

Predicting COPD, Diabetes and Heart Disease Using Machine Learning Algorithm

Kandigati John Prajwal¹, Likitha U K², Kosur Shriyaa Saivijya³, Kavalakuntla Manikanta Reddy⁴, Sushmita Kumari⁵

Student, Computer Science Department, Presidency University, Bangalore, India¹⁻⁴

Assistant Professor, Computer Science Department, Presidency University, Bangalore, India⁵

Abstract: This project mainly contains python, machine learning and web related concepts. The main aim of this project is to predict Chronic obstructive pulmonary disease, Diabetes and Heart disease using machine learning techniques and algorithms. In this project the Doctors/Patients/User can register in the webpage which we developed and again then can login. All their information will store it in the database. Once they login into the webpage they can choose which disease to predict like they will get options COPD, Diabetes and Heart. The selected disease prediction page will open and they can enter the patient details, mail id and the parameters asked in the webpage required to predict that particular disease. Once entering all the details, if they click on predict means they will get final results in the webpage like they have that particular disease or not, stage of the disease and treatment should be taken for that disease will be displayed on the final results. The same predicted results also send to the patient email whatever entered by the user/doctor/patient while entering their details. You can also see the data stored in the database in this project.

Keywords: Department, Presidency Tree, k-nn algorithm.

I. INTRODUCTION

The medical records are the records of the patient's disease occurrence, development, examination, diagnosis, and treatment. Medical records are an important part of medical care, teaching, prevention, research and development. They are the basic support conditions for hospital management. It is difficult to keep the traditional paper medical records. It is also difficult to query the paper medical records. With the development of information technology, it is possible to realize the medical record management through information systems. The electronic medical record system is very important for the digital management of the medical records. It has greatly improved the hospital's management efficiency. The electronic medical records data can be used for medical big data analysis. There are many uncertainties in developing machine-learning (ML)-based products. Due to the gap between research and development, the overall progress becomes slow, and experiences many failures and learning's only to see an initial idea not working or generating no significant revenue. To minimize these drastic effects, there are continuing studies to make a balanced handshake between research and development to shorten the time span of ML-based products from idea generation to deployment. The electronic medical records are different from the paper medical records. They are the digital patient medical records. The electronic medical records can be stored, managed, transmitted, and reproduced by IT devices. They replace the paper medical records. All information on the paper medical records is stored in the electronic medical records. Doctors, patients or other authorized persons can obtain electronic medical records in a complete, accurate and timely manner. Many years ago, some hospitals in Western developed country began to construct the hospital information system (HIS). Electronic medical records (EMR) have then been researched and applied in these countries. The government was vigorously promoting and popularizing the application of EMR. They researched how to transmit the emergency patient's EMR through the Internet. They researched how to record the patient's entire medical procedure through electronic medical records, including medical examinations, test results, X-ray films, CT films, and doctor's advice. In recent years, many Indian hospitals have successively established information systems. It provides technical support for the research and application of electronic medical records in India. An electronic medical record system is proposed in this paper. It implements the storage and query of medical records through information technology. The electronic medical record system can completely record the occurrence, development, examination, diagnosis, treatment and other medical activities of the patient's disease. It can summarize, sort, and analyse the collected data. It can generate patient medical health records in a prescribed format and requirements. It has the advantages of high transmission speed, good sharing, large storage capacity, easy use and low cost. There are many uncertainties in developing machine-learning (ML)-based products. Due to the gap between research and development, the overall progress becomes slow, and experiences many failures and learning's only to see an initial idea not working or generating no significant revenue. To minimize these drastic effects, there are continuing studies to make a balanced handshake between research and development to shorten the time span of ML-based products from idea generation to deployment. This paper demonstrates a three-phase ML product development workflow at One Class. The first phase of the workflow considers the pivotal idea generation for products that involves data reliability assessment, idea prioritization, expectation setting, and building trust among users. The second phase concentrates on several state-of-

the-art strategies for planning and future re-use of several product components. Finally, the actual research and development phase describes the fail-fast method practiced by One Class to learn quickly from failures and act accordingly. The workflow is followed by the company to develop many sophisticated ML-based products successfully within a very short period of time.

The problem statement is to predict the Risk of Diabetic, Heart Disease and COPD Disease from input symptoms, daily habits and lab reports given by the user in order to help patients early detection of disease to help in treatment and to save the patient from higher consequences.

I.LITERATURE SURVEY

A. Yahyaoui, A. Jamil, J. Rasheed and M. Yesiltepe,

With the continuing increase in the number of the deadly diseases that threaten both human health and life, medical Decision Support Systems (DSS) continue to prove their effectiveness in providing physicians and other healthcare professionals with support in clinical decision making. Among these dangerous diseases, diabetes continues to be one of the leading one that has caused several deaths in the world. It is characterized by an increase in blood sugar levels which can have severe effects on other human organs. According to the International Diabetes Federation (IDA), 382 million people are living with diabetes and by 2035; these statistics will double to reach 592 million. In this paper, we propose a DSS for diabetes prediction based on Machine Learning (ML) techniques. We compared conventional machine learning with deep learning approaches. For conventional machine learning method, we considered the most commonly used classifiers: Support Vector Machine (SVM) and the Random Forest (RF). On the other hand, for Deep Learning (DL) we employed a fully Convolutional Neural Network (CNN) to predict and detect the diabetes patients. The proposed system is evaluated on publicly available Pima Indians Diabetes database which consisted of total 768 samples each with 8 features. 500 samples were labelled as non-diabetic while 268 were diabetic patients. The overall accuracy obtained using DL, SVM and RF was 76.81%, 65.38% and 83.67% respectively. The experimental results show that RF was more effective for diabetes prediction compared to deep learning and SVM methods.

B. Keshav Srivastava, Dilip Kumar Choubey

Weighing only 300 grams, Heart is declining the mortality rate at a rapid pace from decades. The major factors that contribute to it are smoking, drinking, unbalanced diet, and many more. Even with these more technical advancements the analysis of the clinical data is a critical challenge. With the use of Machine Learning techniques, it is possible to analyse the data and interpret the cause that led to heart diseases such as Coronary Heart Disease, Arrhythmia, and Dilated Cardiomyopathy. Many researchers are developing IoT enabled hardware to predict these diseases using various ML and DM techniques. In this study, we propose a novel method to determine the disease using Cleveland Heart Disease Dataset by combining the computational power of various ML and DM algorithms and concluded that among all the algorithms, K-Nearest Neighbours gives the highest accuracy of 87%. Along with this, a web app is developed using flask in python with which the user can enter the attributes and predict the heart disease.

Cristóbal Esteban, Javier Moraza, Cristóbal Esteban, Fernando Sancho, Myriam Aburto, Amaia Aramburu, Begoña Goiria, Amaia Garcia Loizaga, Alberto Capelastegi

eEPOC database is composed by daily reports sent by the patients with the following information: heart rate, temperature, oxygen saturation, respiratory rate, steps walked and a questionnaire form about symptoms. According to this, an alarm system composed by three levels of exacerbation (green, yellow and red) is established.

On this data the Random Forests Algorithm was applied to predict when a patient will present a red alarm. We used a 10-fold cross validation to estimate the performance of the model. The development was implemented using the packages Scikit-Learn and Pandas from the programming language Python

Survey of Machine Learning Algorithms for Disease Diagnostic

Meherwar Fatima Maruf Pasha

In medical imaging, Computer Aided Diagnosis (CAD) is a rapidly growing dynamic area of research. In recent years, significant attempts are made for the enhancement of computer aided diagnosis applications because errors in medical diagnostic systems can result in seriously misleading medical treatments. Machine learning is important in Computer Aided Diagnosis. After using an easy equation, objects such as organs may not be indicated accurately. So, pattern recognition fundamentally involves learning from examples. In the field of bio-medical, pattern recognition and machine learning promise the improved accuracy of perception and diagnosis of disease. They also promote the objectivity of decision-making process. For the analysis of high-dimensional and multimodal bio-medical data, machine learning offers a worthy approach for making classy and automatic algorithms. This survey paper provides the comparative analysis of different machine learning algorithms for diagnosis of different diseases such as heart disease, diabetes disease, liver disease, dengue disease and hepatitis disease. It brings attention towards the suite of machine learning algorithms and tools that are used for the analysis of diseases and decision-making process accordingly.

Heart Disease Diagnosis using Machine Learning Algorithm

Recent advances in computing and developments in technology have facilitated the routine collection and storage of medical data that can be used to support medical decisions. However, in most countries, there is a first need for collecting and organizing patient's data in digitized form. Then, the collected data are to be analysed in order for a medical decision to be drawn, whether this involves diagnosis, prediction, course of treatment, or signal and image analysis. In this paper, India centric dataset is used for Heart disease diagnosis. The correct diagnosis performance of the automatic diagnosis system is estimated by using classification accuracy, sensitivity and specificity analysis. The study shows that, the SVM with Sequential Minimization Optimization learning algorithm have better choice for medical disease diagnosis application.

A study on Deep Machine Learning Algorithms for diagnosis of diseases

Dinu A.J., Ganesan R1, Felix Joseph and Balaji

Consider that we are living in a place that is far away from a hospital or do not have sufficient money to cover up the hospital bill or do not have enough time to take off work. In such cases, the disease diagnosis through sophisticated machines would be lifesaving. Scientists had developed numerous artificially intelligent diagnosis algorithms for detecting various diseases like Rheumatoid Arthritis, Cancer, Lung Diseases, Heart Diseases, Diabetic Retinopathy, Hepatitis Disease, Alzheimer's disease, Liver Disease, Dengue Disease and Parkinson Disease. Deep learning uses large artificial neural networks layers having interconnected nodes which can rearrange themselves as and when new information comes in. This technique allows the computers to self-learn on their own without the need for human programming. This paper focuses on recent developments in machine learning which have made significant impacts in the detection and diagnosis of various diseases.

● **REQUIREMENTS**

Hardware:

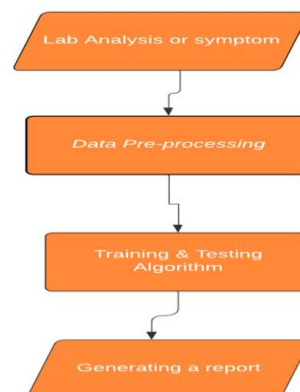
- Laptop
- 4GB RAM(MIN)
- 100MB Free Space in the Disk
- Windows 10 or 10+
- Intel-7 2.4GHZ
- Any Desktop with above Configuration or higher level

Software:

- Python(language)
- PyCharm
- Jupyter Notebook
- Visual studio code
- Dataset of COPD,
- Dataset of Diabetes
- Dataset of Heart diseases.
- Libraries- Sklearn, NumPy, pandas, Email, smtplib

● **EXISTING SYSTEM**

There are many health-related systems where the system predicts the disease from the lab reports and the reports of the doctors. Also, most of the tools are built in doctor usage perspective where a common man is not able to get the full advantage of those systems.



Problem Motivation:

The initial diagnosis of disease results in 93% of cure in any deadly diseases and the disease risk can be predicted using the small symptoms and daily habits of a person. And the most important and common diseases that we are finding in most of the people around the world are Diabetes, Heart Diseases, COPD. And the risk of these diseases can be predicted using simple symptoms we can observe, daily habits and small amounts of lab related data. If we can build a system which can predict the risk of disease then we can limit the higher consequences of the disease.

Project Aims:

- To develop a machine learning prediction algorithm to predict specific diseases based on the user's daily habits symptoms and lab reports.
- To generate a report consisting of the risk of disease name based on the data collected from the user.
- To Send the report via mail from the System automatically to the mail id given by the user.

Drawbacks of Existing System

The Existing System contains the separate mechanisms for three disease predictions and also not achieving good accuracy. In medical domain the most important factor to be driven is Accuracy and if we see the results of different algorithms that in the existing systems used:

- Heart Disease Prediction is predicted using the algorithm KNN with Accuracy of 85%
- Diabetic Disease Prediction is done using the algorithm Random Forest with Accuracy 83.7%
- COPD prediction is done using the algorithm Decision Tree with very less Accuracy 72%

Above numbers show that the algorithms are not efficient enough to be trustworthy or to be used to predict disease perfectly. This shows the unreliability of the existing systems and also since they are separate for each other it's not easy for user also to use it.

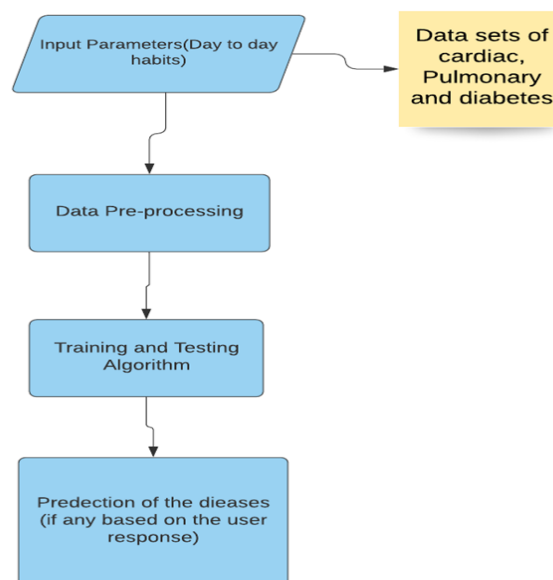
PROPOSED SYSTEM

Fig5: Proposed System

A huge number of medical records are stored in the electronic medical record system. Since the electronic medical record system allows a user to access it through the Internet, a large number of users will search through the Internet. The system needs to be able to support multi-user online search and invocation. The users need to find medical records data quickly and accurately in a huge number of medical records. The server's computing power and database's response speed will directly affect the speed of user data query. It requires the server of the electronic medical record system to have extremely high computing power and database response speed.

The last year the health sector is the most prioritized sector throughout the entire planet. The spread of the virus has created a panic through people. The main aim of the project is to try and detect the problems earlier before the situation are out of hand and take the required precautions according to the solution provided.

The software will take the required assessment related to the disease's (COPD, diabetes and heart) and give the assessment based on the information given by the patient. The details of the patient and the result will be email to the patient as provided in the details.

Datasets has been tested on different algorithms to get the best and accurate model which is lacking in the existing system. So, we tried to maximize the accuracy as much as possible.

The results of the accuracy of the algorithms we've used are represented here:

- Heart Disease dataset Decision Tree algorithm found to be best with accuracy of 90.16%
- We used Gini index as criterion to choose the best split feature with lower Gini index value is chosen as the best split

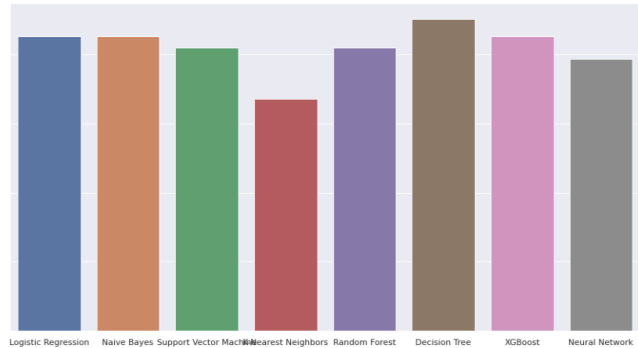


Fig6: Accuracies of different algorithms for Heart Disease Dataset

- The Diabetic dataset we found Decision tree with the best accuracy of 93.97% which is 10% more than the existing algorithm. Here also Gini index is used as the criterion for selecting the split.

Below are the graphical representations of the accuracies of different algorithms trained by diabetic dataset:

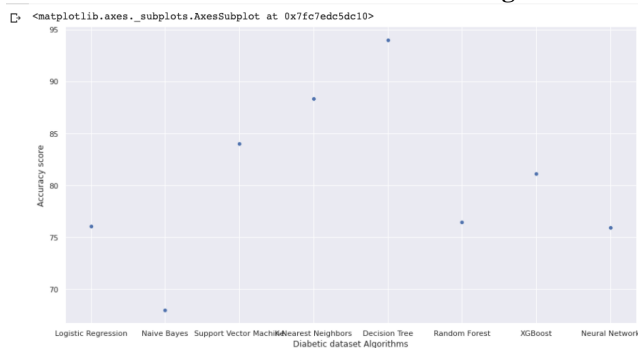


Fig7: Accuracies of different algorithms for diabetic dataset

- COPD dataset the best algorithm that we've found is K-NN Algorithm with accuracy of 94.89% which is almost 15% more than existing algorithms.
- The value of K is 7 neighbours in the model and the metric used is minkowski the value of p is 2 which indicates Euclidean measure is used.

The Graphical representation of different algorithm models on the dataset is given below:

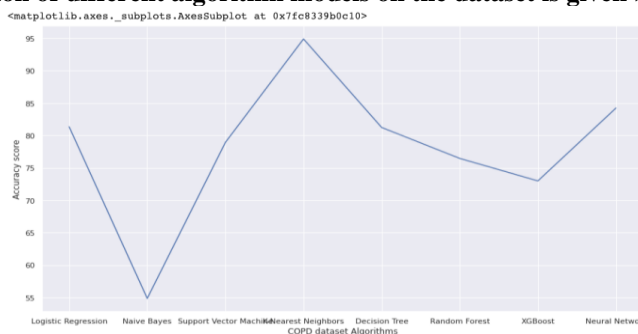


Fig8: Accuracies of different algorithms for COPD Dataset

Advantages of proposed system:

- The accuracy of the all the models is higher to 90% which is very high compared to existing systems and makes it more reliable
- Prediction for 3 diseases is done in single platform which is easier to access

- Good User interface with easily understandable columns which makes the tool more user friendly
- Sending of Email to the given mail address makes the system unique where user don't need to memorize the report he can see whenever he can in email

● IMPLEMENTATION

1. COPD, Heart Disease and Diabetes dataset is taken and loaded.
2. The dataset divided as train data and test data.
3. The data is pre-processed in order to fill the missing values in the dataset, standardization and label encoding etc. to increase the accuracy.
4. The features which are directly affecting the final results were extracted.
5. The model is built using machine learning algorithms like KNN, Decision Tree.
6. The model is trained with the pre-processed data.
7. The model is tested and accuracy is calculated for different ML algorithms.
8. The algorithm with best accuracy is finalized.
9. The finalized model can be used to detect the COPD, Heart Disease and Diabetes with good accuracy.
10. We will create a front-end using Django and Html.
11. The user can register and login in the front end (webpage)
12. After logged in user can choose which disease to predict and after that the relevant parameters required for disease prediction will be collected from the user.
13. Once user entered all the required details and click on the predict button and he will get the disease prediction report and the same thing will be sent to the registered email of the user.

```
text_representation = tree.export_text(dt)
print(text_representation)

--- feature_2 <= 0.50
--- feature_11 <= 0.50
--- feature_12 <= 2.50
--- feature_7 <= 119.50
|--- class: 0
|--- feature_7 > 119.50
|--- feature_12 <= 1.50
|--- feature_0 <= 54.50
|--- class: 0
|--- feature_0 > 54.50
|--- class: 1
|--- feature_12 > 1.50
|--- feature_0 <= 60.00
|--- class: 1
|--- feature_0 > 60.00
|--- feature_0 <= 63.50
|--- feature_7 <= 155.50
|--- class: 0
|--- feature_7 > 155.50
|--- feature_0 <= 61.50
|--- class: 0
|--- feature_0 > 61.50
|--- class: 1
```

Fig10: Textual Representation of Decision tree for heart disease

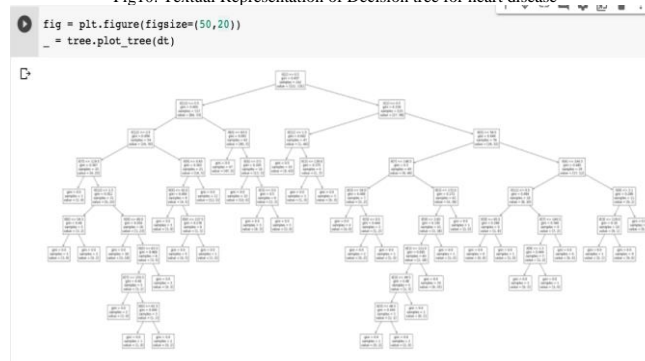


Fig11: Decision Tree for Heart Disease Prediction


```
▼ K Nearest Neighbors
+ Code + Text
[45] from sklearn.neighbors import KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=7)
knn.fit(X_train,Y_train)
Y_pred_knn=knn.predict(X_test)
print(knn)

KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=None, n_neighbors=7, p=2,
weights='uniform')
```

Fig12: Summary of KNN Algorithm used for COPD Prediction

```
for x in range(200):
dt = DecisionTreeClassifier(random_state=x)
dt.fit(X_train,Y_train)
Y_pred_dt = dt.predict(X_test)
current_accuracy = round(accuracy_score(Y_pred_dt,Y_test)*100,2)
if(current_accuracy>max_accuracy):
max_accuracy = current_accuracy
best_x = x

#print(max_accuracy)
#print(best_x)

dt = DecisionTreeClassifier(random_state=best_x)
dt.fit(X_train,Y_train)
Y_pred_dt = dt.predict(X_test)
print(dt)

DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
max_depth=None, max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort='deprecated',
random_state=11, splitter='best')
```

Fig13: Summary of Decision tree used for diabetes

● TESTING

The project was created with an aim to help people try and predict the disease earlier so that they can take the necessary precaution required. We have created the code using PyCharm where the link of the webpage is generated when code is executed.

1. We have created a registration page where each user can create their own personal account by giving the basic and minimal details required.

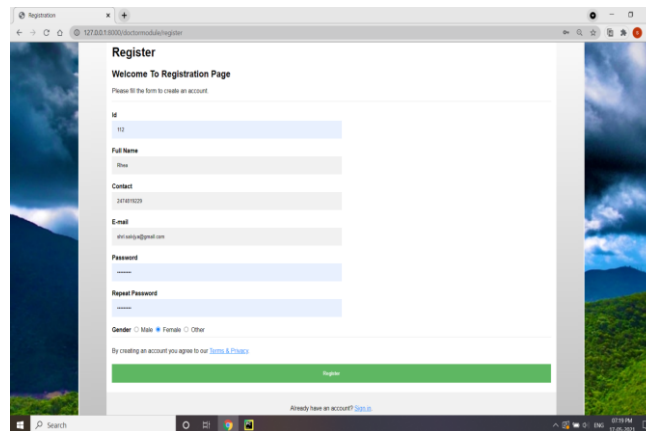


Fig14: Registration Page

2. The login in page of a demo user is show below:

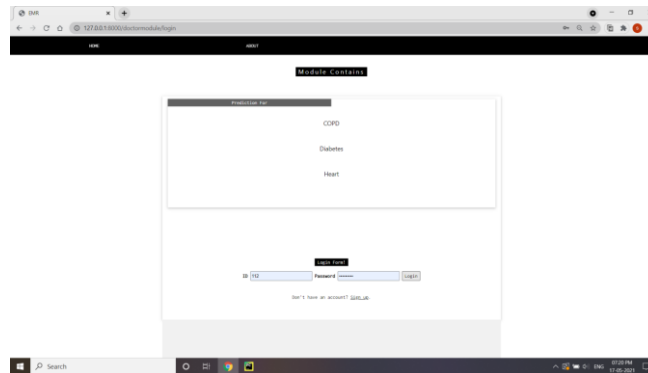


Fig15: Login Page

3. After the login the user can choose the diseases, they want to get a prediction about from the home page as shown below:



Fig16: Home Page

Sample prediction of Diabetes

4. Following the home page with the required detailed for diabetes is shown below:

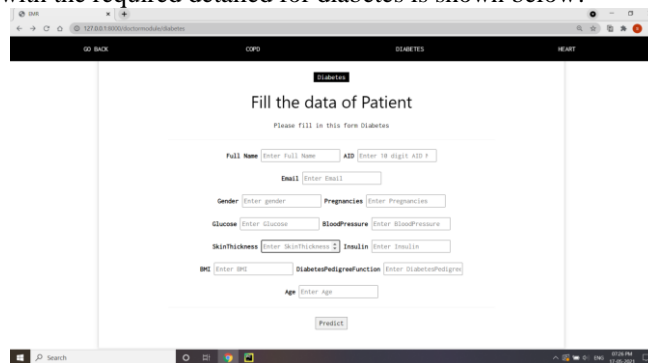


Fig17: Diabetes form

5. After all the boxes are filled in as shown below:

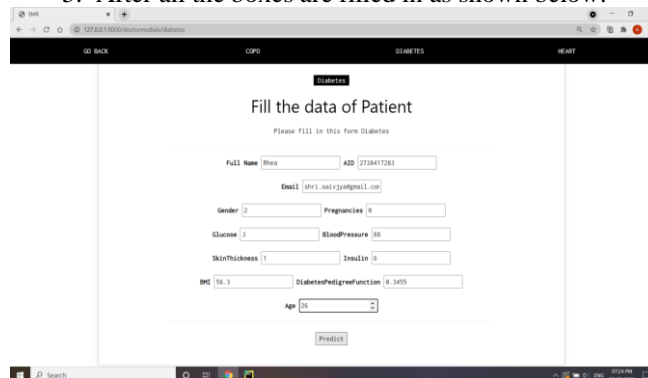


Fig18: Sample filled Diabetes form

6. The code predicts the whether the user is at the risk of the disease or no as shown below:

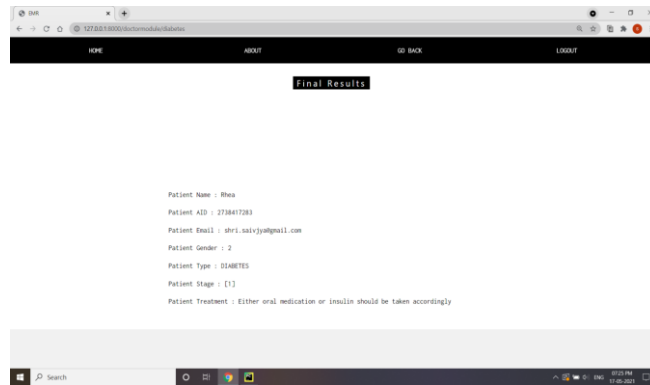


Fig19: Sample Result of Diabetes

7. The predicted results is directly mailed to user for further reference.

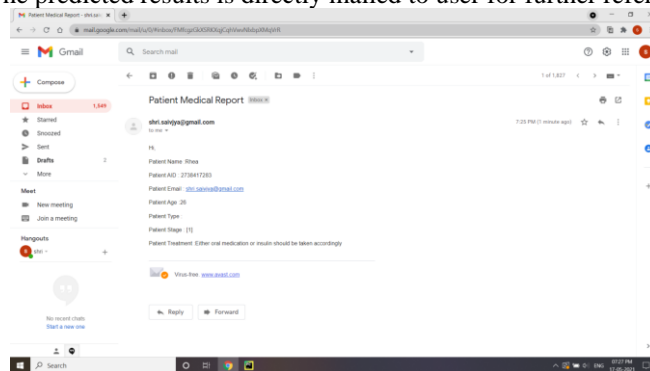


Fig20: Sample Diabetes E-mail report

Sample prediction of COPD (Chronic obstructive pulmonary disease)

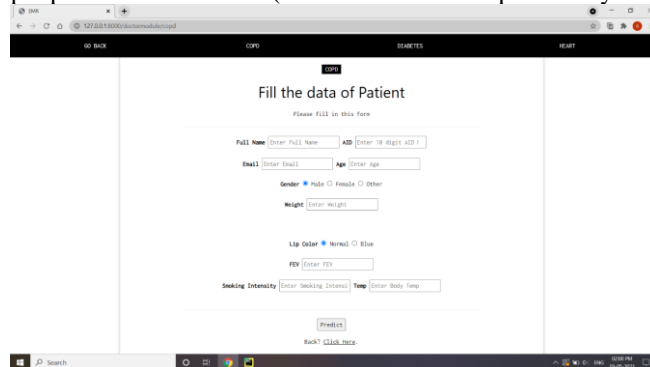


Fig21: COPD Form

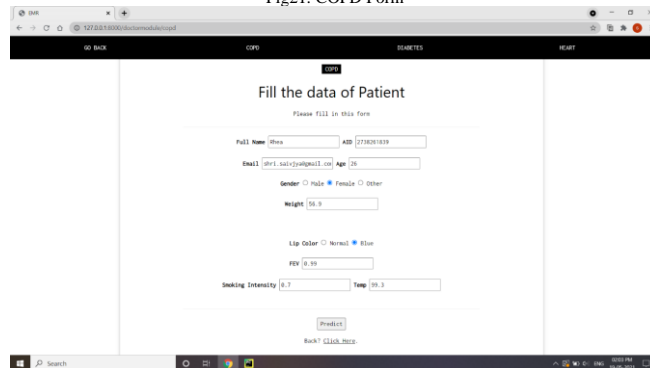


Fig22: Sample COPD Filled Form

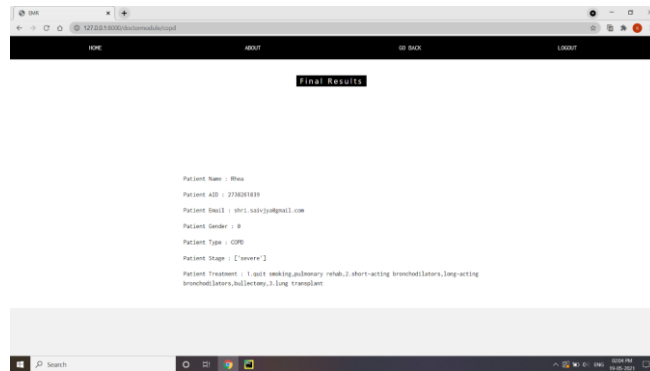


Fig23: Sample Result of COPD

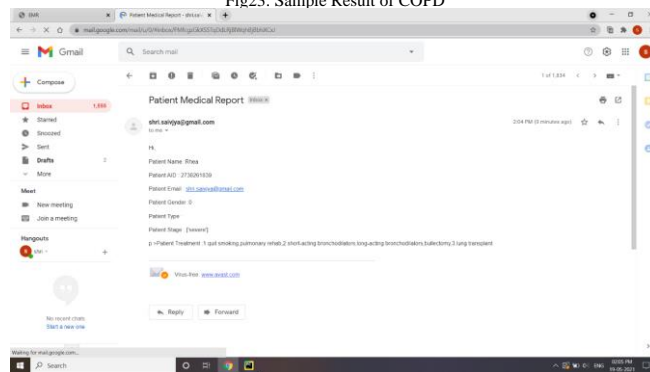


Fig24: Sample COPD E-mail Report
Sample prediction of heart disease

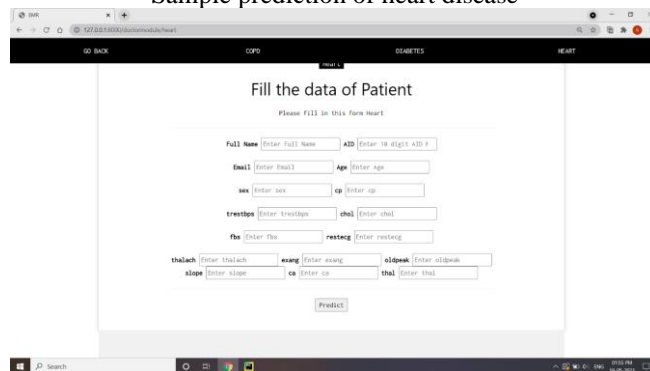


Fig25: Heart Disease Form

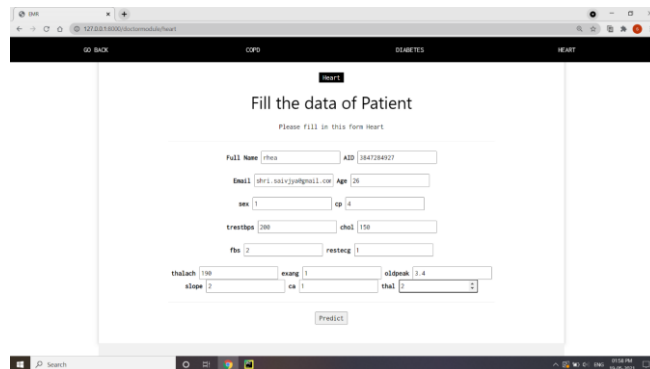


Fig26: Sample Heart Disease Filled Form

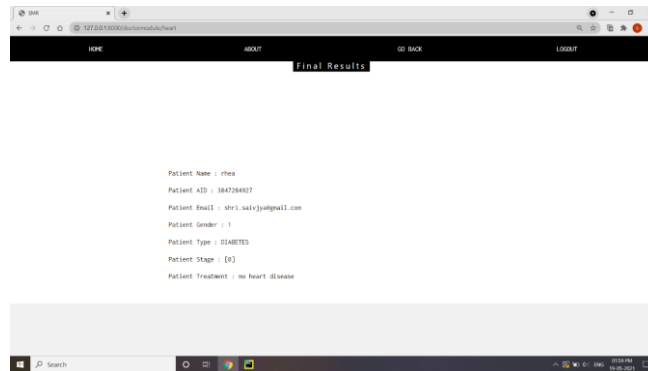


Fig27: Sample Result of Heart Disease

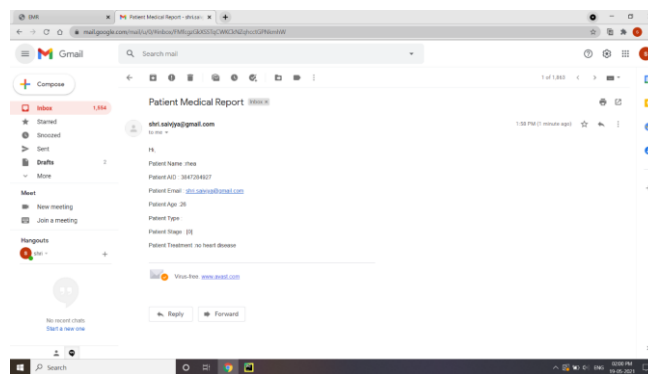


Fig28: Sample Heart Disease E-mail Report

TABLE I. TABLE TYPE STYLES

Table Head	Table Column Head		
	Table column subhead	Subhead	Subhead
copy	More table copy ^a		

^a. Sample of a Table footnote. (Table footnote)

Fig. 1. Example of a figure caption. (Figure caption)

II. CONCLUSION

The selected disease prediction page will open and they can enter the patient details, mail id and the parameters asked in the webpage required to predict that particular disease. Once entering all the details if they click on predict means they will get final results in the webpage like they have that particular disease or not, stage of the disease and treatment should be taken for that disease will be displayed on the final results. The same predicted results also send to the patient email whatever entered by the user/doctor/patient while entering their details. You can also see the data stored in the database in this project. The electronic medical record system realizes the information management of medical records. Doctors can remotely predict the patient medical records via the Internet. In emergency situations, doctors can quickly query the data of the electronic medical records. Paper medical records are usually kept only in the hospital. It is difficult to share paper medical records with other hospitals. The electronic medical record system can be authorized to allow other hospital doctors to view the medical records of the hospital. It achieves the sharing medical records online. It provides technical support for the development of telemedicine. With the development of computer storage technology, the storage capacity of the electronic medical record system database will be massive. Doctors can easily store, retrieve, and browse the medical records through the electronic medical record system. They can also conduct research and statistical analysis through the electronic medical records. This greatly reduces the workload of manual data collection. It has greatly improved the level of medical research. At the same time, the electronic medical record system does not require paper, which saves costs.

Figures and Tables

Positioning Figures and Tables: Place figures and tables at the top and bottom of columns. Avoid placing them in the middle of columns. Large figures and tables may span across both columns. Figure captions should be below the figures; table heads should appear above the tables. Insert figures and tables after they are cited in the text. Use the abbreviation “Fig. 1”, even at the beginning of a sentence.

**ACKNOWLEDGMENT**

First of all, we indebted to the GOD ALMIGHTY for giving me an opportunity to excel in our efforts to complete this project on time.

We express our sincere thanks to our respected dean **Dr. Abdul Sharief**, Dean, School of Engineering, Presidency University for getting us permission to undergo the project.

We record our heartfelt gratitude to our beloved professor **Dr. C. Kalaiarasan**, University Project-II in charge, Associate Dean-Admin, Department of Computer Science and Engineering, Presidency University for rendering timely help for the successful completion of this project.

We are greatly indebted to our guide **Ms. Sushmita Kumari**, Department of Computer Science and Engineering, Presidency University for his/her inspirational guidance, valuable suggestions and providing us a chance to express our technical capabilities in every respect for the completion of the project work.

We thank our friends for the strong support and inspiration they have provided us in bringing out this project.

REFERENCES

- [1] Wang Wanbin. Design and Implementation of Electronic Medical Record Management System. *China Computer&Communication*.2011(7):26-28.
 - [2] Yang Yanli, CaoYan, PangYandan. Electronic medical record management system development based on B/S structure. *Microcomputer and its applications*.31(15):87-89(2012).
 - [3] Li Duo, Fan Qinzhi. Discussion on the unified authentication of digital campus. *Journal o Jilin Normal University (Natural Science Edition)*.8(3):154-156(2012).
 - [4] Zhang Renhui, Wang Xiaoming. Design and realization on management system of electronic medical record. *Microcomputer Information*.25(1-3):267-269(2009)
 - [5] Li Bo. Design and research of hospital electronic medical record management system based on B/S architecture. *Electronic Design Engineering*.25(5):46-49(2017).
 - [6] Yang Lu. Design and Achievement of Emergency EMR Management System. *China Digital Medicine*.12(1):86-88(2017).
 - [7] Y. Yang, A. Wang, X. Zhao, C. Wang, L. Liu, H. Zheng, Y. Wang, Y. Cao, and Y. Wang, "The Oxford shire Community Stroke Project classification system predicts clinical outcomes following intravenous thrombolysis: a prospective cohort study," *There Clin Risk Manga*, vol. 12, pp. 1049–1056, 2016.
 - [8] M. L. Schmitz, C. Z. Simonsen, M. L. Svendsen, H. Larsson, M. H. Madsen, I. K. Mikkelsen, M. Fisher, S. P. Johnsen, and G. Andersen, "Ischemic stroke subtype is associated with outcome in thrombolized patients," *Acta Neurol. Scand.*, vol. 135, no. 2, pp. 176–182, Feb. 2017.
 - [9] P. Sommer, A. Posekany, W. Serles, M. Marko, S. Scharer, E. Fertl, J. Ferrari, W. Lang, M. Vosko, S. Szabo, S. Kiechl, M. Knoflach, S. Greisenegger, Austrian Stroke Unit Registry Collaborators, "Is Functional Outcome Different in Posterior and Anterior Circulation Stroke?," *Stroke*, vol. 49, no. 11, pp. 2728–2732, Nov. 2018.
 - [10] J. F. Meschia, "Addressing the heterogeneity of the ischemic stroke phenotype in human genetics research," *Stroke*, vol. 33, no. 12, pp. 2770–2774, Dec. 2002.
- . Mill Valley, CA: University Science, 1989.