data sets.
712 There are many possible improvements that could be explored to improve the scalability and accuracy of
713 this prediction system. As we have developed a generalized system, in future we can use this system for
714 the analysis of different data sets. The performance of the health's diagnosis can be improved
715 significantly by handling numerous class labels in the prediction process, and it can be another positive
716 direction of research. In DM warehouse, generally, the dimensionality of the heart database is high, so
717 identification and selection of significant attributes for better diagnosis of heart disease are very
718 challenging tasks for future research.

719

720 **REFERENCES**

- 721 1. Abdullah H. Wahbeh, "A Comparison Study between Data Mining Tools over some Classification
722 Methods" (IJACSA) International Journal of Advanced Computer Science and Applications, Special
723 Issue on Artificial Intelligence, vol. 3, no. 2, p 18-26, 2012.
- 724 2. Abhishek Taneja, Heart Disease Prediction System Using Data Mining Techniques; Oriental Journal
725 of computer science & Technology ISSN: 0974-6471 December2013.
- 726 3. Aditya Methaila, Early Heart Disease Prediction Using Data Mining Techniques; CCSEIT, DMDB,
727 ICBB, MoWiN, AIAP pp. 53–59, 2014.
- 728 4. Al-Radaideh "Using data mining techniques to build a classification model for predicting employee's
729 performance", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol.
730 3, No. 2, pp.60-71, 2014.
- 731 5. B. Venkatalakshmi and M. Shivsankar, "Heart disease diagnosis using predictive data mining,"
732 International Journal of Innovative Research in Science, Engineering and Technology, vol. 3, no. 3,
733 pp. 1873–1877, 2014.
- 734 6. Beant Kaur and Williamjeet Singh., "Review on Heart Disease Prediction System using Data Mining
735 Techniques", IJRITCC , pp. 56-72, October 2014.
- 736 7. Blake, C.L., Mertz, C.J. "UCI Machine Learning Databases"
- 737 8. Chaitrali S. Dangare, Sulabha S. Apte, —Improved Study of Heart Disease Prediction System using
738 Data Mining Classification Techniques; International Journal of Computer Applications (0975 – 888)
739 Volume 47– No.10, June 2012
- 740 9. Chaitrali S. Danagre, Sulabha S. Apte, Ph.D, Improved Study of Heart Disease Prediction System using
741 Data mining Classification Techniques, IJCA, June 2012.
- 742 10. Charly, K.: "Data Mining for the Enterprise", 31st Annual Hawaii Int. Conf. on System Sciences, IEEE
743 Computer, 7, 295-304, 2014.
- 744 11. Choi Keunho et al. "Classification and Sequential Pattern Analysis for Improving Managerial
745 Efficiency and Providing Better Medical Service in Public Healthcare Centers" health inform res,
746 pp.67-76, June 2014
- 747 12. D.K, "Classification of women health disease (Fibroid) using decision tree algorithm", International
748 Journal of Computer Applications in Engineering Science Vol.2, Issue 3, pp.84, September2016,]
- 749 13. Fartash. Haghanikhameneh "A Comparison Study between Data Mining Algorithms over
750 Classification Techniques in Squid Dataset" International Journal of Artificial Intelligence, Autumn
751 (October) 2015, Vol. 9, pp66-68.
- 752 14. Fayyad, Piatetsky-Shapiro, Smyth, "From Data Mining to Knowledge Discovery: An Overview", in
753 Fayyad, Piatetsky-Shapiro, Smyth, Uthurusamy, *Advances in Knowledge Discovery and Data Mining*,
754 AAAI Press / The MIT Press, Menlo Park, CA, 2014, pp. 1-34.
- 755 15. Garchchopogh et al, "Application of decision tree algorithm for data mining in healthcare operations:
756 A case study", International Journal of Computer Applications Vol 52 – No. 6, August 2014, pp.567-
757 280.
- 758 16. Global Atlas on Cardiovascular Disease Prevention and Control (PDF). World Health Organization in
759 collaboration with the World Heart Federation and the World Stroke Organization. pp. 3–18. ISBN
760 978-92-4-156437-3 September 2011.
- 761 17. Han, J. and Kamber, M. (2014). *Data Mining: Concepts and Techniques*. fourth Edition, Morgan
762 Kaufmann Publishers, San Francisco Vol. 16, No. 3, 2013, pp. 291-296
- 763 18. Hearty "Analysis of meal patterns with the use of supervised data mining techniques-Artificial Neural
764 Network and Decision Tree", The American Journal of Clinical Nutrition Vol. 18, No. 3, 2013, pp. 192-
765 190, 2015
- 766 19. Hlaudi Daniel Masethe, Mosima Anna Masethe-prediction of Heart Disease using Classification
767 Algorithms; Proceedings of the World Congress on Engineering and Computer Science 2014.
- 768 20. Ho, T. J.: *Data Mining and Data Warehousing*, Prentice Hall, 2016, pp.66-69.
- 769 21. Huang, Li, Su, Watts, & Chen, 2007; Ishibuchi, Kuwajima, Nojima, 2007; Karabatak & Ince, 2009;
770 Shin et al., 2010; Wang & Hoy, 2015, pp256-267.
- 771 22. Jabbar et al "Heart disease prediction system using associative classification and Genetic Algorithm",
772 International Conference on Emerging Trends in Electrical, Electronics and Communication
773 Technologies-ICECIT, 2015

- 774 23. Jyoti Soni et.al. Predictive Data Mining for Medical Diagnosis: An Overview of Heart Disease
775 Prediction; International Journal of Computer Applications (0975 – 8887) Volume 17– No.8, March
776 2011.
- 777 24. Koç et al, “A comparative study of artificial neural network and logistic regression for classification of
778 marketing campaign results”, Mathematical and Computational Applications, Vol. 18, No. 3, 2013, pp.
779 392-398
- 780 25. Mrs.G.Subbalakshmi, “Decision Support in Heart Disease Prediction System using Naive Bayes”,
781 Indian Journal of Computer Science and Engineering. Vol. 3, No. 5, May 2014, pp.227-227-238.
- 782 26. Nakul Soni, Chirag Gandhi, “Application of data mining to health care”, International Journal of
783 Computer Science and its Applications Volume 36– No.10,vol.5 June 2014,
- 784 27. Nidhi Bhatla, Kiran Jyoti, “ An Analysis of Heart Disease Prediction using Different Data Mining
785 Techniques” International Journal of Engineering and Technology Vol.1 issue 8 2012, pp.234-241..
- 786 28. Obenshain, M.K: “Application of Data Mining Techniques to Healthcare Data”, Infection Control and
787 Hospital Epidemiology, 25(8), 690–695, 2014
- 788 29. Patil Dipti “An adaptive parameter for data mining approach for healthcare applications” (IJACSA)
789 International Journal of Advanced Computer Science and Applications, Vol. 3, No. 1, 2014, pp.66.70..
- 790 30. Pushpalata Pujari “ Classification and comparative study of data mining classifiers with feature
791 selection on binomial data set” Journal of Global Research in Computer Science, Vol. 3, No. 5, May
792 2016, pp.675-682.
- 793 31. R. Chitra, Review Of Heart Disease Prediction System Using Data Mining And Hybrid Intelligent
794 Techniques; Ictact Journal On Soft Computing, July 2013,volume: 03, Issue: 04, pp781-785.
- 795 32. R.Wu, W.Peters, M.W.Morgan, “The Next Generation Clinical Decision Support: Linking Evidence to
796 Best Practice”, *Journal of Healthcare Information Management.* 16(4), pp. 50 55, 2016.
- 797 33. S.Asha Rani and Dr.S.Hari Ganesh, “ A comparative study of classification algorithm on blood
798 transfusion” International Journal of Advancements in Research & Technology, Volume 3, Issue 6,
799 June-2014, pp.56-63.
- 800 34. Saichanma et al. "The Observation Report of Red Blood Cell Morphology in Thailand Teenager by
801 Using Data Mining Technique." Advances in hematology, 2014 pp.104-109.
- 802 35. Sakshi and Prof.Sunil Khare “A Comparative Analysis of Classification Techniques on Categorical
803 Data in Data Mining” International Journal on Recent and Innovation Trends in Computing and
804 Communication Vol. 3 Issue: 8,pp.5142 – 5147
- 805 36. Sayad AT, Halkarnikar PP. Diagnosis of heart disease using neural network approach. Int J Adv Sci
806 Eng Technol. 2014;2:88–92.
- 807 37. Setiawan, *et al,* “A Comparative Study of Imputation Methods to Predict Missing Attribute Values in
808 Coronary Heart Disease Data Set”, Journal in Department of Electrical and Electronic
809 Engineering,Vol.21, PP. 266–269, 2008
- 810 38. Shadab Adam Pattekari and Asma Parveen, prediction system for heart disease using naïve bayes,
811 International Journal of Advanced Computer and Mathematical Sciences, 2012, pp.476-484.
- 812 39. Shanthi Mendis; Pekka Puska; Bo Norrving; World Health Organization (2011).
- 813 40. Shelly Gupta et al. “Performance Analysis of Various Data Mining Classification Techniques on
814 Healthcare Data” International Journal of Computer Science & Information Technology (IJCSIT) Vol
815 3, No 4, August 2011,pp.877-892.
- 816 41. Sundar et al. “Performance analysis of classification data mining techniques over heart disease
817 database”, [IJESAT] International Journal of Engineering Science and Advanced Technology,
818 Volume-2, Issue-3,pp. 470 – 478,2013.
- 819 42. Tang, Z. H., MacLennan, J.: *Data Mining with SQL Server 2005*, Indianapolis: Wiley, 2015, pp.445-
820 450.
- 821 43. Tariq O. Fadl Elsid and Mergani. A. Eltahir “An Empirical Study of the Applications of Classification
822 Techniques in Students Database” Int. Journal of Engineering Research and Applications ISSN:
823 2248-9622, Vol. 4, Issue 10(Part - 6), pp.01-10, October 2014
- 824 44. Thuraisingham, B.: “A Primer for Understanding and Applying Data Mining”, IT Professional, 28-31,
825 2015
- 826 45. Umadevi, D.Sundar, Dr.P.Alli, “A Study on Stock Market Analysis for Stock Selection – Naïve
827 Investors’ Perspective using Data Mining Technique”, International Journal of Computer Applications
828 (0975 – 8887), Vol 34– No.3,2011.

- 829 46. V. Manikandan and S. Latha, "Predicting the Analysis of Heart Disease Symptoms Using Medical
830 Data Mining Methods "International Journal of Advanced Computer Theory and Engineering", Vol. 2,
831 Issue. 2,pp.236-240, 2013.
- 832 47. Vikas Chaurasia, et al, Early Prediction of Heart Diseases Using Data Mining Techniques; Caribbean
833 Journal of Science and Technology ISSN 0799-3757, Vol.1,208-217, 2013.
- 834 48. World Health Organization; Cardiovascular Diseases (CVDs) Fact Sheet Reviewed June 2016
- 835 49. Cleveland, Hungary, Switzerland, & VA Long Beach Database: <http://archive.ics.uci.edu/ml/datasets/Heart+Disease>
- 836 50. Nadali, A; Kakhky, E.N.; Nosratabadi, H.E., "Evaluating the success level of data mining projects based on
837 CRISP-DM methodology by a Fuzzy expert system," Electronics Computer Technology (ICECT), 2011 3rd
838 International Conference on , vol.6, no., pp.161,165, 8-10 April 2011
- 839 51. Dangare CS, Apte SS. Improved Study of Heart Disease Prediction System using Data Mining
840 Classification Techniques. *Int J Comput Appl.* 2012;47(10):44–48.

841
842

843

844

845

UNDER PEER REVIEW