

Predicting Presence of Cardiovascular Disease in Patients

Adam Guo(hg426), Karan Maheswari(km857), Krishna Raju(kk992)

December 10th, 2019

1 Intro

In this project, we are concerned with what factors attribute to a person having heart disease. Heart disease is a chronic disease that takes a great amount of resources to treat, and if handled improperly could take away someone's life easily. By figuring out which factors have the most probabilities of leading to heart diseases, we could monitor those factors for patients, and thus prevent heart disease. For this project, we would seek to build a predictive model using data from the 'Cardiovascular Disease dataset' from Kaggle, which contains heart disease information of around 70,000 patients.

1.1 Goal

Our goal for the project is two-fold. First, we want to come up with a machine learning model to determine which factors (features) in the dataset contribute the most to the patient having cardiovascular disease. Second, we want to be able to predict the risk of a new patient having cardiovascular disease given his/her data, and maybe come up with a confidence interval for the prediction. To evaluate our model performance, we want to look at the test error for the model, and also compare the important factors we got from the model to some clinical proven results.

1.2 Data

We are using the '*Cardiovascular Disease Dataset*' from Kaggle. Our data set has 12 features and 70,000 entries. The features in the dataset are of 3 types -

Objective - Factual Information

Examination - Results of medical examination

Subjective - Information given by patient

It is important to distinguish between the three types, so that we can decide how much 'reliance' we can put on each feature. For example, subjective data like alcohol intake might not always be accurate because people might falsely report (knowingly or unknowingly) their medical conditions.

Below we list all our features with their type and some other useful properties:

1. Age | Objective Feature | age | int (days)
2. Height | Objective Feature | height | int (cm) |
3. Weight | Objective Feature | weight | float (kg) |
4. Gender | Objective Feature | gender | categorical code |
5. Systolic blood pressure | Examination Feature | ap_hi | int |
6. Diastolic blood pressure | Examination Feature | ap_lo | int |
7. Cholesterol | Examination Feature | cholesterol | 1: normal, 2: above normal, 3: well above normal |
8. Glucose | Examination Feature | gluc | 1: normal, 2: above normal, 3: well above normal |
9. Smoking | Subjective Feature | smoke | binary |
10. Alcohol intake | Subjective Feature | alco | binary |
11. Physical activity | Subjective Feature | active | binary |
12. Presence or absence of cardiovascular disease | Target Variable | cardio | binary |

2. Procedure

2.1 Data Exploration

There are 70,000 examples and 11 features which describe the target variable.

There isn't any missing data but there are some outliers. We observed this in exploratory data analysis.

```
df.describe()
```

| | age | gender | height | weight | ap_hi | ap_lo | cholesterol | gluc | smoke | alco | acr |
|-------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| count | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 | 70000.000000 |
| mean | 53.837914 | 1.349571 | 164.359229 | 74.205690 | 128.817286 | 96.630414 | 1.366871 | 1.226457 | 0.088129 | 0.053771 | 0.800000 |
| std | 6.766821 | 0.476838 | 8.210126 | 14.395757 | 154.011419 | 188.472530 | 0.680250 | 0.572270 | 0.283484 | 0.225568 | 0.390000 |
| min | 30.000000 | 1.000000 | 55.000000 | 10.000000 | -150.000000 | -70.000000 | 1.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 49.000000 | 1.000000 | 159.000000 | 65.000000 | 120.000000 | 80.000000 | 1.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| 50% | 54.000000 | 1.000000 | 165.000000 | 72.000000 | 120.000000 | 80.000000 | 1.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| 75% | 59.000000 | 2.000000 | 170.000000 | 82.000000 | 140.000000 | 90.000000 | 2.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| max | 65.000000 | 2.000000 | 250.000000 | 200.000000 | 1600.000000 | 1100.000000 | 3.000000 | 3.000000 | 1.000000 | 1.000000 | 1.000000 |

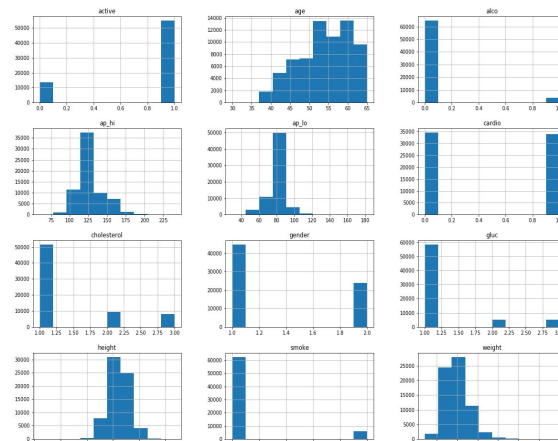
- Here we can see that weight feature (kgs) has a minimum value of 10. For a 30-65 year old person having a weight of 10 kgs is not possible.
- Similarly we see in height a person with a height of 55 cm. This indicates presence of corrupt values
- Another set of features which have corrupt values are ap_hi (Systolic Blood Pressure) and ap_lo (Diastolic Blood Pressure). After carefully studying the data we saw that there are garbage values of negative and very high blood pressure.

There are two possible solutions to this problem: either we can drop the samples with corrupt value or we can replace the value with median/mean values. Since our dataset is large enough and our application is health critical we chose to drop the examples with corrupt entries.

```
df_mask1.describe()
```

| | age | gender | height | weight | ap_hi | ap_lo | cholesterol | gluc | smoke | alco | acr |
|-------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| count | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 | 68630.000000 |
| mean | 19464.711132 | 1.348667 | 164.400685 | 74.129738 | 126.678537 | 81.305872 | 1.364709 | 1.225776 | 0.087964 | 0.053359 | 0.800000 |
| std | 2467.962214 | 0.476552 | 7.961879 | 14.297368 | 16.695001 | 9.465487 | 0.678920 | 0.571647 | 0.283245 | 0.224749 | 0.390000 |
| min | 10798.000000 | 1.000000 | 104.000000 | 35.450000 | 60.000000 | 30.000000 | 1.000000 | 1.000000 | 0.000000 | 0.000000 | 0.000000 |
| 25% | 17658.000000 | 1.000000 | 159.000000 | 65.000000 | 120.000000 | 80.000000 | 1.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| 50% | 19701.000000 | 1.000000 | 165.000000 | 72.000000 | 120.000000 | 80.000000 | 1.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| 75% | 21324.000000 | 2.000000 | 170.000000 | 82.000000 | 140.000000 | 90.000000 | 2.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| max | 23713.000000 | 2.000000 | 207.000000 | 200.000000 | 240.000000 | 182.000000 | 3.000000 | 3.000000 | 1.000000 | 1.000000 | 1.000000 |

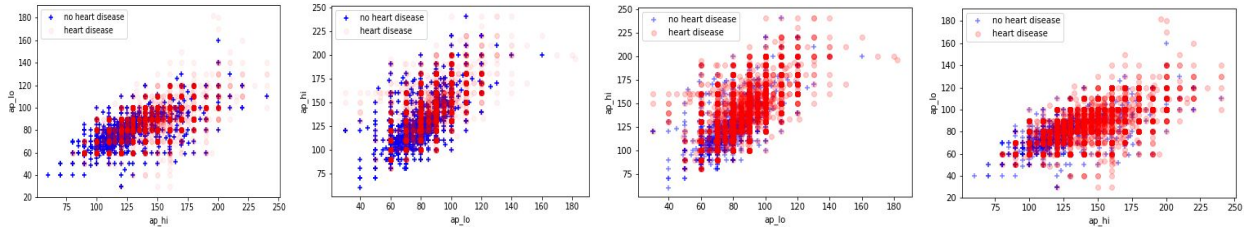
Then we drew the histograms of all the features, which are useful for us to understand each feature and its distribution.



Looking at the target we see that dataset is balanced (equal number of examples for both ‘heart disease’ and ‘no heart disease’) and not skewed.

Next we go on to realise that the number of men and women isn’t uniformly distributed and thus before making any conclusion about correlation of heart disease with gender we should take this fact in account.

It is commonly known that blood pressure is correlated to heart disease, and so we wanted to model the relationship. After drawing some scatter plots :

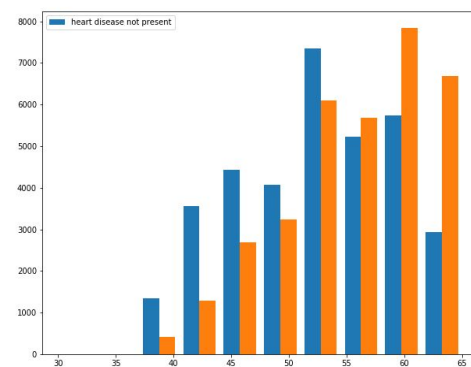


We can see a trend here that high blood pressure is somewhat correlated to the presence of heart disease.

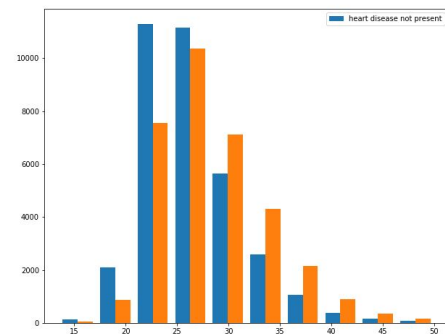
Also we drew a histogram for age with presence of heart disease:

We can clearly see a trend here, as age is increasing people with heart disease increases

It can be seen that at lower age relatively fewer people have heart disease, and at higher age relatively few people do not have heart disease.



We couldn’t find any relation between height and presence of disease or weight and presence of disease so we did a polynomial transformation of height and weight which is Body Mass Index ($\text{height}/\text{weight}^2$). Here again we can very well see a trend in the figure below. At lower BMI fewer people have heart diseases and as obesity increases, the number of heart diseases also increases.



2.2 Feature Engineering

To better prepare our features to be used for our model, we perform feature engineering on some of our model variables. The purpose of doing feature engineering is that some algorithms require features with very specific characteristics to work properly, and sometimes feature engineering can even drastically improve the performance of an existing model. Therefore, it is crucial that we explore some options for feature engineering in our model, to fit better models, and to improve performance. Here, we will break down the feature engineering for our model, by the methods used:

1. Handling Outliers: Having outliers in our data features could prevent the machine learning algorithm from learning the model very well. And therefore, it is crucial that we prevent those outliers from appearing in our data features. Upon reading resources from online, we set thresholds for a few of the data features:
 - a. $50 \leq \text{'ap_hi'} \leq 250$ (systolic blood pressure)
 - b. $40 \leq \text{'ap_lo'} \leq 160$ (diastolic blood pressure)
 - c. $100 < \text{'height'} < 219$
 - d. $\text{'weight'} > 35$

2. **Date Extraction:** Sometimes the data format given in the dataset is not ideally for model input, and therefore we have to do some binning in order to have the right units for the data variable. In our case, the 'age' column of each sample is originally in days, as a double-typed variable (eg. 12125.95). Having this large input feature could potentially cause the model to underperform, and therefore, we perform feature engineering by dividing the 'age' by 365, and then taking the ceiling of the results.
3. **Feature Imputation:** We noticed that 'weight' and 'height' features could be combined into one to reduce the dimension of the input space, so we decided to create an additional variable - 'BMI' that combines the two previous features. The BMI, also known as Body Mass Index, is calculated as the body weight divided by the height (in meters) squared. We kept 2 decimal places after the division, and thresholded the 'BMI' values also.

3 Model Development

3.1 Model Selection

We decided to split the data 85-15, meaning 85% of the entire dataset will go to training dataset, while 15% of the entire dataset will go to testing.

For this model selection portion, we are using Python's Sklearn library, which contains a lot of useful modules for a majority of the machine learning models. Here are a list of models that we attempted:

- Linear model that uses Stochastic Gradient Descent
- Random Forest Classifier
- Logistic Regression
- K-Nearest Neighbor Classifier
- Gaussian Naive Bayes Classifier
- Perceptron Classifier
- Support Vector Machine
- Decision Tree Model

After training the models, we need to evaluate their training performance and decide on which metrics to be used to judge that. We decided to use the classification accuracy (prediction == actual) to compare each model that we trained, because our output label is a 2-class classification label, and the accuracy should be sufficient to do a simple comparison. The mean accuracy of each model is shown here:

We notice here that Random Forest and a Decision Tree Model actually achieve the best training performance. Both methods have similar structure - nodes for representing class label and branch for representing prediction outcome. A decision tree model as opposed to a random forest model will be nice because the decision tree is a white-box model, meaning that it is easily explainable, and one can clearly see where the decision is being split.

Out[19]:

| | Model |
|-------|---------------|
| Score | |
| 96.30 | Random Forest |
| 96.30 | Decision Tree |
| 74.81 | KNN |
| 72.57 | Log Reg |
| 70.68 | Naive Bayes |
| 63.24 | Perceptron |
| 50.62 | SGD |
| 50.54 | SVM |

A random forest model, on the other hand, is just the ensemble method of the decision tree model. It builds multiple decision trees given random subsets of the entire dataset, and eventually average the results. Averaging here means each tree at the end will vote for the classification label, and the label that received most votes by simple majority will be the prediction. Compared to a decision tree model, random forest model doesn't overfit as often, because by subsetting the dataset and averaging the results, one effectively lower the variance in the final prediction. And therefore, to avoid potentially overfitting the dataset and to better be able to generalize the model to real world data, **we decide to use a random forest model**. We accept the inherent trade-off of using a black-box model, i.e, our model might not be interpretable but we could expect better results than a simpler white box model like decision trees. Our choice is also incentivised by the fact that, the basic corelation between our features and the risk of heart

disease is already known. Therefore, it would not be of much value for us to build a white-box model which would simply reaffirm already known predictors of heart disease. Instead, we attempt to build a model which can outperform many of these simple models and do a better job at predicting heart disease.

3.2 Cross Validation & Overfitting

Given the high training accuracy that we are seeing from above, we were a little skeptical of the model's performance on real data. To give us a sense of how the model will perform on real data, we decided to use a 10-fold cross validation on the training dataset using the random forest model, to check for potential overfitting. Here is the result we get from the 10-fold cross validation:

```
Scores of model 2 with X_train 2: [0.69704264 0.69222834 0.69893398 0.70323246 0.70392022 0.6966466
0.69406707 0.70106639 0.71448228 0.71517028]
Mean of scores2: 70.16790259900283
Std Deviation of scores2: 0.007452267115960629
```

As we can see from the output, the mean accuracy of the 10 models is only around 70.17%, which is much lower than the training accuracy we saw earlier (>90%). Therefore, there is clearly some overfitting that we are doing earlier with the training of the random forest. Here we will describe how we battled overfitting.

Normally, the methods we discussed in class to reduce overfitting were to either

1. Get more data to fit the model
2. Reduce the amount of features
3. Use a less complex model

We couldn't possibly get more data to fit our model, so we chose two methods to reduce the overfitting. The first was to use less features in our model - we decided to drop additionally 'Alco', 'Smoke', and 'Active' from our model. The second was actually a little unconventional - because random forest is an ensemble method that meant to reduce variance, we decided to do some hyper-parameter tuning (eg. to train more trees in the forest) to reduce overfitting.

Hyperparameters, different from model parameters, are configurations external to the model and cannot be estimated or trained from the data. They are usually manually set by the person running the model, and given the predictive modeling problem, they can be fine-tuned in order to better the performance of the model. For our random forest model, we decided to perform a grid search, in order to find the optimal hyperparameter that results in the best prediction.

Here are some of the hyperparameters we can specify for our random forest model:

- max_depth: the maximum depths of the trees in the forest
- criterion: the scoring function used to construct the trees (gini or entropy)
- n_estimators: number of trees to train in the forest
- max_features: the maximum amount of features in each tree
- max_leaf_nodes: the maximum amount of leaf nodes in each tree
- min_samples_split: the minimum amount of samples required to split an node

Here are some possible values (grids) of the hyperparameters that we hope to tune:

What the grid searching is doing is that it iteratively trains the model based on unique combinations of those configurations, then it notes down the accuracy and variance of each model. At the end it returns the set of

configuration that gave us the best results before. During grid search, in fact a cross-validation process is used to make sure the performance difference is resulted from the configurations, but not the trained weights(threshold) of the model.

```
param_grid = {
    'max_depth' : [5,10,15,20,30],
    'criterion' : ['gini', 'entropy'],
    'n_estimators' : [50,100,200,400,500,700,1000],
    'max_features': ['auto', 5, 8],
    'max_leaf_nodes': [100,200,300,500,700]
    'min_samples_split' : [20,50,100,200]
}
```

Here is the result returned from grid search:

```
GridSearchCV(cv='warn', error_score='raise-deprecating',
             estimator=RandomForestClassifier(bootstrap=True, class_weight=None, criterion='entropy',
                                              max_depth=20, max_features='auto', max_leaf_nodes=200,
                                              min_impurity_decrease=0.0, min_impurity_split=None,
                                              min_samples_leaf=1, min_samples_split=2,
                                              min_weight_fraction_leaf=0.0, n_estimators=700, n_jobs=-1,
                                              oob_score=False, random_state=1, verbose=0, warm_start=False),
             fit_params=None, iid='warn', n_jobs=-1,
             param_grid={'min_samples_split': [20, 50, 100, 200]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring=None, verbose=0)
```

We can notice that the best `n_estimators`, the number of trees in the forest was returned to be 700, much larger than what we used to train the model before(100). Also, notice that the `max_depth`, maximum depth of the tree is set to 20, further restricting the complexity of the decision trees in the forest. We will use this set of hyperparameters to retrain our model and get some results.

4 Testing

4.1 Retrain & Testing Model

Given the hyperparameters set returned from the grid search, we retrained our random forest model, which resulted in a training accuracy of 74.47%. This lowered accuracy is expected because we don't want our model to fit every bit of details of the training data too well, but instead want it to be able to have good performance for data we haven't seen before.

We later ran the model on the test dataset.

Here are the accuracy numbers that we achieved:

Train Accuracy: 74.47

Test Accuracy: 73.56523433693852

The test accuracy is 73.565%, higher than from what we saw before (68%) and close to that of training. This is good because we have relatively high training and testing accuracy, and because they are close together, it could mean that our model would perform well on real world data. However, as discussed before, this also resulted in lowered training accuracy, and potentially more bias in the model, and we understand that the **bias-variance-tradeoff** is in effect here.

4.2 Model Performance Comparison

Given our increased bias, we wanted to see how other previous efforts have fared in tackling the same classification tasks, especially using other machine learning models.

We found one instance of a team solving the same task of cardiovascular disease prediction using a Neural Network solution. The best test accuracy they achieved after tuning and optimization was 72%. Here is a link to their project: <https://www.kaggle.com/camiloqq/nn-to-predict-cardiovascular-diseases>

We also noticed from online that the highest accuracy to have achieved on this dataset was around 72% or 73%, and our solution could potentially be the new gold standard. Therefore, we decided to stick with our models, features, and the hyperparameters that we chose.

4.3 Further Results Investigation

Because our prediction task involves a binary classification in the medical field, we also want to look at statistical measures specific for this field, such as sensitivity and specificity. Here are the definition of both terms:

- Sensitivity: also the recall or true positive rate, is the proportion of true positives that are correctly identified. In our case, it will be the percent of heart disease patients who are actually predicted as such.

- **Specificity:** also the true negative rate, is the proportion of true negatives that are correctly identified. In our case, it will be the percent of healthy people who are correctly predicted as such.

Both of those values are important, because in most settings, there are threshold requirements for sensitivity and specificity for medical classification tests to achieve. In fact, classification tests with both values above 90% are considered to have high credibility, although there are often tradeoffs between the two values.

To look at the sensitivity and specificity for our model, we calculated the true positives and true negatives, and for references we also included the false negatives and false positives. We then calculated the specificity and sensitivity and printed the results:

```

Classification Report :

```

| | | | | precision | recall | f1-score | support |
|--|--------------|------|------|-----------|--------|----------|---------|
| | 0 | 0.71 | 0.78 | 0.75 | 5136 | | |
| | 1 | 0.76 | 0.69 | 0.72 | 5127 | | |
| | micro avg | 0.74 | 0.74 | 0.74 | 10263 | | |
| | macro avg | 0.74 | 0.74 | 0.73 | 10263 | | |
| | weighted avg | 0.74 | 0.74 | 0.74 | 10263 | | |

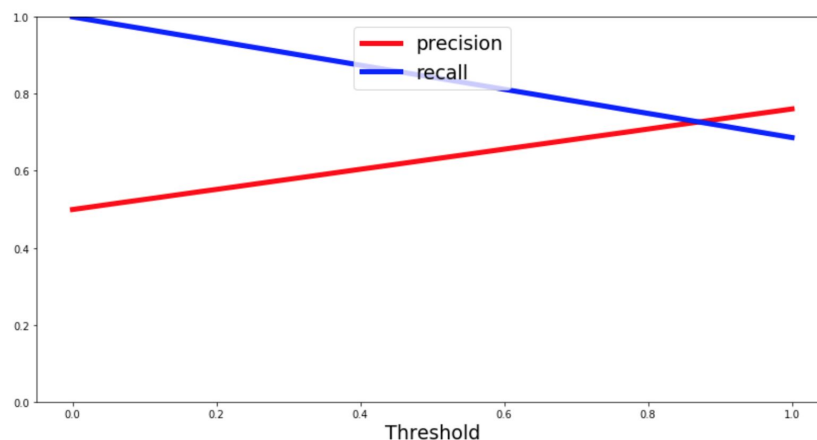
```

True Positive : 3520
True Negative : 4030
False Positive : 1106
False Negative : 1607
Specificity : 0.7846573208722741
Sensitivity : 0.6865613419153501

```

As we can see from the output above, the sensitivity of the model is 68.66%, and the specificity of the model is 78.47%. The sensitivity value is actually a little disappointing - the low value indicates that we are not doing a good job avoiding false negatives in our prediction. And in this context, it means that our model is not doing a good job predicting that heart disease patients actually have heart disease, and so that the disease is overlooked. This could potentially result in the patient not getting immediate attention to his/her health and getting the right treatment.

One way to combat this low sensitivity value is set a different threshold for classifying a patient as having heart disease. If we make the threshold lower, then we are more likely to predict someone having heart disease, regardless of the person actually having it or not. Here we showed a graph describing how the sensitivity and specificity change with the moving threshold values (recall is sensitivity):



We see that there is clearly a tradeoff for both values, and setting the threshold lower could potentially benefit our model by increasing the sensitivity but decreasing the specificity.

5. Discussion

5.1 Results and Limitations

In this part we will discuss the limitations to our modeling implementation in solving the heart disease classification task. Specifically, we will speak about this classification task in general, whether this model can be deployed, and it being a weapon of math destruction.

Solving the heart disease classification task has been attempted by many people in the past, using various machine learning algorithms. The fact that currently there are no existing implementations of such predictive models in the medical field showed that the problem space is hard to tackle. There are only so many data variables available about the patients that can be collected, and there are so many privacy and other concerns about obtaining more data. In an ideal world, it will be useful to have more features (eg. sweat level, heart disease history) to help predict the outcome.

The model that we have, achieved 73% accuracy at its best. Although this looks relatively low, our results are as good as those achieved by other groups using the same data. Yet, the model still cannot be deployed in the real world because of the sensitivity and specificity values that we have seen. The 69% and 78% are well below the required threshold to be used medically. This means that if our model is deployed, the false positives and false negatives result from the wrong predictions will either delay the heart disease patients of immediate treatment or give financial burden to people who are actually healthy.

Additionally, our model is potentially a weapon of math destruction. The section below will speak to this further.

5.2 Weapons of Math Destruction

In this section we discuss whether our model potentially fits the category of “Weapons of Math Destruction” based on how it is defined by Cathy O’Neil in her book “Weapons of Math Destruction”. There are a few questions that really help us categorize a model as a WMD. We pose these questions to our model and give our answers below.

Q) Is the model opaque or not interpretable?

A) Sadly, the answer is yes. As discussed earlier, in section 3.1, we chose a random forest model and accepted the tradeoff in favour of better performance at the cost of interpretability.

Q) Do the results of the model affect a large population (possibly in a negative way)?

A) Potentially yes. The ideal place of deploying such a model, would be as a preliminary check for heart disease that individuals can use when deciding whether or not go to a cardiologist. Such a deployment would surely affect a lot of people. Wrong results, could lead to unnecessary visits to the cardiologist or worse could deter patients who need medical attention from not seeking it.

Q) Does it create a feedback loop?

A) The answer is debatable, but here is our take - No it does not create a feedback loop. This is because a yes prediction would warrant a visit to the cardiologist who would further assess the situation, and a no prediction would mean the person is unlikely to have heart disease and so does not need to recheck using the model again. Note that this does not imply that wrong predictions are not harmful or dangerous, only that they don’t create a feedback loop.

Based on the above criteria, our model would fit under the category of a WMD. Thus, we would need to be extremely cautious of deploying our model into the real world.

5.3 Fairness

One part of fairness, which concerns false positive and false negatives has already been discussed in sections 4.3 and 5.1. Here we concern ourselves with fairness related to discrimination. Interestingly, determining fairness for the

given problem is complicated. This is because general notions of fairness, pertaining to algorithms not discriminating based on race, gender, age do not apply to this problem. This is primarily because race, gender and age have a significant correlation with heart disease and therefore it would make sense to use them as features. We cannot have fairness in terms of “disparate impact” either, because some groups are more likely to have heart disease than others. For example, men are more likely to have heart disease than women and black men are more likely to have heart disease than white or asian men. The fairness that is most applicable in our setting would be ‘individual fairness’, this would mean that an individual with similar features would get a similar output. It is generally hard to classify two instances as ‘similar’ without a strong understanding of what ‘similarity’ would mean in this context. For example, all things equal would patients with different exercise levels be considered similar or different? If we go ahead and assume similarity to mean exactly equal, then our model has ‘individual fairness’ by virtue of being a random forest. This is because if the features are exactly the same, every decision tree would output the same label for both of the instances and thus the final prediction would be similar. However, to get a better sense of fairness, we would have to loosen our definition of fairness and then compare predictions. This could only be done with the help of a medical professional who could identify two cases as similar.