

Secure Data Analytics for Heart Disease Prediction

Rahul S. Belli, Prof. Ajay Nadargi

Department of Computer Engineering Sinhgad Institute of Technology, Maharashtra, India

ABSTRACT: Technological rapid growth in biomedical applications generate high volume of personal data. This biomedical data raises privacy concerns as it reveals sensitive data such as health status and peoples living style. Information generated by biomedical mobile applications need to keep private. The proposed system keeps private data locally on mobile and only data required for heart disease prediction is uploaded to server. Data analytics for Heart disease prediction is implemented using two algorithms Logistic Regression and Nave Bayes. This paper proposes a rule based model to compare the accuracies of applying rules to the individual results of nave bayes and logistic regression on the mHealth application database in order to present an accurate model of predicting heart disease.

KEYWORDS: Data Analytics, Nave Bayes, Logistic Regression, mHealth application

I. INTRODUCTION

Data mining is a process of extracting useful information from large amount of data set. The resulted data is used for prediction purpose so that it can be beneficial for various purposes like improving business process, finding causes of diseases likewise. Different techniques involved in data mining are classification, clustering, association etc. Data mining has immense applicability in diverse area like Biomedical Analysis, Telecom Industry, Intrusion Detection System, Financial Data Analysis etc. In biomedical analysis, different algorithms are used to design model for healthcare prediction. The algorithms used are Nave Bayes, SVM, logistic regression etc.

Around 17.3 billion people die in the world for year of 2008[9]. In spite of the fact that cardiovascular diseases are controllable by taking cautions which were analyzed using different data mining techniques. Last two decades, a lot of research is going on in health care industry to find out the causes of various diseases as precaution is always better than prevention. Cardiovascular diseases includes heart failure, cardiomyopathy, coronary heart disease etc and the common causes for these diseases are diabetes, smoking, high cholesterol, hypertension. Data set used for the data analytics is downloaded from UCI Machine Learning Repository[10].

Motivation behind our proposed system are as follows:

1. Around 17.3 billion people died in the world for year of 2008
2. Precaution is always better than prevention
3. Security to private health data is inevitable

The proposed model consists of two approaches for disease prediction, i.e, logistic regression and naive bayes along with data is secured with encryption. Our objective is to suggest best suitable approach for heart disease prediction with providing security to health data set.

II. REVIEW OF LITERATURE

Now a days multiple people can use forward biomedical sensors and mobile application. That technology generate large amount of biomedical data which include some personal information of patient about lifestyle. So in this paper personal information can be kept secretly and logistic regression technique can be used to predicate disease prediction and treatment[1].

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 7, July 2017

In today's environment large companies or hospital can be store the patients sensitive information on host data centre. An efficient algorithm / techniques are used for designing predictive model for disease diagnosis and treatment. In this paper the information can be keeping locally and use Homomorphic encryption. Homomorphic encryption can be handle enumeration of such encryption in which can not be decryption and not need of decryption key[3].

The Support Vector Machine (SVM) is a state-of-the-art classification method introduced in 1992. The SVM classifier is widely used in bioinformatics due to its high accuracy, ability to deal with high-dimensional data such as gene expression, and exhibity in modeling diverse sources of data. In this paper, PPSVC techniques can be used to tackle the privacy violation problem of the classification model of the SVM[5].

MediNet is a system that can be used to developed to personalize the self healthcare process for patients with diabetes and cardiovascular disease using a mobile phone network. It can be use current and past information from monitoring devices to recommendations. It can provide for the uniqueness of each patient by personalizing its recommendations based on personal level characteristics of the patient, as well as groups of patients share that characteristics[2].

Cloud-assisted mobile health (CAM) can be monitoring mobile communications and cloud computing technologies to provide feedback decision support, has been considered as a revolutionary approach to improving the quality of healthcare service while lowering the healthcare cost. CAM, which can effectively provide security for privacy of clients and the intellectual property Of mHealth service providers.[4]

The healthcare industry collects huge amounts of healthcare data which, unfortunately, are not “mined” to discover hidden information for effective decision making. Discovery of hidden patterns and relationships often goes unexploited. Advanced data mining techniques can help remedy this situation. This paper describes about a prototype using data mining techniques, namely Naïve Bayes and WAC (weighted associative classifier). This system can answer complex “what if” queries which traditional decision support systems cannot. Using medical profile 0073 such as age, sex, blood pressure and blood sugar it can predict the likelihood of patients getting a heart disease. It enables significant knowledge, e.g. patterns, relationships between medical factors related to heart disease, to be established. It can serve a training tool to train nurses and medical students to diagnose patients with heart disease. It is a web based user friendly system and can be used in hospitals if they have a data ware house for their hospital. Presently we are analyzing the performances of the two classification data mining techniques by using various performance measures.[11]

III. SYSTEM ARCHITECTURE

For heart disease prediction, this paper proposes a combination of models which is shown in Fig 1. This secure data analytics approach is divided into five modules involving mHealth application, Preprocessing, Training Model, Applying Logistic Regression, applying Naive Bayes, and Result Comparison and Heart Disease prediction. Patient database is collected from mHealth application and also taken from UCI repository for training purpose.

Proposed System:

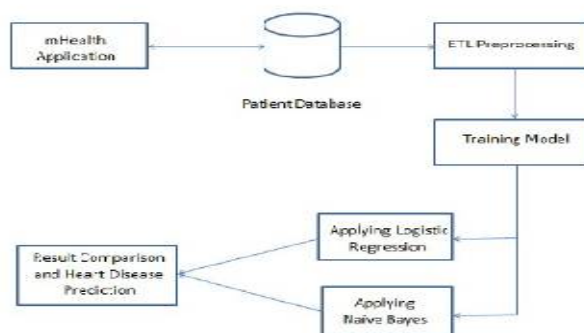


Fig 1: Proposed System

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 7, July 2017

The above mentioned five models are described as below.

A. mHealth Application

Mobile health technologies are rapidly growing which includes wearable devices and embedded sensors. As technology growing rapidly, mHealth application development using mobile health technology is also growing in a same way. mHealth applications record all day to day activity of individual and inform if any cautions need to take. But it comes with privacy concern.

B. ETL Preprocessing

ETL stands for Extract, Transform, Load. It is a tool which is combination of three functions which is mentioned above. It is used to get data from one database and transfer to other database. Data preprocessing is a data mining technique used to transform sample raw data into an understandable format. Real world collected data may be inconsistent, incomplete or contains an error. This paper proposes ETL and preprocessing combined together to process on existing patient data and load it into mHealth server database.

Following are data processing techniques.

1. Data Cleaning: Resolve inconsistency and eliminate noise in data.
2. Data Integration: Incorporate data from different sources into one rational source such as data warehouse.
3. Data Transformation: In data transformation, data transform from one source to another source. It involves following terms.
 1. Normalization
 2. Aggregation
 3. Generalization

C. Training the Model

Each of the two models has been trained using different methods. For logistic regression, the first step to training is to find the significant attributes by calculating their individual P-values. As a rule of thumb, if it is below 0.05, only then is the attribute significant. The Hosmer-Lemeshow test is also required to check for goodness fit of the model. The corresponding P-value must abide by a 5% level. Naive Bayes is a classification method which works on Bayes theorem which assumes independence among attributes. The basic assumption is that presence of a particular feature in a class unrelated to presence of any other. Naive Bayes requires less amount of data as compared to other method for estimation of parameters for classification.

D. Applying Logistic Regression and Naive Bayes

For prediction of categorical dependent variable outcome, from set of independent variables [8]. Logistic regression is mainly used for prediction and also calculating the probability of success. Logistic Regression involves fitting an equation of the form to the data [8].

$$y = e^{\Lambda(b_0 + b_1 * x)} / (1 + e^{\Lambda(b_0 + b_1 * x)}) \dots\dots\dots(1)$$

x = input values

y = output values

b₀ = bias or intercept term

b₁ = coefficient for single input value (x)

The Naive Bayes classifier is based on Bayes theorem with independent assumptions between predictors. Continuous values associated with each class are distributed according to a Normal distribution [6]. Naive Bayes classification algorithm is based on Bayes theorem.

$$P(A_k/B) = P(A_k \cap B) / (P(A_1 \cap B) + P(A_2 \cap B) + \dots + P(A_n \cap B)) \dots\dots\dots(2)$$

A_k = set of mutually exclusive events

B = any event from the sample space such that P(B) > 0

Here, In proposed system, Bayes theorem can be written in below given way.

$$P(D/S) = P(D/S) * P(D)/P(S) \dots\dots\dots(3)$$

D = Disease

S = Symptom

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 7, July 2017

IV.METHODOLOGY

A. Logistic Regression

For prediction of categorical dependent variable outcome, from set of independent variables[8]. Logistic regression is mainly used to for prediction and also calculating the probability of success. Logistic Regression involves fitting an equation of the form to the data[8].

$$y = \frac{e^{(b_0 + b_1 x)}}{1 + e^{(b_0 + b_1 x)}}$$

x=input values

y=output values

b0 = bias or intercept term

b1=coefficient for single input value (x)

B. Naïve Bayes

The Naive bayes classifier is based on Bayes theorem with independent assumptions between predictors Continuous values associated with each class are distributed according to a Normal distribution [6]. Naive Bayes classification algorithm is based on Bayes theorem.

$$P(A_k | B) = \frac{P(A_k \cap B)}{P(A_1 \cap B) + P(A_2 \cap B) + \dots + P(A_n \cap B)}$$

A_k = set of mutually exclusive events

B = any event from the sample space such that P(B)>0

C. Naïve Bayes

Here, In proposed system, Bayes theorem can be written in below given way.

$$P(D/S) = \frac{P(D/S) P(D)}{P(S)}$$

D = Disease

S= Symptom

D. PSR Algorithm:

1. Positional method
2. Substitutional method
3. Random pattern

1. Positional method

This method is used for generating a key not for encryption. The generated key will be in next method for encryption.

Given: - File to be encrypted, a key from user (0 ≤ key ≤ 9), set of symbols

We'll use all above given + File length to generate a new key (we'll append given key with a delimiter).

E.g.: - Generated key => xyz###key

2. Substitutional method

Here we'll use 1st user key to divide the content of file to be encrypted into several blocks and then for each block, we'll use generated key to split the content of each block into array.

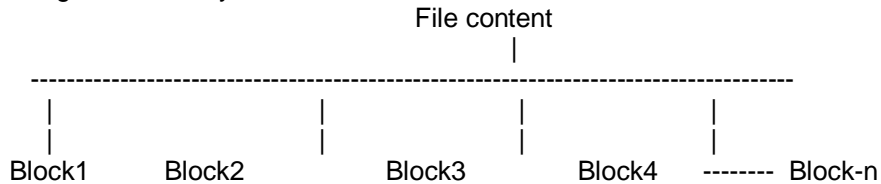
International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 7, July 2017

E.g.: - If user key = n, then we'll divide the File content into n-blocks.



For each block, we'll use below method to encrypt its content.

E.g.: - If Generated key=>xyz###key, then we'll use its length to split the block content into array.

	x	y	z	#	#	key
C1	—	—	—	—	—	—
C2	—	—	—	—	—	—
C3	—	—	—	—	—	—
·	—	—	—	—	—	—
·	—	—	—	—	—	—
Cn	—	—	—	—	—	—

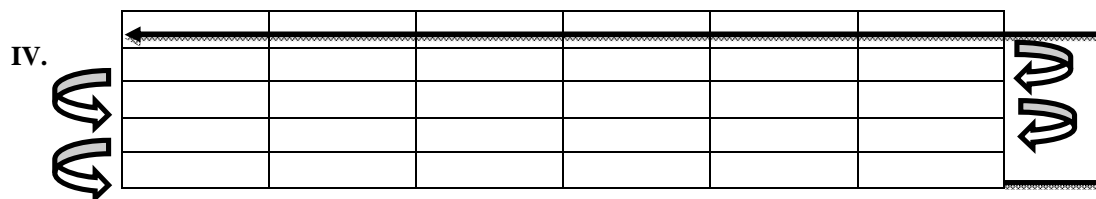
Now, we'll encrypt content of each block with respect to generated key character. With this, we'll take diagonal values and concatenate them.

Let's say, the O/P of this method for block-1 is—**abcdefg**, then we'll be having n outputs.

B1 -	—	—	—	—	—
B2 -	—	—	—	—	—
B3 -	—	—	—	—	—
B4 -	—	—	—	—	—
B5 -	—	—	—	—	—
Bn -	—	—	—	—	—

3. Random Pattern

Here, we'll shift content of blocks from one block to its next block by the value of user key.



V.RESULT COMPARISON AND HEART DISEASEPREDICTION

This modules includes the comparison of results from LogicalRegression and Naive Bayes algorithm. Also this systemfind out the accuracy percentage from both the algorithms.

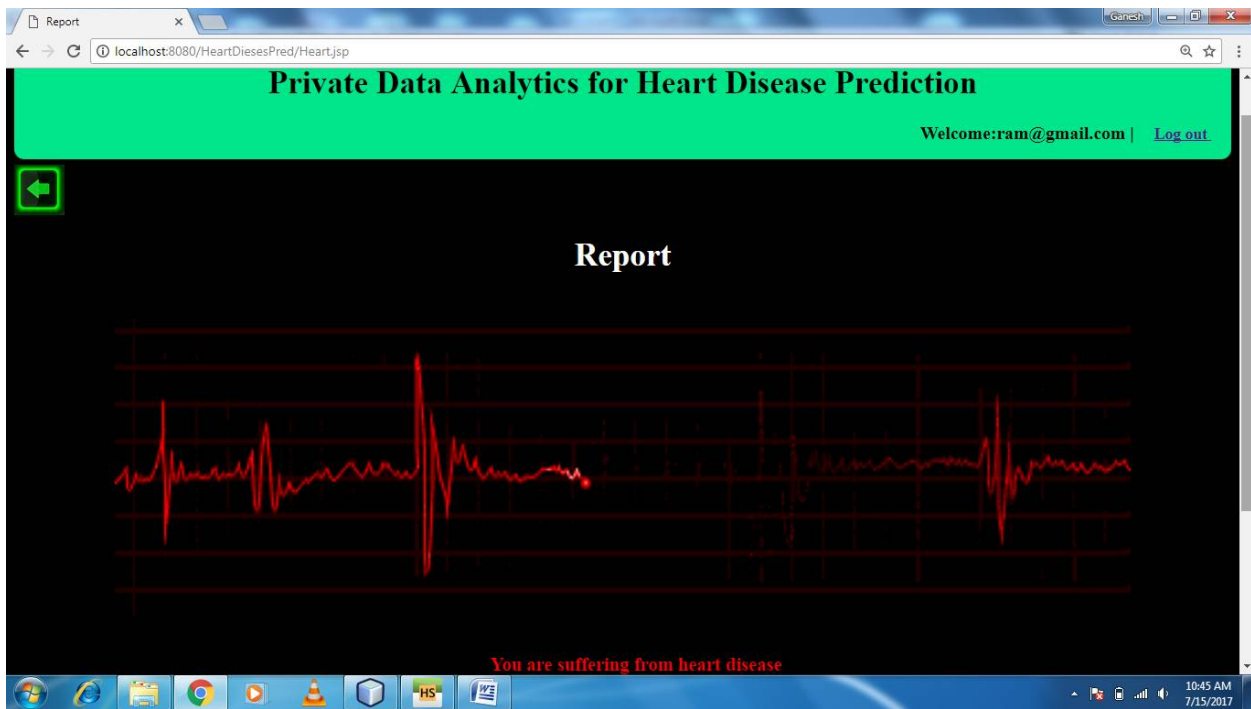


Fig2: Report For heart dieses



Fig3:Result Page for not heart dieses

International Journal of Innovative Research in Science, Engineering and Technology

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirset.com

Vol. 6, Issue 7, July 2017

VI. CONCLUSION

This system proposes confidential scheme for predicting heart disease using two different models, Naive Bayes and Logistic Regression. As identified through survey, it is a need to have combinational approach to increase the accuracy of prediction for heart disease. In our research attempt, the future work may also involve the development of a tool to predict the risk of disease of a prospective patient. Thus, our future work would be also incorporate other data mining techniques, such as Clustering and Association Rules etc.

REFERENCES

- [1] YanminGong ,Yuguang Fang , YuanxiongGuo ,Private Data Analytics on Biomedical Sensing Data Via Distributed Computation 1545-5963 (c)2015 IEEE.
- [2] P. Mohan, D. Marin, S. Sultan, and A. Deen, Medinet: personalizing the self-care process for patients with diabetes and cardiovascular disease using mobile telephony, in Engineering in Medicine and Biology Society, 2008. EMBS 2008. 30th Annual International Conference of the IEEE. IEEE, 2008, pp. 755-758.
- [3] J. W. Bos, K. Lauter, and M. Naehrig, Private predictive analysis on encrypted medical data, Journal of biomedical informatics, vol. 50, pp. 234-243, 2014.
- [4] H. Lin, J. Shao, C. Zhang, and Y. Fang, Cam: cloud-assisted privacy preserving mobile health monitoring, Information Forensics and Security, IEEE Transactions on, vol. 8, no. 6, pp. 985-997, 2013.
- [5] K.-P. Lin and M.-S. Chen, On the design and analysis of the privacy preserving svm classifier, Knowledge and Data Engineering, IEEE Transaction on, vol. 23, no. 11, pp. 1704-1717, 2011.
- [6] R. Agrawal and R. Srikant, Privacy-preserving data mining, in ACM Sigmod Record, vol. 29, no. 2. ACM, 2000, pp. 439-450.
- [7] Shalabi, L.A., Z. Shaaban and B. Kasasbeh, Data Mining: A Preprocessing Engine, J. Comput. Sci., 2: 735-739, 2006.
- [8] Mythili T., DevMukherji, Nikita Padalia, and Abhiram Naidu, A Heart Disease Prediction Model using SVM-Decision Trees Logistic Regression (SDL), International Journal of Computer Applications (0975 8887), Volume 68 No.16, April 2013.
- [9] <http://www.world-heart-federation.org/cardiovascular-health/global-factsmap/>
- [10] Robert Detrano 1989 Cleveland Heart Disease Database V.A. Medical Center, Long Beach and Cleveland Clinic Foundation.
- [11] N. adityasundar, P. pushpalatha, M. ramachandra "Performance analysis of classification data mining techniques over heart disease data base" [IJESAT] international journal of engineering science & advanced technology issn: 2250-3676 volume-2, issue-3, 470 – 478