# Diagnosing Heart Disease

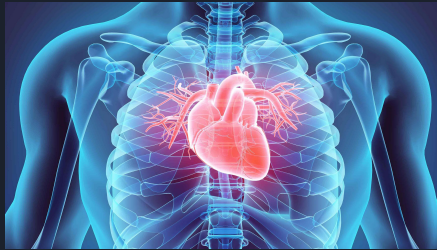University Circle, Inc (UCI) Cleveland Hospital

By Isabella Lindgren
Part Time Data Science Student

Hello everyone,
Today we are covering a topic I find extremely interesting and I hope you do too. We will be covering the topic of Heart Disease and how we can diagnose patients using Machine Learning.

# Background

- 647,000 Americans die from heart disease each year
- The leading cause of death for men, women, and most racial and ethnic groups in the US
- Costs the US ~$219 billion each year



Data from the CDC -
- About 647,000 Americans die from heart disease each year—that's 1 in every 4 deaths
- Heart disease is the leading cause of death for men, women, and people of most racial and ethnic groups in the United States.
- Heart disease costs the United States about $219 billion each year from 2014 to 2015. This includes the cost of health care services, medicines, and lost productivity due to death.

# Objective/Business Value

What factors have the greatest impact on heart disease diagnoses?

Which patients should be considered 'high risk?'

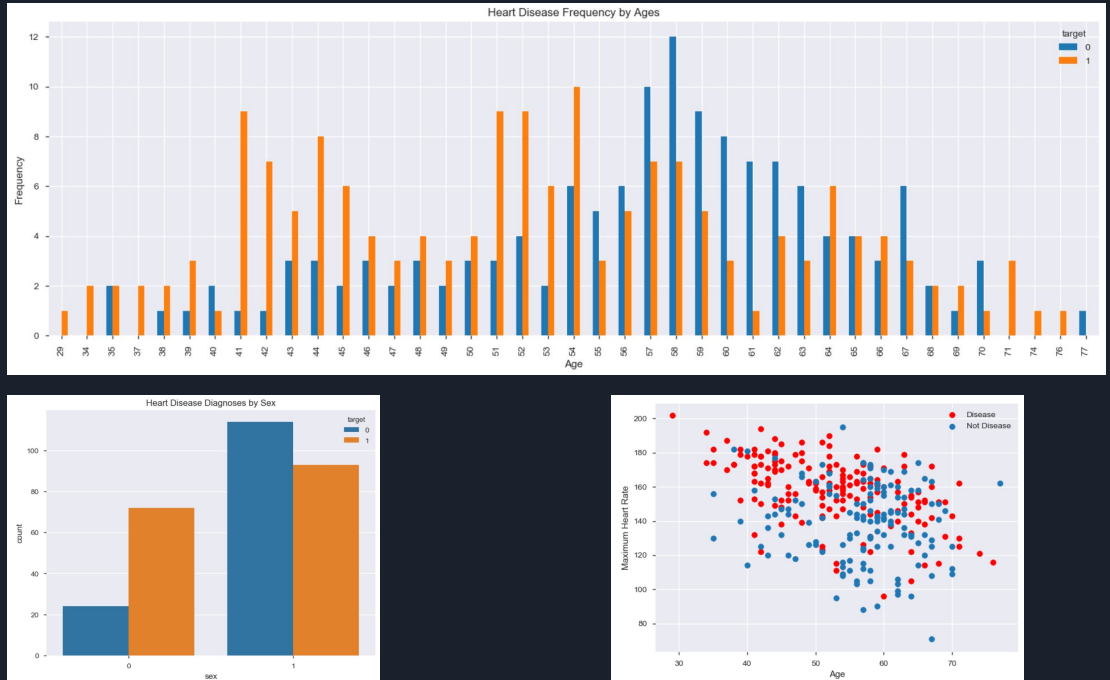How can we treat these high risk patients in a cost and time effective manner?

- Our objective is to determine which factors have the greatest impact on determining a heart disease diagnoses and use these factors to create a model to classify patients as 'high risk' or 'low risk'
- The findings from this model can help determine a more efficient way to allocate funding and to improve the health care system in regards to heart disease diagnosis and treatment.
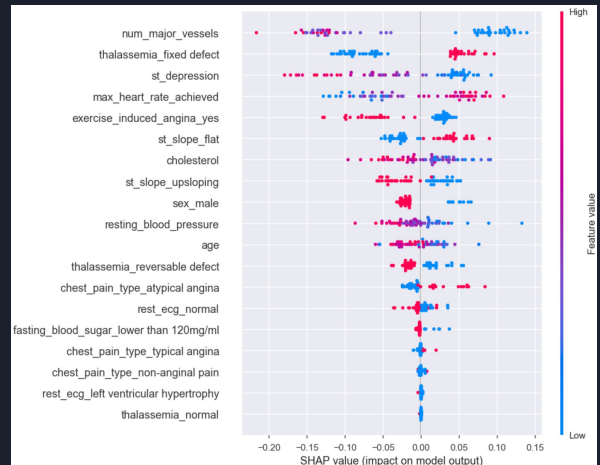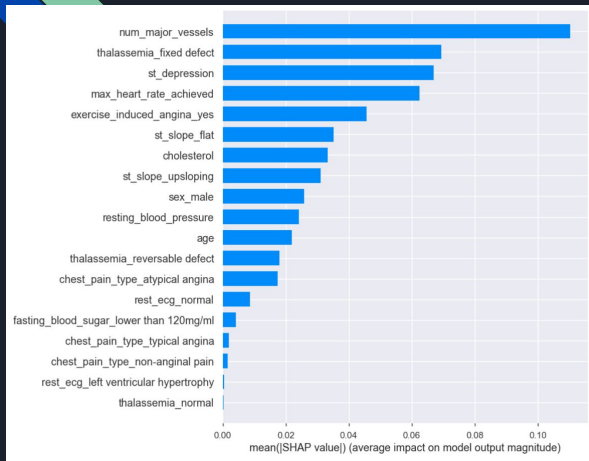
## Model

- Heart Disease UCI data from Cleveland Hospital
- Data: 303 patients, 20 attributes
- Diagnose heart disease with 84% Sensitivity, 73% Specificity

---

- Heart Disease UCI data from Cleveland Hospital
  Included data from 303 patients and 20 attributes including age, sex, chest pain type, resting blood pressure, cholesterol, blood sugar, testing results, etc.

- Our final model could determine with 84% Sensitivity and 73% Specificity whether a patient should be diagnosed with heart disease.

- In this dataset, we can see a trend of increasing frequency of heart disease with age
- Women were more likely to be diagnosed with heart disease compared to men in this dataset which is contradictory to the overall data of the united states. The CDC states that men are more likely to develop heart disease and at earlier ages than women.
- We can also see that the maximum resting heart rate is generally higher in patients who have heart disease, which makes sense since many heart conditions are caused by decreased blood flow, thus the heart works harder to distribute enough oxygen throughout the body.
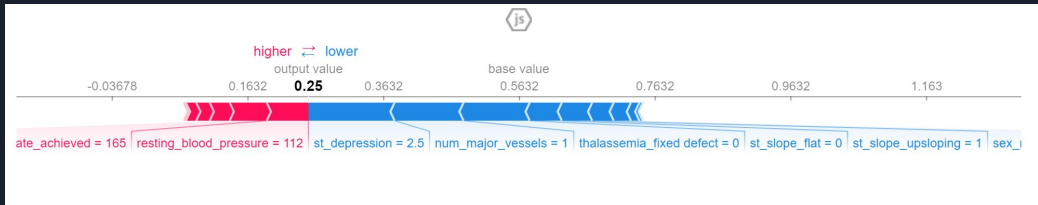
Top Features

- The features with the highest impact in our model were: number of major blood vessels, fixed Thalassemia defect, a depression in an ECG test, maximum heart rate, excercise induced angina,
- Number of major vessels had the most significant impact on our model

On the right
- The features with a clear divide were strong classifying features

# Example: Patient 1



# Example: Patient 12



- Ex1: Patient 1: ST- depression on ECG test, higher number of vessels, No Thalassemia worked in their favor. Their output value was below baseline, classifying them as 'low risk'
- T**halassemia** is a blood disorder passed down through families (inherited) in which the body makes an abnormal form or inadequate amount of hemoglobin
- Ex2: Patient 12: Low # major blood vessels, no ST-depression on ECG, high cholesterol and a Thalassemia fixed defect classified them as 'high risk'
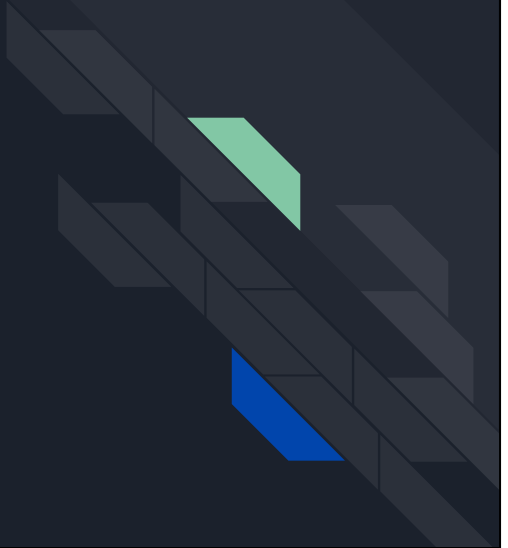
## Implications for this model

- Reduce invasive, unnecessary tests for low risk patients
- Preventative care for patients who are high risk
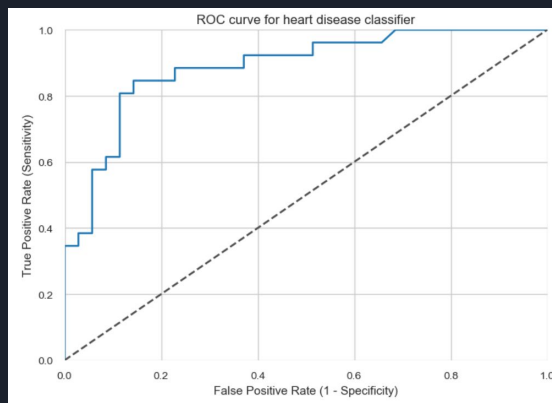- Assess the frequency of appointments depending on risk level

-

## Next Steps

- Use a larger data set across all states
- Time series data to look at trends
- Include race/ethnicity in data to ensure diversity and accurate representation of US population
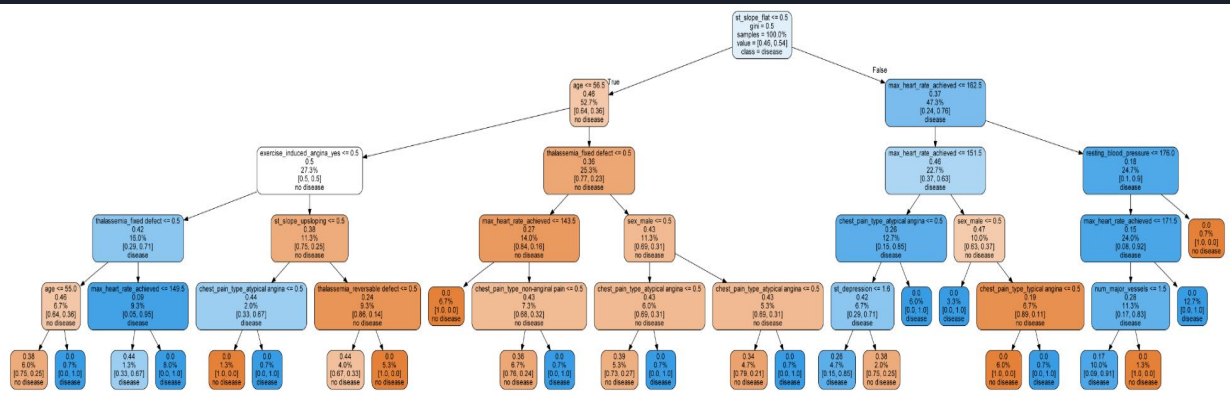- Principal Component Analysis on data

Thank you

# Appendix

| Weight | Feature |
|---|---|
| 0.0820 ± 0.0359 | thalassemia_fixed defect |
| 0.0361 ± 0.0321 | st_slope_flat |
| 0.0295 ± 0.0245 | age |
| 0.0230 ± 0.0161 | sex_male |
| 0.0230 ± 0.0334 | cholesterol |
| 0.0230 ± 0.0445 | max_heart_rate_achieved |
| 0.0197 ± 0.0525 | num_major_vessels |
| 0.0164 ± 0.0000 | resting_blood_pressure |
| 0.0098 ± 0.0161 | thalassemia_reversable defect |
| 0.0066 ± 0.0334 | exercise_induced_angina_yes |
| 0.0066 ± 0.0161 | chest_pain_type_atypical angina |
| 0.0066 ± 0.0161 | st_slope_upsloping |
| 0.0033 ± 0.0131 | chest_pain_type_non-anginal pain |
| 0 ± 0.0000 | thalassemia_normal |
| 0 ± 0.0000 | fasting_blood_sugar_lower than 120mg/ml |
| 0 ± 0.0000 | rest_ecg_left ventricular hypertrophy |
| 0 ± 0.0000 | rest_ecg_normal |
| 0 ± 0.0000 | chest_pain_type_typical angina |
| -0.0033 ± 0.0382 | st_depression |



ROC curve for heart disease classifier

```
total=sum(sum(confusion_matrix))

sensitivity = confusion_matrix[0,0]/(confusion_matrix[0,0]+confusion_matrix[1,0])
sensitivities['RFC'] = sensitivity
print('Sensitivity : ', sensitivity )

specificity = confusion_matrix[1,1]/(confusion_matrix[1,1]+confusion_matrix[0,1])
specificities['SVC'] = specificity
print('Specificity : ', specificity)

Sensitivity :  0.84375
Specificity :  0.7241379310344828
```