

# Heart Failure Prediction

Kartik Chhipa, Rushil Ashish Shah, Ruthvik K

B20CS084, B20AI036, B20AI037

**Abstract** – *This paper reports our experience with building a classifier that predicts if a person will have heart failure or not depending on a given set of features. We are given a dataset containing several feature vectors that provide complete information about the person. We are supposed to apply machine learning based models for classification and predict if that person will have heart failure or not.*

## I. INTRODUCTION

Cardiovascular diseases kill approximately 17 million people globally every year, and they mainly exhibit myocardial infarctions and heart failures. Heart failure is a common event caused by CVDs and this dataset contains 12 features that can be used to predict mortality by heart failure. Most cardiovascular diseases can be prevented by addressing behavioral risk factors such as tobacco use, unhealthy diet and obesity, physical inactivity and harmful use of alcohol using population-wide strategies. People with cardiovascular disease or who are at high cardiovascular risk (due to the presence of one or more risk factors such as hypertension, diabetes, hyperlipidaemia or already established disease) need early detection and management wherein a machine learning model can be of great help.

### **Datasets**

We have used the dataset present on the kaggle for this project.

It contains 12 columns and 918 Rows .

The columns present in the dataset are: Age, Sex, ChestPainType, RestingBP, Cholesterol, FastingBS, RestingECG, MaxHR, ExerciseAngina, Oldpeak, ST\_Slope and HeartDisease

1. Explanation of Columns in the dataset
2. Age - age of the person in years
3. Sex - sex of the person (M: Male, F: Female)
4. ChestPainType - chest pain type
5. (TA: Typical Angina, ATA: Atypical Angina, NAP: Non-Anginal Pain, ASY: Asymptomatic)
6. RestingBP: resting blood pressure (mm Hg) Cholesterol: serum cholesterol (mm/dl)
7. FastingBS: fasting blood sugar (1: if FastingBS > 120 mg/dl, 0: otherwise)
8. RestingECG: resting electrocardiogram results (Normal: Normal, ST: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV), LVH: showing probable or definite left ventricular hypertrophy by Estes' criteria)
9. MaxHR: maximum heart rate (Numeric value between 60 and 202)
10. ExerciseAngina: exercise-induced angina (Y: Yes, N: No)
11. Oldpeak: oldpeak = ST (Numeric value measured in depression)
12. ST\_Slope: the slope of the peak exercise ST segment (Up: upsloping, Flat: flat, Down: downsloping)
13. HeartDisease: output class (1: heart disease, 0: Normal)

## II. METHODOLOGY

### **OVERVIEW**

There are various classification algorithms present out of which we shall implement the following

- Random Forest Classification
- KNN
- Gradient Boosting
- Gaussian Naive Bayes
- XGBoost
- Neural Network

## Exploring the dataset and pre-processing

On counting the number of NULL values in the train dataset , it was found that there are no NULL values present. The column data types for latent vectors are float and for glasses field is integer which suits our requirements. Hence , there is no need for general pre-processing and the dataset is well-built.

## Implementation of classification algorithms

*A better approach for implementing classifiers is to objectively search different values for model hyperparameters and choose a subset that results in a model that achieves the best performance on a given dataset. This is called **hyperparameter optimization** or **hyperparameter tuning** and here we use **RandomizedSearchCV** for the same.*

- **Random Forest Classifier** : Random Forest Classifiers use boosting ensemble methods to train upon various decision trees and produce aggregated results. It is one of the most used machine learning algorithms.

After applying **RandomizedSearchCV** for both the data's i.e original one and the one after removing the outliers we get the best parameters for the classifier and now we implement the models with those parameters.

Best parameters obtained for the data with outliers:

roc\_auc\_score for the model with outlier data : 0.9461979913916787

roc\_auc\_score for the model without outliers : 0.9572649572649573

As expected we can clearly see the increase in Accuracy of the model after removing the outliers.

- **KNN (*k* - nearest neighbors)** : KNN are supervised algorithms which classify on the basis of distance from similar points. Here *k* is the number of nearest neighbors to be considered in the majority voting process.

After applying **RandomizedSearchCV** for both the data's i.e original one and the one after removing the outliers we get the best parameters for the classifier and now we implement the models with those parameters.

roc\_auc\_score for the model with outlier data : 0.9189383070301291

roc\_auc\_score for the model without outliers : 0.9262912482065997

As expected we can clearly see the increase in Accuracy of the model after removing the outliers.

- **Gradient Boosting** : Gradient Boosting Classifiers are a group of machine learning algorithms that combine many weak learning models together to create a strong predictive model. Decision trees are usually used when doing gradient boosting. Gradient boosting models are becoming popular because of their effectiveness at classifying complex datasets

After applying **RandomizedSearchCV** for both the data's i.e original one and the one after removing the outliers we get the best parameters for the classifier and now we implement the models with those parameters.

roc\_auc\_score for the model with outlier data : 0.9530726924916308

roc\_auc\_score for the model without outliers : 0.9634971509971509

As expected we can clearly see the increase in Accuracy of the model after removing the outliers.

- **Gaussian Naive Bayes** : GNB is a type of Naive Bayes classifier that assumes that the distribution of data is gaussian and classifies data based on this assumption.

After applying **RandomizedSearchCV** for both the data's i.e original one and the one after removing the outliers we get the best parameters for the classifier and now we implement the models with those parameters.

roc\_auc\_score for the model with outlier data : 0.9282639885222381

roc\_auc\_score for the model without outliers : 0.9601139601139601

- **XGBoost Classifier** : XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements Machine Learning algorithms under the Gradient Boosting framework. It

*provides a parallel tree boosting to solve many data science problems in a fast and accurate way*

After applying RandomizedSearchCV for both the data's i.e original one and the one after removing the outliers we get the best parameters for the classifier and now we implement the models with those parameters.

roc\_auc\_score for the model with outlier data : 0.9458393113342898

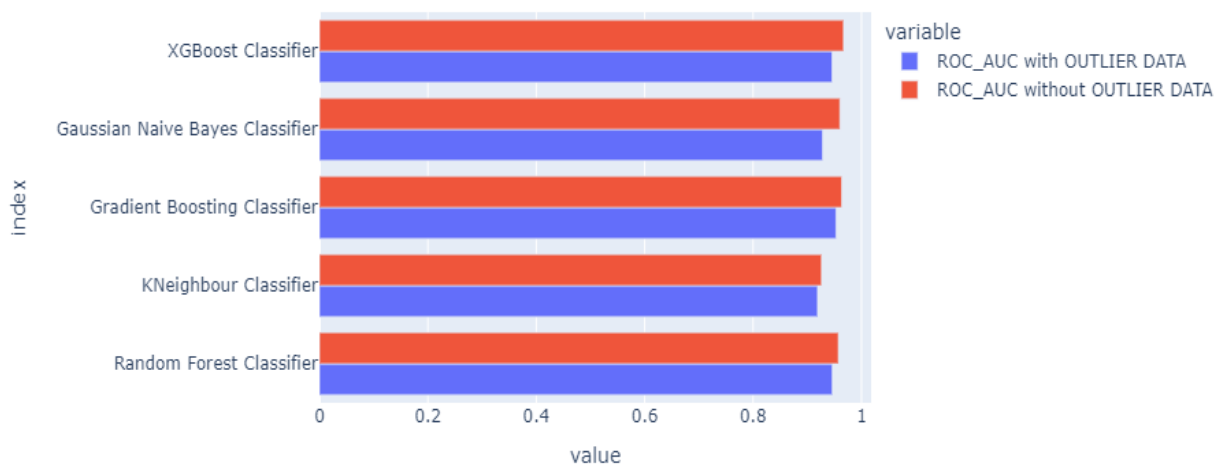
roc\_auc\_score for the model without outliers : 0.9668803418803419

#### IV. EVALUATION OF MODELS

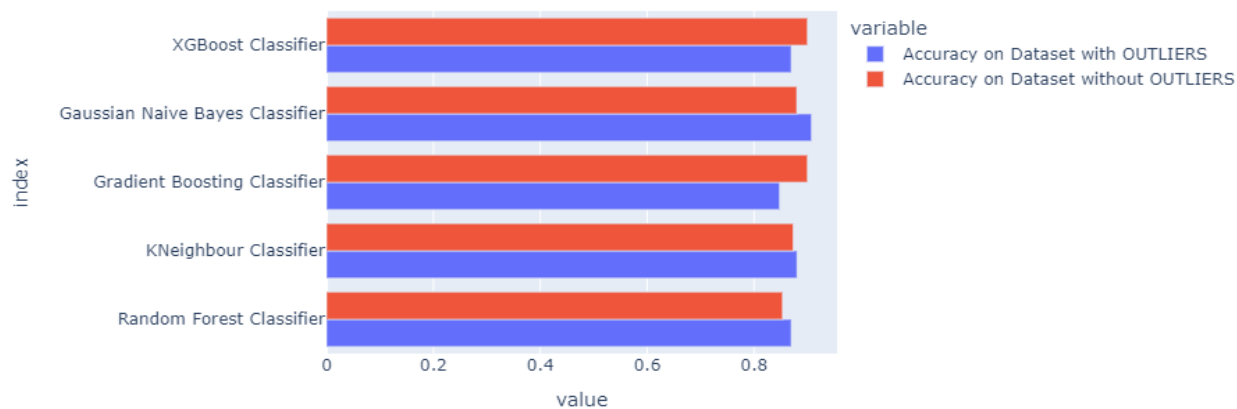
The models implemented were evaluated using techniques like - Classification report : precision , recall , f1 score and support , Confusion matrix , ROC plots , accuracy score and cross validation scores.

	ROC_AUC with OUTLIER DATA	ROC_AUC without OUTLIER DATA
Random Forest Classifier	0.946198	0.957265
KNeighbour Classifier	0.918938	0.926291
Gradient Boosting Classifier	0.953073	0.963497
Gaussian Naive Bayes Classifier	0.928264	0.960114
XGBoost Classifier	0.945839	0.966880

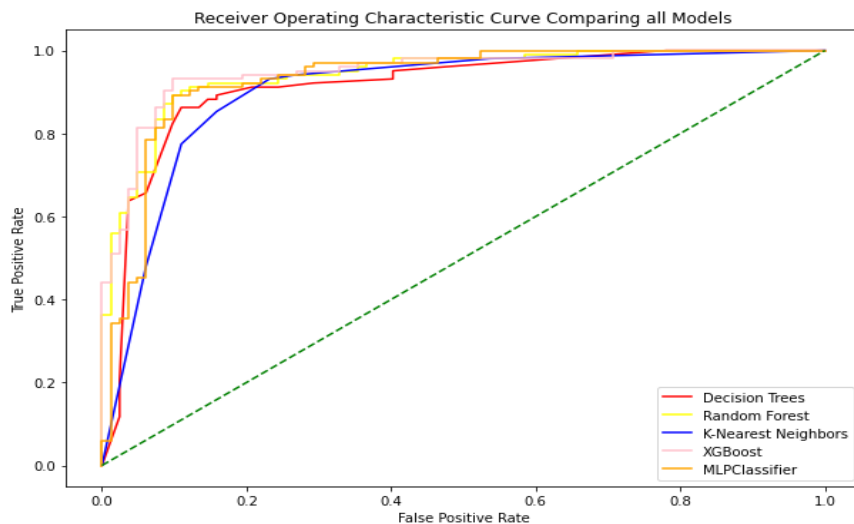
#### Visualizing roc\_auc scores of the models implemented



## Visualizing the Results of Pipeline



### ROC - Plot



## VI. RESULTS AND ANALYSIS

The table shows that all classifiers have had nearly equally efficient performance. Also we have observed that the results have significantly improved after applying exploratory data analysis. We can see a significant amount of increase in roc\_auc\_score after removing the outliers from the data for all the classifier models that we have implemented.

XGBoost has turned out to be the best classifier when evaluated on the basis of roc\_auc score which is equal to 96.68. It is marginally better than Gradient boosting algorithm. Kneighbour classifier has come out to be the classifier with least roc\_auc score equal to 92.6.

Link to the webapp : <https://share.streamlit.io/kartikchhipa01/heartfailureprediction/main.py>

## CONTRIBUTIONS

The learning and planning was done as a team. The individual contributions are as given

- Kartik Chhipa (B20CS084): Exploratory Data Analysis (EDA), Web app, Compiling individual Work.
- Rushil Ashish Shah(B20AI036) : Artificial Neural Network, Pipeline, Gaussian Naive Bayes, Report
- Ruthvik K (B20AI037): Random Forest Classifier, KNeighbors, Gradient Boosting, *XGBoost*, Report.

#### REFERENCES

- [1] K Nearest Neighbor | KNN Algorithm | KNN in Python & R (analyticsvidhya.com)
- [2] Pattern Classification -Book by David G. Stork, Peter E. Hart, and Richard O. Duda
- [3] Understanding AUC - ROC Curve | by Sarang Narkhede | Towards Data Science
- [4] Link for image dataset : Heart Failure Prediction | Kaggle