

Analysis of Heart Disease using in Data Mining Tools Orange and Weka

Sarangam Kodati¹

¹ Sri Satya Sai University of Technology and Medical Science

Received: 9 December 2017 Accepted: 2 January 2018 Published: 15 January 2018

Abstract

Health care is an inevitable task to be done in human life. Health concern business has become a notable field in the wide spread area of medical science. Health care industry contains large amount of data and hidden information. Effective decisions are made with this hidden information by applying patient; however, with data mining these tests could be reduced. But there is a lack of analyzing tool according to provide effective test outcomes together with the hidden information, so and such system is developed using data mining algorithms for classifying the data and to detect the heart diseases. Data mining acts so a solution by many healthcare problems. Naïve Bayes, SVM, Random Forest, KNN algorithm is one such data mining method which serves with the diagnosis regarding heart diseases patient. This paper analyzes few parameters and predicts heart diseases, thereby suggests a heart diseases prediction system (HDPS) based total on the data mining approaches

Index terms— data mining, weka, orange, heart disease, data mining classification techniques.

1 Analysis of Heart Disease using in Data Mining

Tools Orange and Weka Sarangam Kodati ? & Dr. R. Vivekanandam ? Abstract-Health care is an inevitable task to be done in human life. Health concern business has become a notable field in the wide spread area of medical science. Health care industry contains large amount of data and hidden information.

Effective decisions are made with this hidden information by applying patient; however, with data mining these tests could be reduced. But there is a lack of analyzing tool according to provide effective test outcomes together with the hidden information, so and such system is developed using data mining algorithms for classifying the data and to detect the heart diseases. Data mining acts so a solution by many healthcare problems. Naïve Bayes, SVM, Random Forest, KNN algorithm is one such data mining method which serves with the diagnosis regarding heart diseases patient. This paper analyzes few parameters and predicts heart diseases, thereby suggests a heart diseases prediction system (HDPS) based total on the data mining approaches.

Keywords: data mining, weka, orange, heart disease, data mining classification techniques.

I.

2 Data Mining

Data mining is concerned together with the method of computationally extracting unknown knowledge from vast sets of data. Extraction of useful knowledge from the enormous data sets and providing decision-making results for the diagnosis or remedy of diseases is very important. Data mining can stand used to extract knowledge by analyzing and predicting some diseases. Health care data mining has a large potential according to discover the hidden patterns among the data sets about the medical domain. Various data mining methods are available with their suitability dependent on the healthcare data. Data mining applications in health care can have a wonderful potential and effectiveness. It automates the process of finding predictive information in large databases. Disease prediction plays an important role in data mining. Finding of heart disease requires the performance of some

43 tests on the patient. However, use of data mining techniques can reduce the number of tests. This reduced test
44 set plays a significant role in performance and time. Health care data mining is an important task because it
45 allows doctors to see which attributes are more important for diagnosis such as age, weight, symptoms, etc. This
46 will help the doctors diagnose the disease more efficiently. Knowledge discovery in databases is the method of
47 finding useful information and patterns into data. Knowledge discovery within databases can be do using data
48 mining. It makes use of algorithms after extract the information and patterns derived by the knowledge discovery
49 in databases process. Various stages of knowledge discovery in databases process are highlighted in Fig. ??.

3 Fig. 1: KDD Process

51 Various stages concerning knowledge discovery of databases method are described as follows. In Selection stage,
52 that obtains the different data resources. In preprocessing stage, it removed the unwanted missing and noisy
53 data and furnished the clean data which execute format in accordance including a common format of transform
54 stage. Then data mining techniques are applied according to get desired output. Finally into the between the
55 signification stage, that will present the result after end user in a meaningful manner. instead of predicting
56 the discrete class the outcome is a numeric value. c) Association rule mining: The association and patterns
57 between the some attributes are extracted or from its attributes, rules are created. The rules and patterns are
58 used predicting the categories or classification of the test data. d) Clustering: The grouping of similar instances
59 into clusters takes place. The challenges or drawbacks considering this type of machine learning is that we have
60 according to first identify clusters and assign a new instance according to these clusters [8].

4 II.

5 Data Mining Techniques

63 Out of this four types of learning methods, we need to identify the algorithm as performs better. The application
64 of data mining methods depends on the types of data which is fitted to be used in the techniques, or solving data
65 mining troubles depend on the types of data to stand used and the selection about data mining technique which
66 is most suitable for the data used.

6 III.

7 Machine Learning

69 Machine learning (ML), employed as like a method in data science, is the process of programming computers after
70 learning from past experiences (Mitchell, 1997Machine Learning seeks to develop algorithms to that amount learn
71 out of data directly with little or no human intervention. Machine Learning algorithms perform a range of tasks
72 such so like prediction, classification, or decision making. Machine Learning stems from artificial intelligence
73 research and has become an essential aspect of data science. Machine learning begins with input so a training
74 data set. In this phase, the Machine Learning algorithm employs the training dataset after learning from the data
75 and structure patterns. The learning phase outputs a model so much is used by way of the testing phase. The
76 testing phase employs any other dataset, applies the model from the training phase, and results are presented for
77 analysis. The overall performance regarding the test dataset demonstrates the model's ability in conformity with
78 performing its task against data. Machine learning extends beyond a statically coded set regarding statements
79 into statements, so a lot are dynamically generated based as regards the input data.

80 IV.

8 Open Source Softwares

82 Open source has, in the minds regarding many, come to be synonymous with free software ??Walters, 2007).
83 Open source software is software where the development then the source code are made publically available and
84 designed after denying everyone the right according to exploit the software ??Laurent, 2004). Open source general
85 refers in conformity with the source code concerning the application being freely and openly available because
86 of modifications. Two such examples of open source licenses are the GPL, or general people consent (GNU.org,
87 2015a), then GNU(GNU.org, 2015b). Anyone be able to develop extensions then customizations about open
88 source software; though, charging a fee for certain things to do is typically prohibited by using a public license
89 agreement whereby any modifications to the source code automatically become public domain. Communities
90 emerge around software with developers worldwide extending open source software.

91 V.

9 Heart Diseases

93 The highest mortality in both India and abroad is due to heart disease. So it is vital time to check this death
94 toll by correctly identifying the disease between initial stage. The matter becomes a headache for all medical
95 doctors both in India and abroad. Nowadays doctors are adopting many scientific technologies and methodology
96 for both identifications or diagnosing not only the common disease but also many fatal diseases. The successful
97 treatment is continually attributed to right and accurate diagnosis. Doctors may also sometimes fail to take

98 accurate decisions while diagnosing the heart disease about a patient, therefore heart disease prediction systems
99 which use machine learning algorithms assist in such cases to get accurate results [1].

100 **10 VI.**

101 **11 Heart Disease Dataset**

102 The dataset used for this work is from UCI Machine Learning repository from which the Cleveland heart disease
103 dataset is used. The dataset has 303 instance and 76 attributes. However, only 14 attributes are used of this
104 paper. These 14 attributes are the consider factors for the heart disease prediction [8]. Even though it has 303
105 instances as only 297 are completed and the remaining rows contained missing values and removed out of the
106 experiment.

107 **12 VII. Overview of Data Mining Tools**

108 Data mining has a wide number of applications ranging from marketing and advertising about goods, functions
109 and products, artificial intelligence research, biological sciences, crime investigations to high-level government
110 intelligence. Due to its widespread usage and complexity involved in building information mining applications,
111 a vast number of Data mining tools hold been developed over decades. Every tool has its advantages and
112 disadvantages. [6] Within data mining, there is a group of tools that have been developed by a research community
113 and data analysis enthusiasts; he are provided free of the price using one on the existing open-source licenses. An
114 open-source development model means that the tool is a result of a community effort, not necessarily supported
115 by a single

116 **13 Global Journal of Computer Science and Technology**

117 Volume XVIII Issue I Version I 18 Year 2 018 () C organization but alternatively the result regarding contributions
118 from an international and informal development team. This development style affords a means on incorporating
119 the various experiences Data boring gives many excavation techniques according to extract data from databases.
120 Data mining tools predict future trends, behaviors, allowing business according to make proactive, knowledge-
121 driven decisions. The development and application concerning data mining algorithms require the use of very
122 powerful software tools. As the number of accessible tools continues by grow the choice of the most suitable tool
123 becomes increasingly difficult. [6] The top 6 open source tools available because data mining is briefed as below.

124 Data mining tools like Weka and Orange are used to perform various data mining techniques. The first step
125 of the methodology consists of selecting a number of available open source data mining tools in accordance with
126 being tested. Many open data mining tools are available for free on the Web. After surfing the Internet, some
127 tools were chosen; including the Waikato Environment for Knowledge Analysis (WEKA) durability and Orange
128 Canvas.

129 **14 VIII.**

130 **15 Weka**

131 The Waikato Environment for Knowledge Analysis (WEKA) [7] is an open source software and machine learning
132 toolkit introduced by Waikato University, New Zealand. WEKA helps several standard data mining tasks as
133 data preprocessing, clustering, classification, regression, visualization and feature selection New algorithms can
134 also be implemented the usage concerning WEKA with existing data mining and machine learning techniques.
135 WEKA gives a number sources because loading data, which include files, URLs then databases. It helps file
136 formats include WEKA's own ARFF format, CSV, Lib SVMs format, and C4.5's format. Many evaluation
137 criteria are also provided of WEKA certain as confusion matrix, precision, recall, true positive and false negative,
138 etc. Some of the advantages of WEKA tool includes Open source, platform independent and portable, graphical
139 user interface and contains a very vast collection of different data mining algorithms.

140 **16 IX.**

141 **17 Orange**

142 Orange is an open source machine learning technology or data mining software. Orange can be used for explorative
143 data analysis and visualization [3]. It gives a platform for experiment selection, predictive modeling, and
144 recommendation systems and can be used of genomic research, biomedicine, bioinformatics, and teaching. Orange
145 is always preferred when the factor of innovation, quality, or reliability is involved[10], [4].

146 X.

147 **18 The Comparative Study**

148 The methodology of the study constitutes regarding collecting a set of free data mining and knowledge discovery
149 tools according to be tested, specifying the data sets to be used, and selecting a set of classification algorithm

150 according to test the tools' performance. Demonstrates the overall methodology followed for fulfilling the goal of
151 its research.

152 **19 b) Recall**

153 It is defined as the average probability of complete retrieval. $\text{Recall} = \frac{\text{True positives}}{\text{True positives} + \text{False negative}}$
154

155 **20 c) Navie Bayes**

156 When the dimensionality of the inputs is high, the Naïve Bayes Classifier method is particularly suited. The
157 problem including the Naïve Bayes Classifier is so that assumes all attributes are independent on each other
158 which in general cannot be applied. Naive Bayes is harder to debug and understandable [2]. Naive Bayes
159 used into robotics and computer vision. In naive Bayes, decision tree perform poorly. Comparative analysis
160 of precession and recall analyzing for heart disease data sets precession in Orange 82.4% and Recall 80.6%. In
161 WEKA precession 83.7% and Recall 83.7 %.Compare to Orange tool and WEKA, weka is best precession and
162 Recall.

163 **21 d) Support Vector Machine**

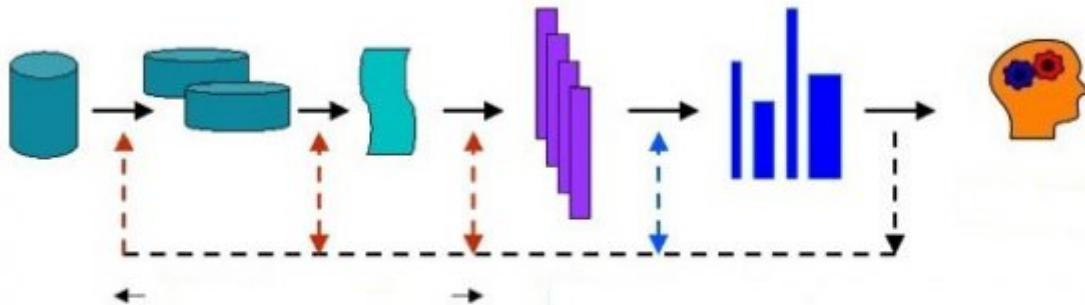
164 Support Vector Machines proved themselves to be very fine into a variety of pattern classification tasks and
165 accordingly received a great deal of attention recently. Support vector machine is a supervised machine learning
166 technique. The SVM algorithm predicts the occurrence about heart disease by ability on plotting the disease
167 predicting attributes regarding the multidimensional hyperplane or classifies the classes optimally by creating
168 the approach between two data clusters [5]. This algorithm attains high accuracy by the use regarding nonlinear
169 features called kernels. Comparative analysis of precession and recall analyzing for heart disease data sets
170 precession in Orange 81.7% and Recall 70.5%. In WEKA precession 81.8% and Recall 81.9 %.Compare to
171 Orange tool and WEKA, weka is best precession and Recall.

172 **22 e) Random Forest**

173 Random Forest is essentially an ensemble of unpruned classification trees. It gives excellent performance
174 concerning a number about practical problems, largely because such is not sensitive to noise in the dataset,
175 and it is not subject to overfitting. It works fast and generally exhibits a substantial performance improvement
176 over many other tree-based algorithms. Random forests are built by combining the predictions on a number of
177 trees, each of which is trained within isolation. There are three main choices to stand performed when constructing
178 a random tree. Comparative analysis of precession and recall analyzing for heart disease data sets precession in
179 Orange 77.9% and Recall 73.4%. In WEKA precession 81.8% and Recall 81.9 %.Compare to Orange tool and
180 WEKA, weka is best precession and Recall.

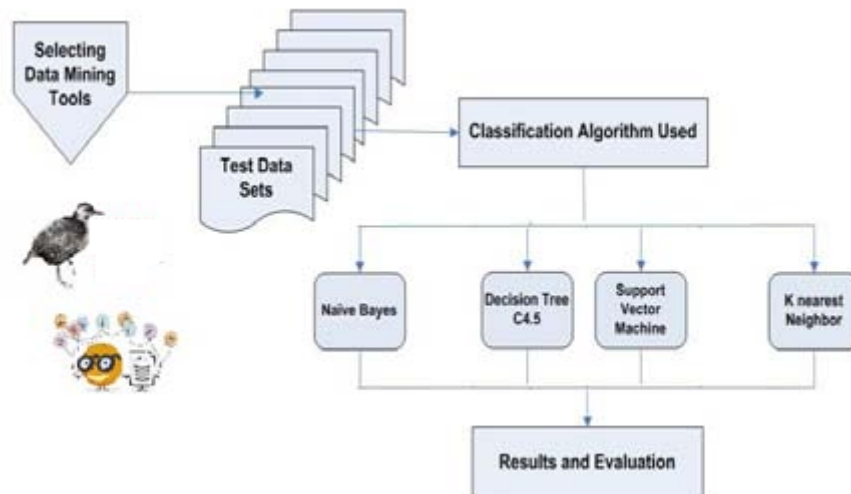
181 **23 f) KNN Classifier**

182 K-nearest neighbor is a sophisticated approach for classification that finds a group of K objects in the training
183 documents that are close to the test value. To classify an unlabeled object, the distance between it object and
184 labeled object is computed and it's K nearest neighbors are identified. Classification accuracy commonly depends
185 of the choice value of K and will be better than that of using the nearest neighbor classifier [9]. For vast data sets,
186 K can be larger to reduce the error. Choosing K can be done experimentally, where a number concerning patterns
187 taken out from the training set can be categorised using the remaining training patterns for different values over
188 k. The value of K which gives the least error in classification will be chosen. If same class is shared in various of
189 K-nearest neighbors, then per-neighbor weights of as class are added together, and the resulting weighted sum is
190 used as the likelihood score of that class with respect to the test document [8].Comparative analysis of precession
191 and recall analyzing KNN for heart disease data sets precession in Orange 58% and Recall 54.7%. In WEKA
192 precession 75.3% and Recall 75.2 %. Compare to Orange tool and WEKA weka is best precession and Recall.
193



2

Figure 1: Fig. 2 :



3

Figure 2: Fig. 3 :

Data Fusion

Data Classification

- -
noising

D

The most frequently used Data Mining techniques are specified below: a) Classification

a set of classified examples (training set) and uses it for training the algorithms. With the trained algorithms, classification of the test data takes place based over the patterns and r

Author ? : Research Scholar, Department of Computer Science and Engineering, Sri Satya Sai University of Technology and Medical Science, Sehore, Bhopal, Madhya Pradesh, India.

e-mail: k.sarangam@gmail.com

Author ? : Professor, Director in Muthayamal Engineering College, Namakkal, India.

b) Numeric predication: This is a variant of classification learning with the excepti

© 2018 Global Journals

Figure 3: Raw Data Target Data Preprocessed Data Transformed Data Patterns Knowledge
Data preprocessing Pattern Recognition Interpreting Results Knowledge

1

Figure 4: Table 1 :

showing best classification

Algorithm classification Average	Precession in Orange	Recall in Orange	Precession in WEKA	Recall in WEKA
Naïve base classifier	0.824	0.806	0.837	0.837
SVM or Support Vector Machine	0.817	0.705	0.84	0.8365
Random Forest	0.779	0.734	0.818	0.819
1BK or K-Nearest Neighbor	0.58	0.547	0.753	0.752

Figure 5:

-
- 194 [Mikut and Wiley (2011)] , Ralf Mikut , Markus Reischl Wiley . *Interdisciplinary Reviews: Data Mining and*
195 *Knowledge Discovery* September/ October 2011. 1 (5) p. .
- 196 [Gosain and Kumar (2009)] ‘Analysis of health care data using different data mining techniques’. A Gosain , A
197 Kumar . *IAMA 2009. International Conference on*, 2009. July 2009. 6 p. . (Intelligent Agent & Multi-Agent
198 Systems)
- 199 [Majali et al. ()] ‘Data mining techniques for diagnosis and prognosis of cancer’. J Majali , R Niranjana , V Phatak
200 , O Tadakhe . *International Journal of Advanced Research in Computer and Communication Engineering* 2015.
201 4 (3) p. .
- 202 [Iyer et al. ()] ‘Diagnosis of diabetes using classification mining techniques’. A Iyer , S Jeyalatha , R Sumblay .
203 *IJDKP* 2015. 5 (1) p. .
- 204 [Tan ()] ‘Neighbor-weighted K-nearest neighbor for unbalanced text corpus’. S Tan . [http://orange.biolab.](http://orange.biolab.si/)
205 [si/](http://orange.biolab.si/) *Expert Systems with Applications* 2005. 28 (4) p. .
- 206 [Orange Data Mining Library Documentation Release] *Orange Data Mining Library Documentation Release*, 3.
207 (Orange Data Mining)
- 208 [Prerana et al. ()] ‘Prediction of Heart Disease Using Machine Learning Algorithms-Naïve Bayes, Introduction
209 to PAC Algorithm, Comparison of Algorithms and HDPS’. T H M1 Prerana , N Shivaprakash , C2 , N3
210 Swetha . *International Journal of Science and Engineering* 2347-2200. 2015. 3 (2) p. .
- 211 [Kirkby ()] *WEKA Explorer User Guide for version 3-3-4*, R Kirkby . 2002. University of Weikato