# DAYANANDA SAGAR COLLEGE OF ENGINEERING
# COMPUTER SCIENCE & ENGINEERING

## Minor Project- Report
## Aug-2021-2022

Course Faculty:  Dr. Vinothini C
Course Name & Code: Cloud & Big Data Laboratory with Minor Project & 18CS7DLCBL
Semester: 7
Date: 19-1-2022

| TITLE OF THE PROJECT | Heart Disease Analysis | | |
|---|---|---|---|
| | | | |
| STUDENT NAME | Gautam Kumar | Rithvik K Bhat | Vignesh K |
| USN | 1DS18CS710 | 1DS18CS732 | 1DS18CS744 |
| INDIVIDUAL CONTRIBUTION | AWS<br>Random Forest Algorithm | Decision Tree Algorithm<br>Random Forest Algorithm | AWS<br>Decision Tree Algorithm |
| GUIDE | Mr. Ravichandra H | | |
| | | | |
| PROJECT ABSTRACT | Heart Disease prediction is one of the most complicated tasks in medical field. Data science plays a crucial role in processing huge amount of data in the field of healthcare. As heart disease prediction is a complex task, there is a need to automate the prediction process to avoid risks associated with it and alert the patient well in advance. This project uses the heart disease data set available on the AWS cloud S3 platform. The project predicts the chances of heart disease and classifies patient's risk level by implementing two data mining techniques they are Decision Tree and Random Forest. Thus, this project analyses the performance of the two machine learning algorithms. The trial results verify that Random Forest algorithm has achieved the highest accuracy compared to Decision Tree. | | |
| PLATFORM USED (H/W & S/W TOOLS TO BE USED | Jupyter Notebook & Amazon Web Services S3 | | |
| | | | |

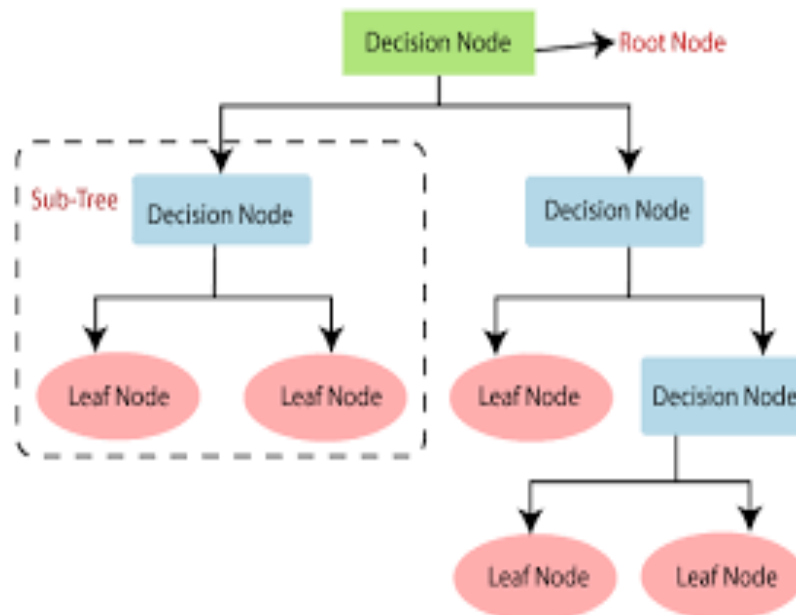| | |
|---|---|
| **INTRODUCTION** | Human heart is the main part of the human body which regulates blood flow throughout our body. Any irregularity to the heart can cause distress in other parts of the body and can be classified as a heart disease. Today, heart disease is one of the primary reasons for the occurrence of most deaths. Heart disease may occur due to unhealthy lifestyle, smoking, alcohol and high intake of fat which may cause hypertension.<br><br>The main challenge in today's healthcare is provision of best quality services and accurate diagnosis. The accuracy in management of a disease lies in the proper time of detection of that disease. The project makes an attempt to detect these heart diseases at an early stage to avoid disastrous consequences.<br>Records of large set of medical data are available for analyzing and extracting valuable knowledge from it. Data mining techniques are used for extracting valuable information from the large data available. The medical database consists of discrete information which makes decision making a complex and tough task. Machine Learning (ML) which is a subfield of data mining handles large scale well-formatted dataset efficiently. In the medical field, machine learning can be used for diagnosis, detection and prediction of various diseases. The main goal of our project is to provide a tool for doctors to detect heart disease at an early stage. This in turn will help to provide effective treatment to patients and avoid severe consequences. ML plays a very important role to detect the hidden discrete patterns and analyze the given data. This project shows the performance analysis of two ML techniques: Decision Tree and Random Forest for predicting heart disease at an early stage. |
| | |
| **DESIGN** | <br>Fig. 1 –  Decision Tree Algorithm |

Fig.2 – Random Forest Algorithm

| | |
|---|---|
| PROJECT SOURCE CODE LINK (GITHUB/ GOOGLE DRIVE) | https://github.com/gautamK007/Heart-Disease-Analysis |
| | |
| CONCLUSION /FUTURE ENHANCEMENT | Predicting heart disease is one of the most difficult problems in medicine. In the realm of healthcare, data science is critical for examining massive amounts of data. Because predicting cardiac illness is a difficult undertaking, it is necessary to automate the process in order to avoid the risks connected with it and to inform the patient well in advance. This project makes use of the AWS cloud S3 platform's heart disease data set. The study uses two data mining approaches, Decision Tree and Random Forest, to predict the likelihood of heart disease and classify the risk level of patients. As a result, the performance of the two machine learning algorithms is examined in this study. The trial results show that the Random Forest algorithm gives the best results. For future enhancements, the project can be used to a variety of machine algorithms. Using visuals to detect cardiac disease could be a future improvement as well. |
| | |

UI SCREENSHOTS