# Early Detection of Heart Disease

## Metis Classification Project

# Learning the heart way

- Heart disease is the leading cause of death in the United States

- ~659,000 people in the United States die from heart disease each year

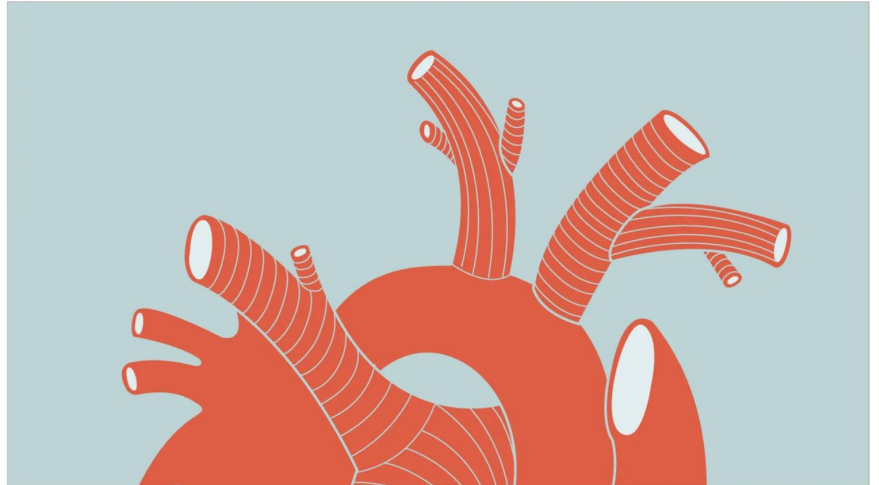- Every 40 seconds an American will have a heart attack



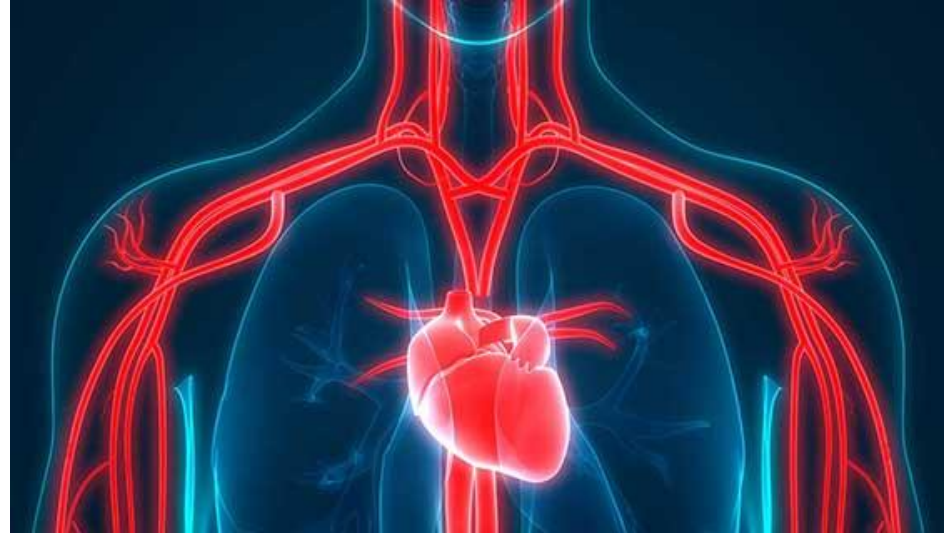Illustration by Ruth Basagoitia and Maya Chastain

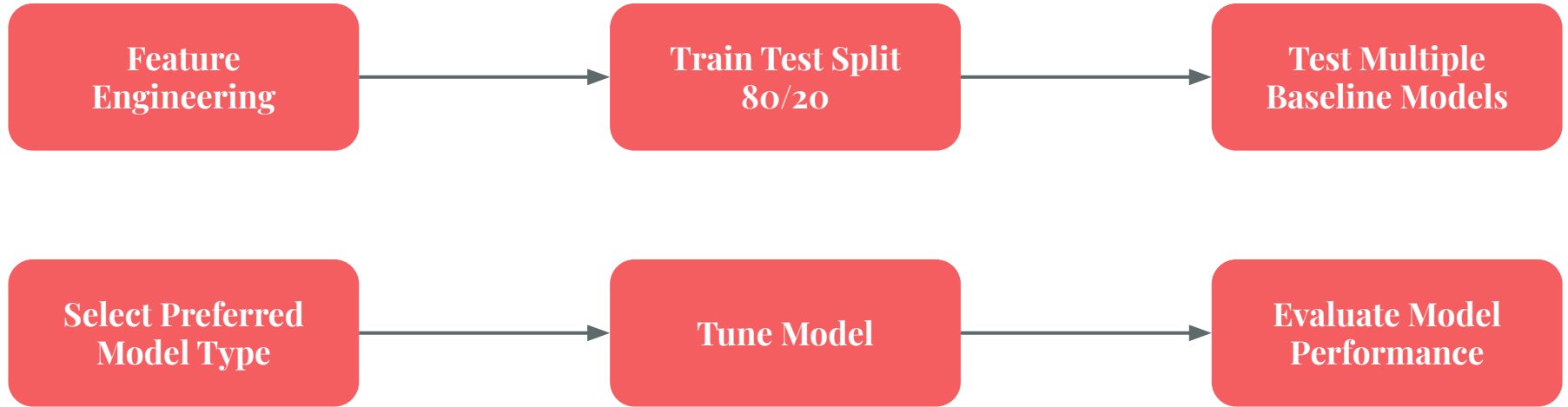https://www.cdc.gov/heartdisease/facts.htm

# Background

**GOALS**

- Implement classification model to predict likelihood of heart disease or a heart attack based on other health and lifestyle characteristics

- Identify who to target for early detection and intervention

- Understand feature importance to inform screening

**DATA**

- Obtained from Kaggle

- Comes originally from the CDC

- 22 different attributes for 253,680 individuals

# Modeling Approach

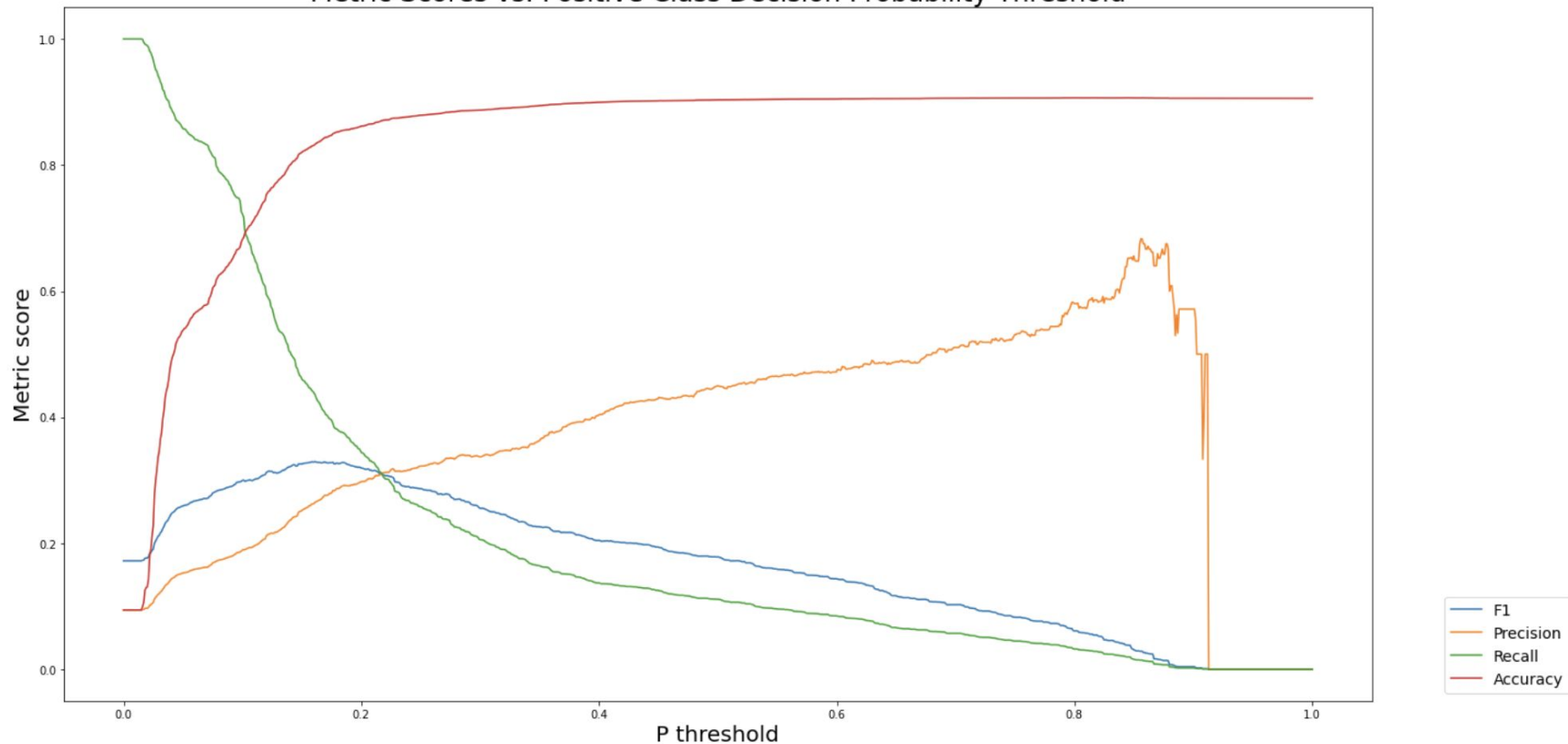| Feature Engineering | → | Train Test Split 80/20 | → | Test Multiple Baseline Models |
|:---:|:---:|:---:|:---:|:---:|

| Select Preferred Model Type | → | Tune Model | → | Evaluate Model Performance |
|:---:|:---:|:---:|:---:|:---:|

# Baseline Model Comparison

| Model | Neg Log Loss | Precision | Recall |
|---|---|---|---|
| Logistic Regression | -0.573144 | 0.439604 | 0.135823 |
| Random Forest | -0.585770 | 0.430109 | 0.130722 |
| XGB | -0.580914 | 0.450486 | 0.137719 |
| KNN | -0.566398 | 0.449111 | 0.136477 |

# Baseline Model Comparison

| Model | Neg Log Loss | Precision | Recall |
|---|---|---|---|
| Logistic Regression | -0.573144 | 0.439604 | 0.135823 |
| Random Forest | -0.585770 | 0.430109 | 0.130722 |
| XGB | -0.580914 | 0.450486 | 0.137719 |
| KNN | -0.566398 | 0.449111 | 0.136477 |

## Metric Scores vs. Positive Class Decision Probability Threshold

Best F1 score 0.330 at decision threshold >= 0.161
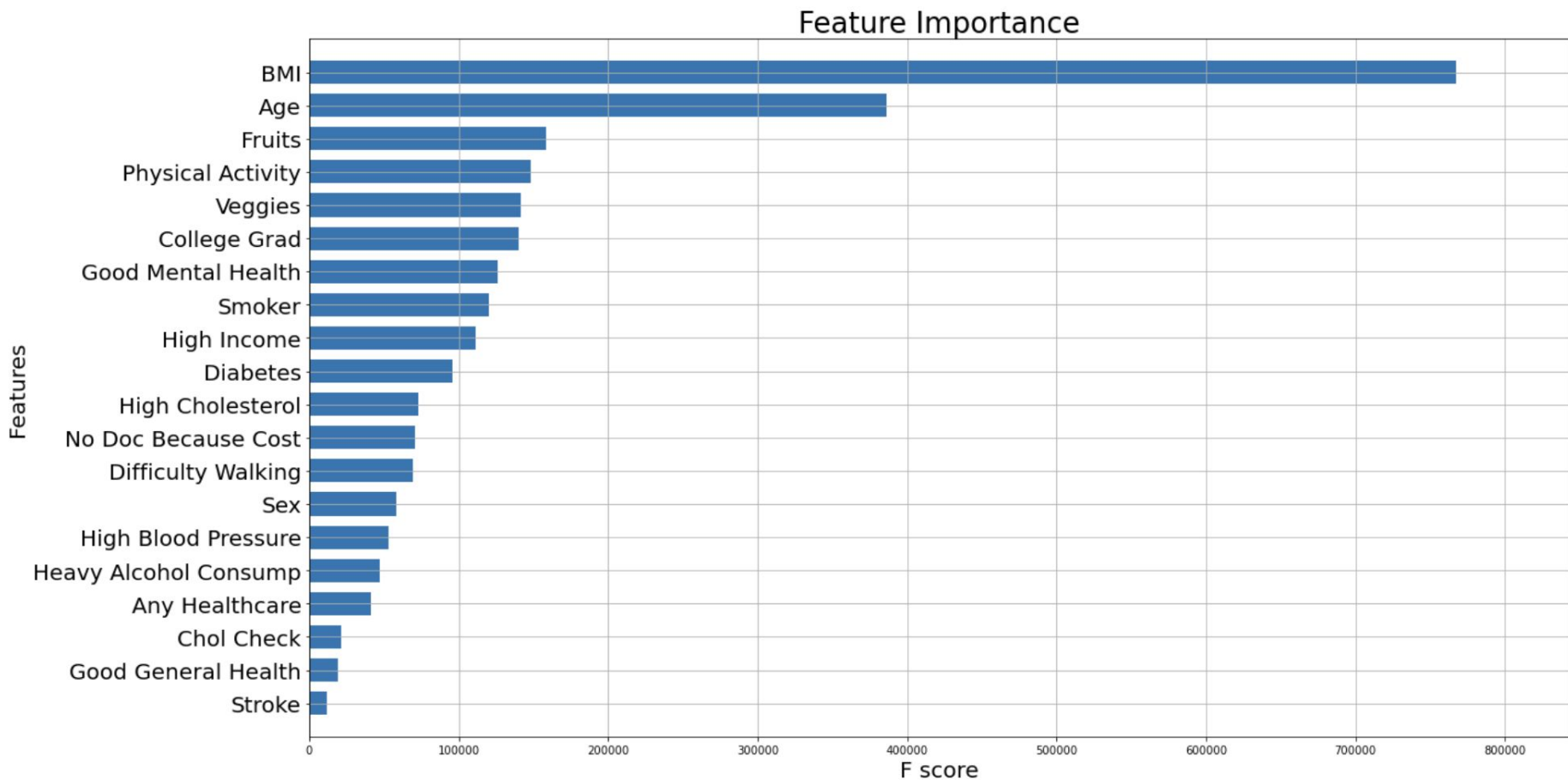
Confusion Martix

## Baseline Model

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 0.99      | 0.41   | 0.58     | 45957   |
| 1.0          | 0.15      | 0.97   | 0.25     | 4779    |
|              |           |        |          |         |
| accuracy     |           |        | 0.46     | 50736   |
| macro avg    | 0.57      | 0.69   | 0.42     | 50736   |
| weighted avg | 0.91      | 0.46   | 0.55     | 50736   |

## Tuned Model

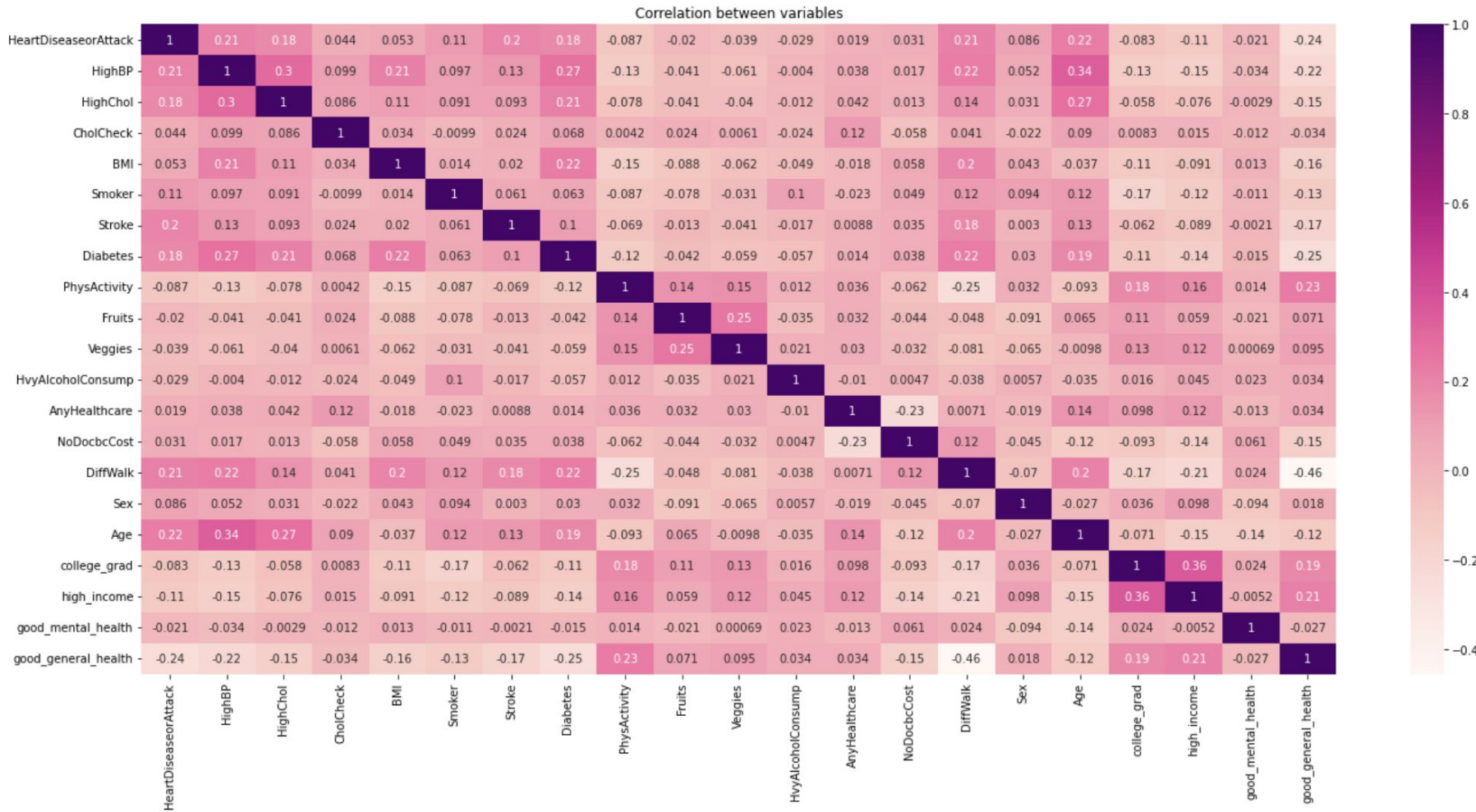|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0.0          | 0.98      | 0.63   | 0.76     | 45957   |
| 1.0          | 0.19      | 0.85   | 0.31     | 4779    |
|              |           |        |          |         |
| accuracy     |           |        | 0.65     | 50736   |
| macro avg    | 0.58      | 0.74   | 0.54     | 50736   |
| weighted avg | 0.90      | 0.65   | 0.72     | 50736   |

Feature Importance

# Conclusion

**KEY TAKEAWAYS**

- BMI and age are key features

- Minimizing false negatives comes with the cost of increasing false positives

**NEXT STEPS**

- What age should screening begin?

- What is the cost of a false positive and a false negative?

- Could ensembling improve prediction performance?

# Appendix

Correlation between variables

```python
#create dummy variable to indicate whether someone has a college degree
df['college_grad']= (df['Education'] == 6 ).astype(int)

#create variable to indicate whether someone makes more than $50,000 annually
df['high_income']= (df['Income'].isin([8,9]).astype(int))

#create variable to indicate whether people reported having poor mental health in more than 15 of the past 30 days
few_days = []
few_days.extend(range(1,15))
few_days.append(88)
df['good_mental_health']= (df['MentHlth'].isin(few_days).astype(int)) #88 is "None"

#create variable to indicate whether people report health as good/very good/excellent (1) vs fair/poor (0)
df['good_general_health']= (df['GenHlth'].isin([1,2,3]).astype(int))
```

```
y_test   rf and xgb agree
0.0      True                     30214
         False                    15743
1.0      False                     3319
         True                      1460
```

| | variables | vif |
|---|---|---|
| **0** | HighBP | 2.302320 |
| **1** | HighChol | 2.033872 |
| **2** | CholCheck | 21.481795 |
| **3** | BMI | 16.239262 |
| **4** | Smoker | 1.933629 |
| **5** | Stroke | 1.104151 |
| **6** | Diabetes | 1.416915 |
| **7** | PhysActivity | 4.581396 |
| **8** | Fruits | 3.022398 |
| **9** | Veggies | 5.701816 |
| **10** | HvyAlcoholConsump | 1.082131 |
| **11** | AnyHealthcare | 19.026871 |
| **12** | NoDocbcCost | 1.168536 |
| **13** | DiffWalk | 1.687309 |
| **14** | Sex | 1.872135 |
| **15** | Age | 9.822311 |
| **16** | college_grad | 2.112582 |
| **17** | high_income | 1.953899 |
| **18** | good_mental_health | 1.309966 |
| **19** | good_general_health | 7.375602 |