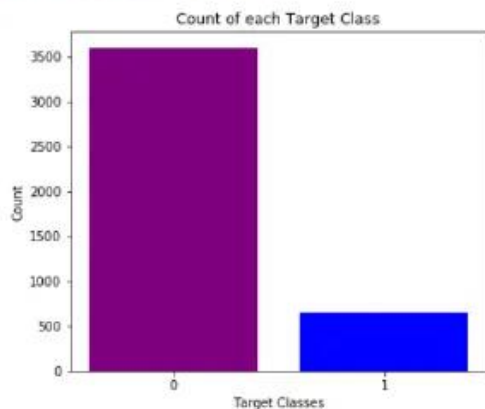# Visualizing and Predicting Heart Diseases with an Interactive Dash Board
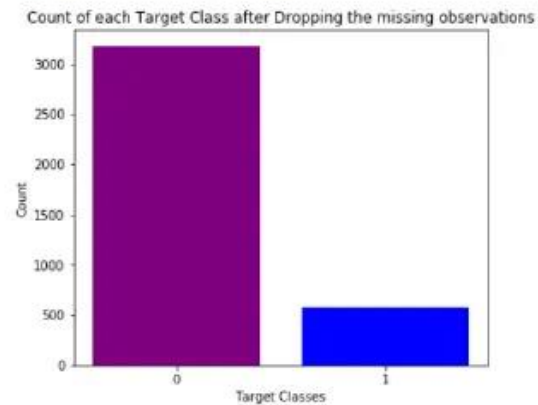
## Problem Solution Fit

### Data Preparation:

Since the dataset consists of 4240 observations with 388 missing data and 644 observations to berisked for heart disease, two different experiments were performed for data preparation. First, wechecked by dropping the missing data, leaving with only 3751 data and only 572 observationsrisked for heart disease.





This leads to reduced number of the observations providing irrelevant training to our model. So,we progressed with imputation of data with the mean value of the observations and scaling themusing SimpleImputer and StandardScaler modules of Sklearn

## Exploratory Analysis Document:

Correlation Matrix visualization Before Feature Selection shows



It shows that there is no single feature that has a very high correlation with our target value.Also, some of the features have a negative correlation with the target value and some have positive. The data was also visualized through plots and bar graphs.

## Discussion On Results:

When performing various methods of feature selection, testing it was found that backwardelimination gave us the best results among others. The various methods tried were BackwardElimination with and without KFold, Recursive Feature Elimination with Cross Validation. Theaccuracy that was seen in them ranged around 85% with 85.5% being maximum. Though bothmethods gave similar accuracy but it was seen that in Backward Elimination we found that thenumber of misclassifications of True Negative was more and it was observed that the accuracy hadmore variance compared to RFEV. The precision of Backward Elimination and RFEV are 84%and 86% respectively. And the recalls are 0.99

and 1 respectively. The precision and recall also shows that the number of misclassifications is less in RFECV than in Backward Elimination.

| Evaluation Metrics | Backward Elimination | RFECV |
|---|---|---|
| Accuracy | 83% | 85% |
| Recall | 0.99 | 0.99 |
| Precision | 0.84 | 0.86 |

Table 3: Comparison between the feature selection models after training and testing through LogisticRegression model

## Contributions:

| Task \ Members | Tharun R | Raman KB | Pavan kumar B | Kamal R |
|---|---|---|---|---|
| Data Imputation and Scaling | ■ | | | ■ |
| Data Cleaning | | | ■ | |
| Exploratory Analysis | | ■ | | ■ |
| Feature Selection | ■ | | ■ | ■ |
| Building Model | ■ | | ■ | |
| Result analysis and Accuracy Test | | ■ | | |
| Documentation | ■ | ■ | ■ | ■ |

Table 4: Work Division

## Conclusion:

The early prognosis of cardiovascular diseases can aid in making decisions on lifestyle changes inhigh risk patients and in turn reduce the complications,

which can be a great milestone in the field of medicine. This project resolved the feature selection i.e. backward elimination and RFECV behind the models and successfully predict the heart disease, with 85% accuracy. The model used was Logistic Regression. Further for its enhancement, we can train on models and predict the types of cardiovascular diseases providing recommendations to the users, and also use more enhanced models