

Course: IST 707 Applied Machine Learning  
Title: Heart Disease Prediction

## Project Description:

Objective of the project is to understand the fundamental components that are the major causes of a heart attack which could help lower the risk of heart disorders. The goal is to use machine learning to predict if a person would have heart disease based on dataset's major key performance indicators (KPIs) like smoking, alcohol consumption, sex, stroke, age category, sleep time, BMI, mental health, physical health, etc.

## Dataset Description:

The dataset for this study was taken from Kaggle in the form of a csv file. The dataset originally was from the CDC and is a major part of the Behavioral Risk Factor Surveillance System (BRFSS). This was the most recent dataset (as of February 2022) at the time which included data from 2020. It consisted of 401,958 rows and 279 columns. The dataset was reduced to about 18 attributes and 319,795 transactions.

## Data Exploration:

We observed that there were no null values in the dataset and there were 4 attributes which were numerical, and all others were categorical. We then checked the dataset for duplicate values, and we found that there are 18,078 records which were duplicates. But after analyzing the dataset we went ahead without removing the duplicates as we assumed that there can be people with the exact same statistics. Our target variable was the column HeartDisease. In the dataset, we saw that the percentage of people having a heart disease was about 8.6%. This implied that the dataset was imbalanced and planned on dealing with the bias using some sampling techniques.

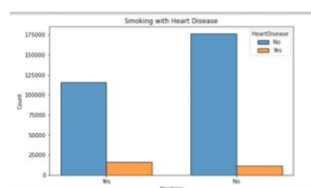


Figure 1: People having heart disease in the smoking category

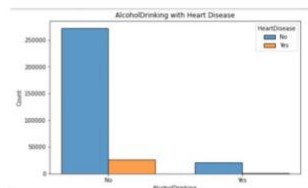


Figure 2: People having heart disease in the alcohol drinking category

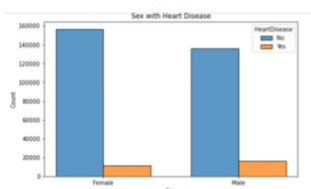


Figure 3: People having heart disease in the gender category

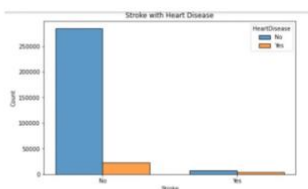


Figure 4: People having heart disease in the alcohol drinking category

We observed that, the proportion of people having heart disease is proportionally more in people who smoke, consume alcohol, and have had a stroke in the past. Also that the number of heart attacks in females was much more compared to males.

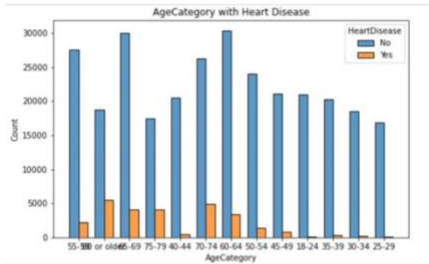


Figure 5: People having heart disease in different age categories

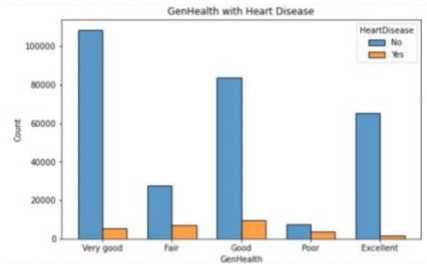


Figure 6: People having heart disease in the different general health categories

We also observed that at the age group of 60 and older and 70-74 had the highest proportion of heart diseases

and heart diseases had occurred proportionally more in people having poor general health.

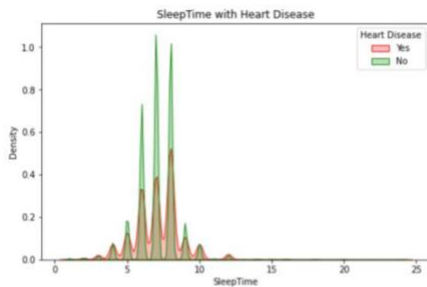


Figure 7: SleepTime density comparison for our target variable HeartDisease

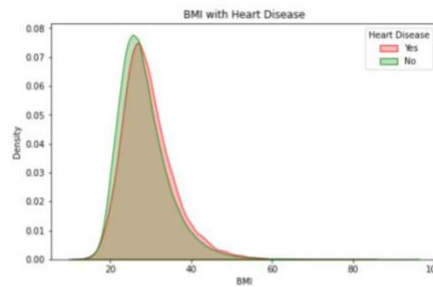


Figure 8: BMI density comparison for our target variable HeartDisease

Here we observed that people who slept less than the average number of hours of sleep had chances of developing or having a heart

disease. People having higher BMI were prone to having a heart disease.

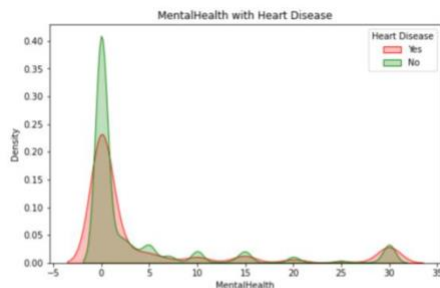


Figure 9: MentalHealth density comparison for our target variable HeartDisease

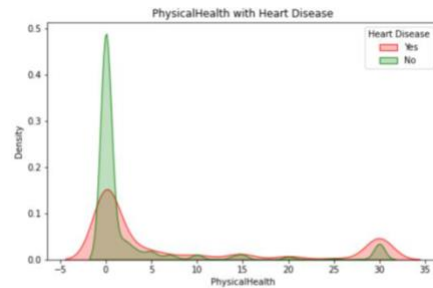


Figure 10: PhysicalHealth density comparison for our target variable HeartDisease

Here we observed that, people with physical and mental health problems showed similar characteristics. Greater the number of days a

person had these problems, more was the possibility of them developing or having a heart disease.

		Heart Disease		Percentage
		No	Yes	
Race	American Indian/Alaskan Native	4660	542	10.42%
	Asian	7802	266	3.30%
	Black	21210	1729	7.54%
	Hispanic	26003	1443	5.26%
	Other	10042	886	8.11%
	White	222705	22507	9.18%

We also observed that the greater number of heart diseases was seen in American Indian/ Alaskan Native followed by the white race in the dataset.

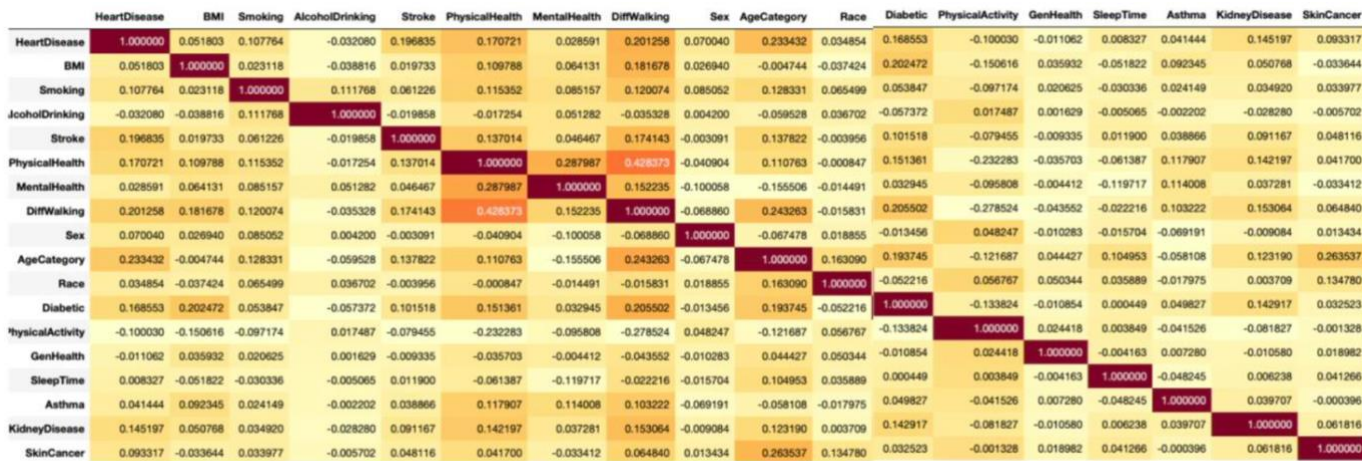


Figure 11: Correlation heat map

Here we observed that as the age of a person increased, the probability of getting the heart disease also increased. People with better physical health were less prone to getting a heart disease. People having a kidney disease, stroke and diabetes were more prone to getting a heart disease. More number of days a person had difficulty walking, more were the chances of getting a heart disease.

### Data Modelling:

#### Data Pre-Processing:

All categorical variables were encoded using Label Encoder into numerical values so that they can be in machine-readable form. We then split the dataset into train and test sets with a train size of 70 percent and a test size of 30 percent.

As mentioned earlier, we used 3 sampling techniques to overcome bias: Undersampling, Oversampling and SMOTE Oversampling.

- For Undersampling, we reduced the size of our majority class such as it is twice the size of our minority class.
- For Oversampling, we increased the size of our minority class to match the size of our majority class.
- For SMOTE oversampling, just like oversampling, we increased the size of our minority class to match the size of our majority class.

#### Models used:

We used 6 algorithms to build our models: K-Nearest Neighbors, Decision Tree, Naïve Bayes, Random Forest, XGBoost Classifier and Logistic Regression

We then used 3 evaluation metrics namely accuracy, precision and recall to compare these models to find our best fit model.

Sampling	Accuracy	Precision	Recall
Unsampled	0.90	0.33	0.88
Undersampled	0.82	0.28	0.66
Oversampled	0.88	0.42	0.99
Oversampled (SMOTE)	0.83	0.33	0.97

Table 2: K-Nearest Neighbors performance metrics

Sampling	Accuracy	Precision	Recall
Unsampled	0.86	0.22	0.25
Undersampled	0.82	0.33	0.99
Oversampled	0.99	0.93	0.99
Oversampled (SMOTE)	0.99	0.98	0.96

Table 3: Decision Tree performance metrics

Sampling	Accuracy	Precision	Recall
Unsampled	0.84	0.27	0.46
Undersampled	0.81	0.24	0.54
Oversampled	0.79	0.23	0.59
Oversampled (SMOTE)	0.71	0.19	0.71

Table 4: Naive Bayes performance metrics

Sampling	Accuracy	Precision	Recall
Unsampled	0.91	0.50	0.09
Undersampled	0.84	0.29	0.54
Oversampled	0.74	0.21	0.76
Oversampled (SMOTE)	0.70	0.18	0.71

Table 5: Logistic Regression performance metrics

Sampling	Accuracy	Precision	Recall
Unsampled	0.90	0.36	0.11
Undersampled	0.87	0.40	0.99
Oversampled	0.99	0.93	0.99
Oversampled (SMOTE)	0.99	0.96	0.98

Table 6: Random Forest performance metrics

Sampling	Accuracy	Precision	Recall
Unsampled	0.91	0.52	0.09
Undersampled	0.84	0.30	0.66
Oversampled	0.75	0.23	0.83
Oversampled (SMOTE)	0.84	0.26	0.46

Table 7: XGBoost performance metrics

We did not entirely rely on accuracy as there was bias in the dataset. We considered Recall as our main evaluation metric. Comparing all models' performance, we deduced that Tree based models performed better for this data. Our best models out of all the models we ran were Decision Tree and Random Forest.

## Insights:

- Exploratory analysis gave us some of the attributes which were used as factors.
- Correlation heat map helped us find the relation between the attributes affecting the target variable.
- To overcome bias, we had used sampling techniques.
- Random Forest and Decision Tree gave us the best results.

- After Hyperparameter tuning our best model – Decision Tree for all types of sampled data: undersampled, oversampled and SMOTE oversampled data, the following 2 decision models were the best:
  1. SMOTE Oversampled Decision Tree Model with a recall of 0.9461
  2. Oversampled Decision Tree Model with a recall of 0.9272

### Specific contribution to the project:

- Cleaning data and recognizing useful features for our analysis.
- Sampling
- Machine learning modelling in the project.

### Learning outcomes from the project:

- Better understanding of sampling techniques and how to implement them.
- Use of various models in our project and finalizing our best model with the use of evaluation metrics.
- Relevant hyper-parameter tuning of our best fit model to get more accurate results.