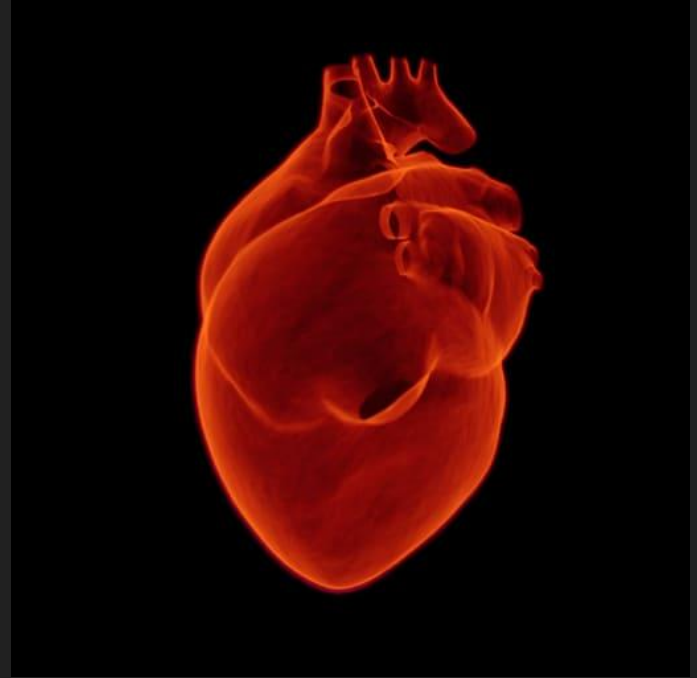


Predicting Heart Disease

Using ML and the UCI Heart Disease Dataset

Impact

- Leading cause of death in the US
- 1 in every 5 deaths (US) is a direct result of the condition
- Someone dies of it every 30-45 seconds
- Costs the country between 2 and 3 billion dollars every year (i.e., healthcare services, medicines, lost productivity due to death, etc.)



Identification

- Identification and prevention of HD is one of the most important topics in the healthcare industry
- Primary prevention models are not accurate enough to be effective.
- Low estimates put their accuracy at just under 60%
- High estimates place them around 80%
- ML models are being used to increase the accuracy of heart disease identification and improve treatment outcomes



Project Goal

- Come up with a supervised learning model that could outperform the current gold standard in primary prevention models
- Use the UCI Heart Disease Dataset
- Predict heart disease with greater than 80% accuracy



**American
Heart
Association®**

Dataset

- UCI Heart Disease Dataset
- Sixteen columns of data and 920 observations
- Four sub-datasets merged to form this dataset, and each sub-dataset corresponds to a different hospital
 - Cleveland
 - Hungary
 - Switzerland
 - VA Long Beach
- Data excluded from modeling process
 - Three predictive features that were missing a third or more of their values
 - About 6% of the patients that were missing 40% or more of their data
- Final dataset had 865 patients and 10 predictive features (five continuous and five binary features)

Continuous

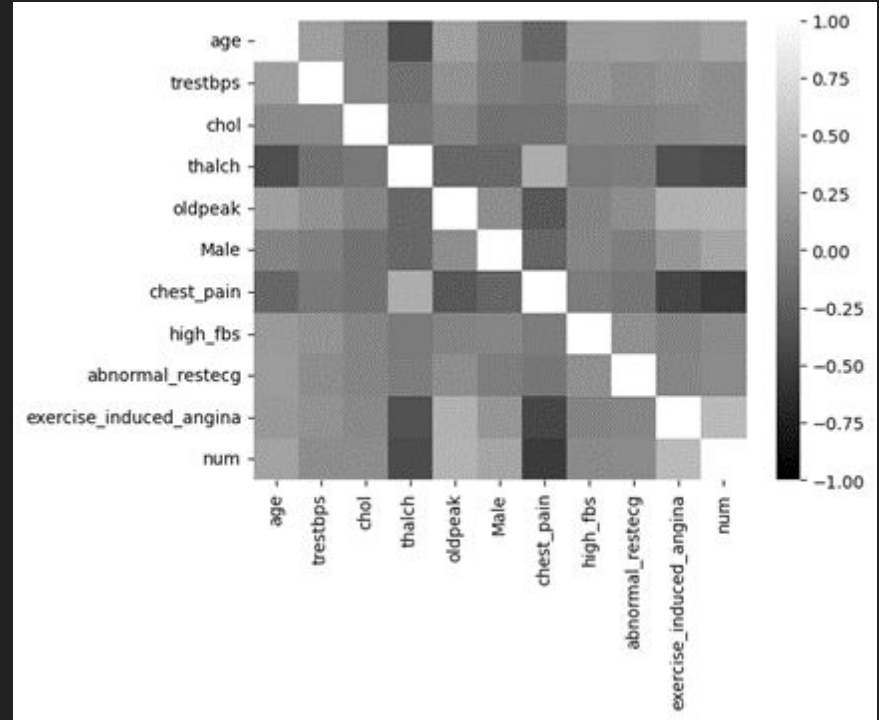
- Age - Age of the patient in years
- Trestbps - Resting blood pressure (in mm Hg) on admission to the hospital
- Chol - Serum cholesterol (in mg/dl)
- Thalch - Maximum heart rate achieved
- Oldpeak - ST depression induced by exercise relative to rest

Binary

- Male - Whether the patient is a male
- Chest-pain - Whether the patient is experiencing chest pain
- High-fbs - Whether fasting blood sugar is above 1200 mg/dl
- Abnormal_restecg - Whether the electrocardiogram results were abnormal
- Exercise_induced_angina - Whether the patient had exercise induced angina

Exploratory Data Analysis

- Project focused on classifying patients as (1) “has heart disease” and (0) “does not have heart disease”
- Bonferroni corrected t-tests of (continuous features) showed significant differences between groups at the .05 level for all five features
- Bonferroni corrected chi-squared tests (binary features) showed significant differences at the .05 level for all but one of the features (high_fbs)
- Thalch and chest pain were the only two features not positively correlated with “has heart disease”



Model Training - Preprocessing

- Data was first split into training (70%) and testing (30%) sets
- Missing values were imputed with either the mean or the median.
- If necessary, features were then scaled



Model Training - Initial Evaluation

- Five different models were evaluated during training:
 1. Logistic Regression
 2. AdaBoost with a Decision Tree estimator
 3. AdaBoost with an SVM Estimator
 4. AdaBoost with a Logistic Regression Estimator
 5. Random Forest
- Hyperparameters were selected using either GridSearchCV or RandomizedSearchCV (n_iter = 180)
- One of each model was selected for the final model selection process
- Hyperparameters for the best models were determined by the highest cross-validation F1 score during training

$$\mathbf{F1\ Score} = \frac{2}{\left(\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}} \right)}$$

$$\mathbf{F1\ Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Model Training - Final Selection

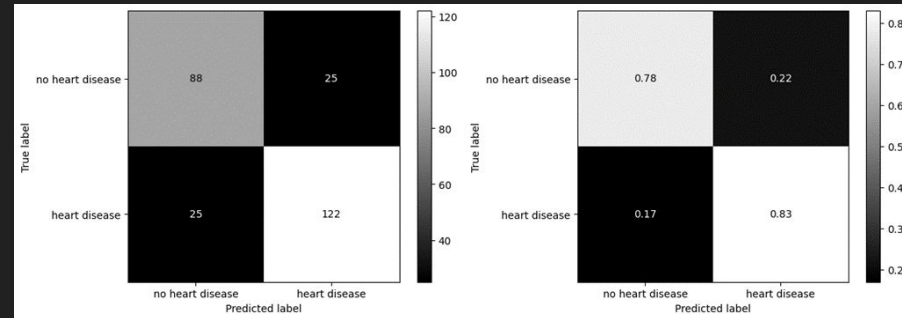
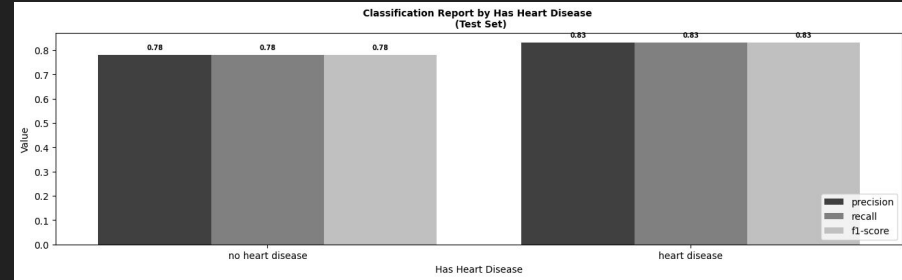
- Model with the highest F1 score on the testing set was chosen as the best model
- Adaboost model with Decision Tree estimator (F1 score = 0.83), narrowly outperformed the other Adaboost models
- Final model did not use any scaling and missing values were imputed using the median

```
Hyper-parameters
-----
AdaBoostClassifier(estimator=DecisionTreeClassifier(ccp_alpha=0.01, max_depth=4,
                                                    max_features='log2',
                                                    min_impurity_decrease=0.0001,
                                                    min_samples_leaf=5,
                                                    min_samples_split=4),
                  learning_rate=0.01, n_estimators=300))]

verbose: False
simpleimputer: SimpleImputer()
standardscaler: None
adaboostclassifier: AdaBoostClassifier()
simpleimputer__add_indicator: False
simpleimputer__copy: True
simpleimputer__fill_value: None
simpleimputer__keep_empty_features: False
simpleimputer__missing_values: nan
simpleimputer__strategy: median
adaboostclassifier__algorithm: SAMME.R
adaboostclassifier__base_estimator: deprecated
adaboostclassifier__estimator_ccp_alpha: 0.01
adaboostclassifier__estimator_class_weight: None
adaboostclassifier__estimator_criterion: gini
adaboostclassifier__estimator_max_depth: 4
adaboostclassifier__estimator_max_features: log2
adaboostclassifier__estimator_max_leaf_nodes: None
adaboostclassifier__estimator_min_impurity_decrease: 0.0001
adaboostclassifier__estimator_min_samples_leaf: 5
adaboostclassifier__estimator_min_samples_split: 4
adaboostclassifier__estimator_min_weight_fraction_leaf: 0.0
adaboostclassifier__estimator_random_state: None
adaboostclassifier__estimator_splitter: best
adaboostclassifier__estimator: DecisionTreeClassifier()
adaboostclassifier__learning_rate: 0.01
adaboostclassifier__n_estimators: 300
adaboostclassifier__random_state: None
```

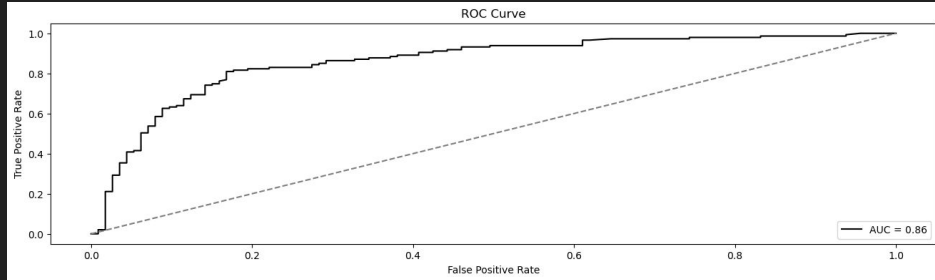
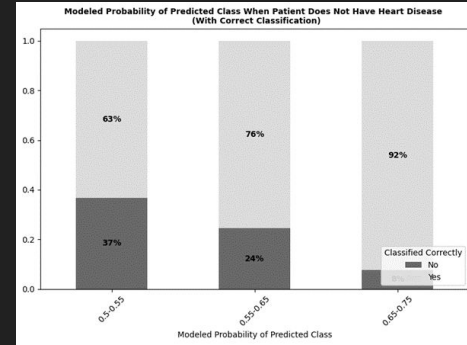
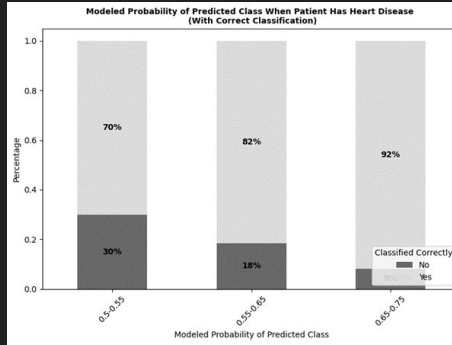
Testing Set - Further Analysis

- Model was better at classifying heart disease patients than those without the condition
- Better to make this type of error (Type I) since patients are better off taking precautionary measures with their lifestyle if they are a fringe case



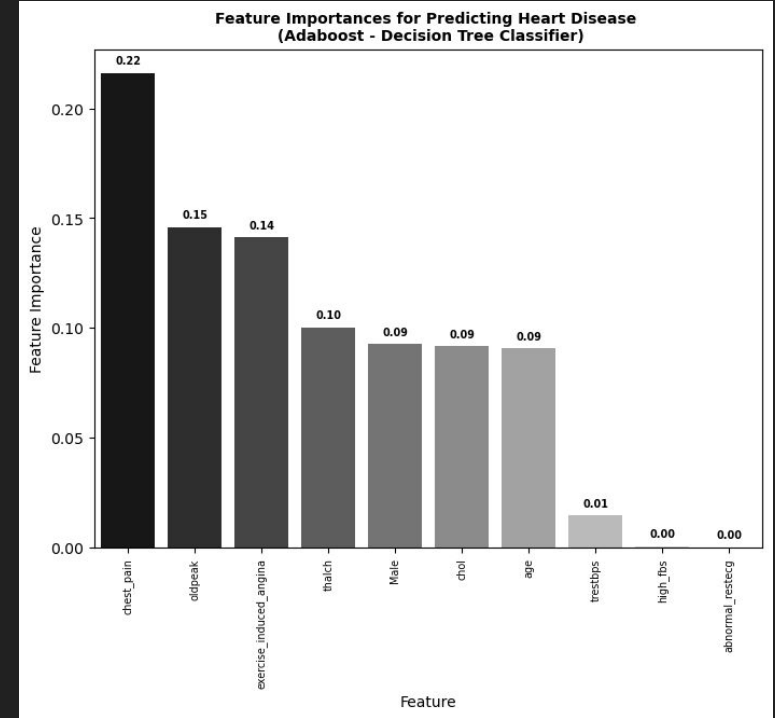
Testing Set - Further Analysis (continued)

- Model performed better when it was more confident in a prediction
- Effect was more pronounced for cases when the patient did not have heart disease
- Plotting the true positive rate against the false positive rate showed that it was doing a really good job overall (despite differences between classes)



Feature Importances

- Random Forest estimator gave insight into which features were most important (ordered below):
 1. Whether or not a patient had chest pain
 2. ST depression induced by exercise relative to rest (oldpeak)
 3. Whether a patient had exercise induced angina
- Maximum heart rate achieved (thalach), sex, cholesterol, and age were also important features



Final Thoughts and Future Research

- Model met initial goal of predicting heart disease with $> 80\%$ accuracy
- Was better at identifying cases when patients had heart disease than when they did not
- Ideas for future research:
 1. Add other features related to exercise (2 of the top 3 features were exercise-related)
 2. Incorporate three predictive features that were dropped at the beginning of the project
 3. Add data from other hospitals
 4. Explore whether the predicted probabilities could be used to come up with a risk score



Thank you.