

Under the supervision of:

Samsung Innovation Campus

Presented By:



Wessam Salah Walid

<https://www.linkedin.com/in/wessamwalid/>



Yousr Ashraf Hegy

<https://www.linkedin.com/in/yousrhejy>



Yasmin Salah Ahmed

<https://www.linkedin.com/in/yasmin-salah-49967a21a>

HEART DISEASE PREDICTION

Supervised By:

Instructor: Dr. Doaa Mohamed

Facilitator: Eng. Shaimaa



AGENDA



1

Problem Statement

Address the problem we need to find a solution for

2

Key Findings

Demonstrate findings discovered upon analysis

3

Preprocessing

Describe our process in making our data suitable for conducting analysis

4

Models & Evaluation

Models tested and comparison

1

PROBLEM STATEMENT & DESCRIPTION

What's our problem and what do we need to achieve?





PROBLEM STATEMENT:

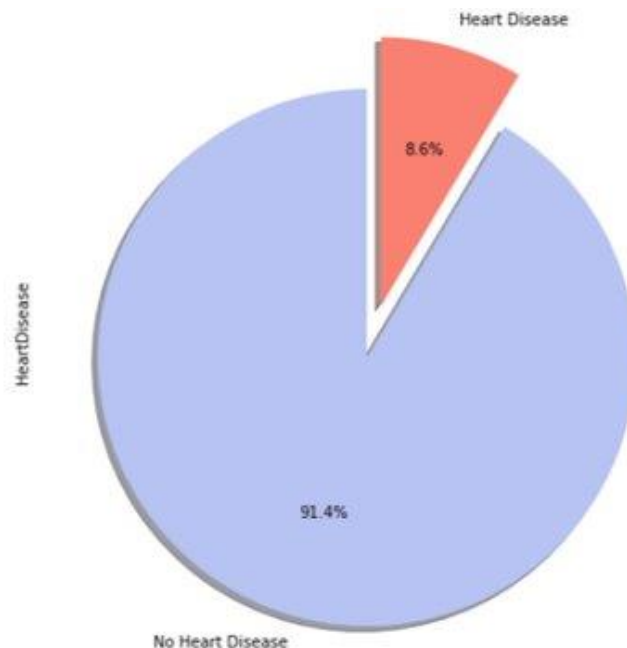
Cardiovascular diseases are one of the main causes of global death, taking an estimated 17.9 million lives each year, which represents about 31% of global deaths. Nowadays more people tend to depend on convenience food due to our packed schedules. Meanwhile, it has been reported that heart diseases are in constant increase the past few years and many are unaware of this.



DATA DESCRIPTION



- ❖ **Our Data Comes from** 2020 annual CDC survey data of 400k adults related to their health status, which conducts annual telephone surveys to gather data on the health status of U.S. residents.
- ❖ **The set consists** of 319795 rows and 18 columns.
- ❖ **The data set is** highly imbalanced and doesn't contain any missing values as we have a large ratio for people without HD than others and this indicates imbalanced data.
- ❖ **Our Target is** to predict if the patient has Heart Disease or not.
- ❖ **About** 27373 of Adults have Heart disease and 292422 doesn't have Heart Disease.



DATA DESCRIPTION

- 1** Heart Disease:
| Have you ever had a heart attack?
| (Yes / No)
- 2** BMI: Body Mass Index
|
- 3** Smoking :
| Have you ever smoked? (Yes / No)
|
- 4** Alcohol Drinking :
| Have you ever drank alcohol (Yes / No)



- 5** Stroke :
| Have you ever had a stroke? (Yes / No)
- 6** Physical Health :
| How many days during the past 30 days was
| your physical health not good? (0-30 days)
- 7** Mental Health :
| How many days during the past 30 days
| was mental health not good? (0-30 days)
- 8** DiffWalking:
| Difficulty walking or climbing stairs
| (Yes / No)
- 9** Sex :
| Male or Female

DATA DESCRIPTION

10

Age Category :
Thirteen-level age category

11

Physical Activity :
Doing physical activity or exercise during the past 30 days other than their regular job (Yes / No)

12

Race :
Ethnicity

13

Diabetic :
Ever told you had diabetes? (Yes / No)

14

GenHealth :
General Health

15

SleepTime:
Hours of sleeping in 24-hour period

16

Asthma :
Ever told you had asthma? (Yes / No)

17

Kidney Disease :
Ever told you had kidney disease? (Yes / No)

18

Skin Cancer :
Ever told you had skin cancer? (Yes / No)

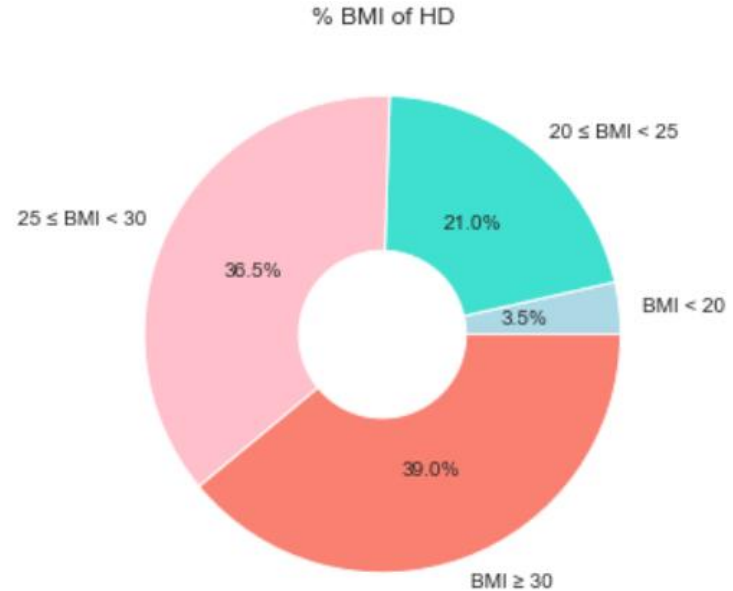
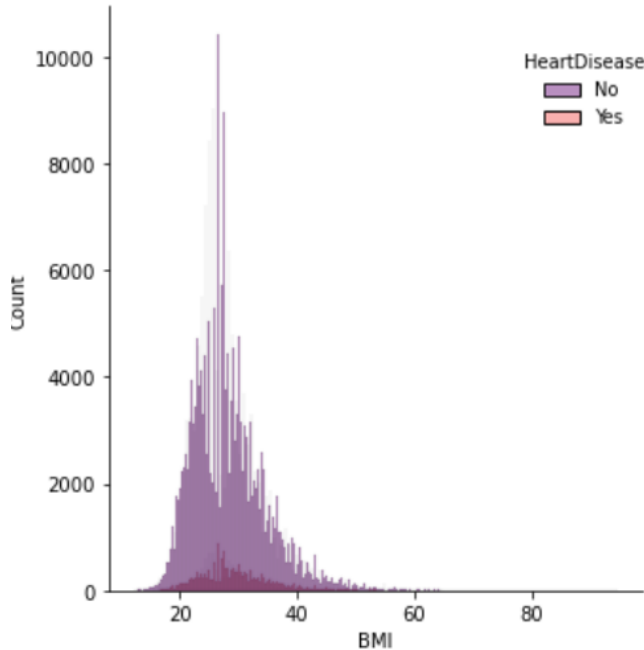
A close-up photograph of a person's chest and neck area. A white, circular medical electrode is attached to the skin on the left side of the chest. A thin, grey cable is connected to the electrode. The person's mouth and chin are visible at the top left of the frame.

Key Findings

**A PICTURE ALWAYS
REINFORCES
THE CONCEPT**



Heart Disease and BMI

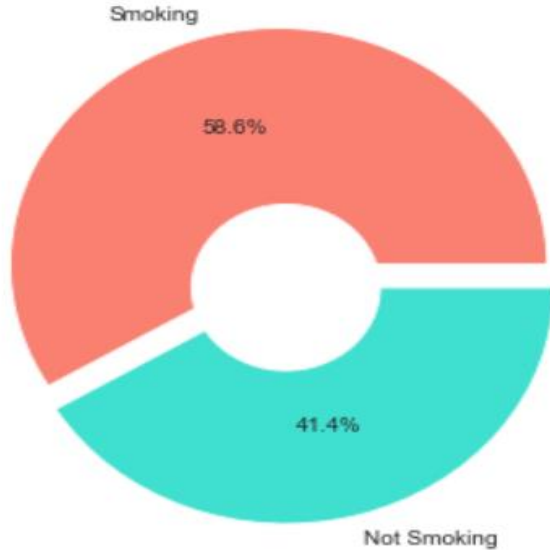


As the body mass index increases, the risk of getting the disease increases which is also stated by British Heart Foundation.

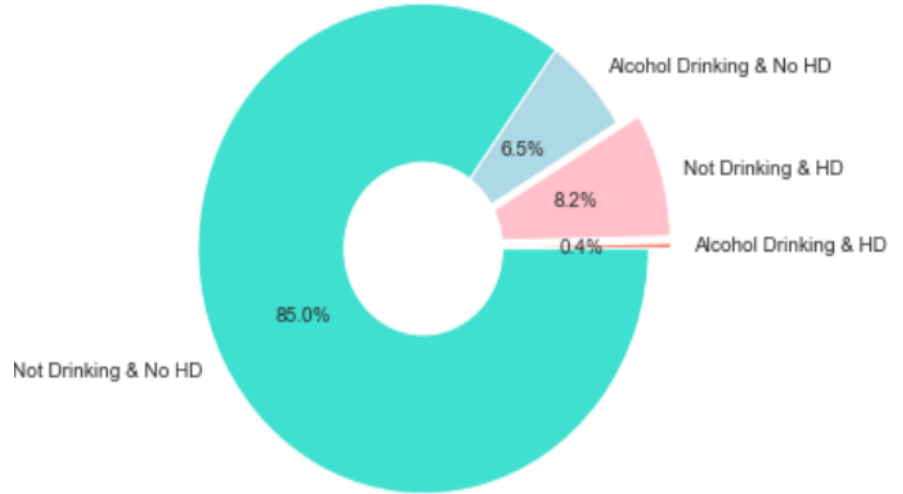
(<https://www.bhf.org.uk/informationsupport/riskfactors/obesity#:~:text=How%20does%20obesity%20increase%20the%20risk%20of%20heart%20and%20circulatory%20diseases%3F&text=Excess%20weight%20can%20lead%20to,lead%20to%20a%20heart%20attack.>)

Smoking and Alcohol Drinking

% Smokers Among Heart Disease



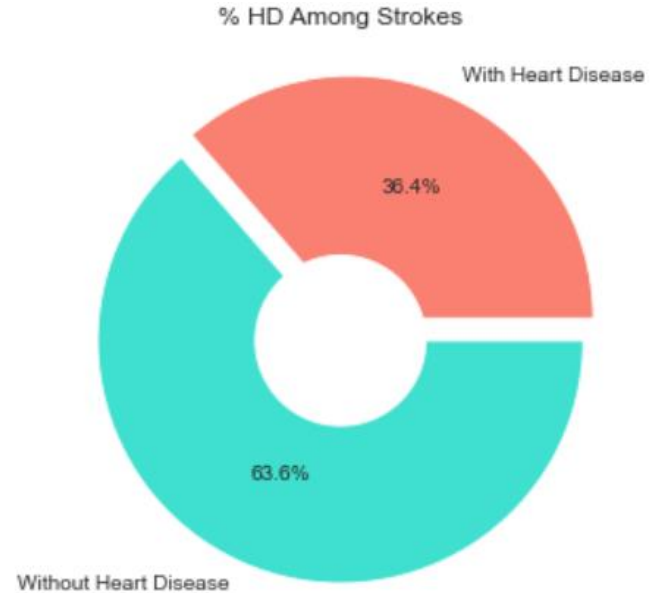
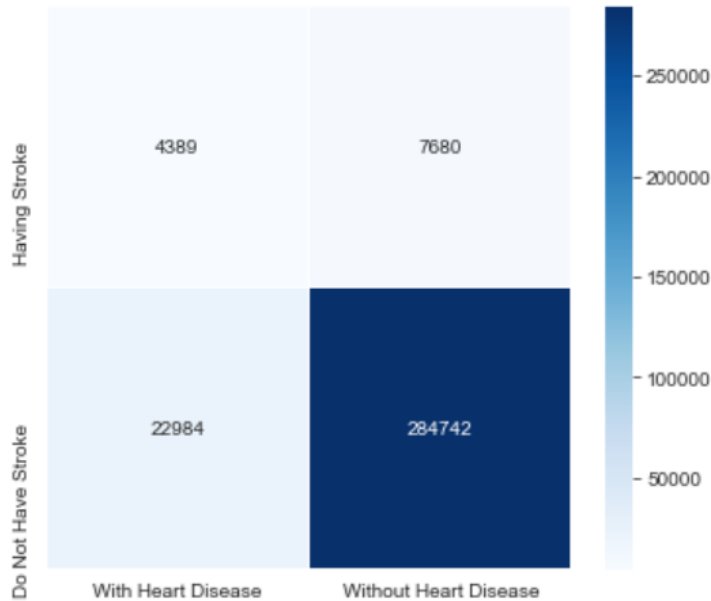
% Alcohol Drinking & HD



People who smoke are most vulnerable to toxic effects on health and cardiovascular systems which is proved by FDA. That's besides BHF that stated the dangerous effect of drinking alcohol on patients.

(<https://www.fda.gov/tobacco-products/health-effects-tobacco-use/how-smoking-affects-heart-health>)(<https://www.bhf.org.uk/informationsupport/heart-matters-magazine/medical/effects-of-alcohol-on-your-heart>)

Strokes and Heart disease



Here are the people to have HD without strokes are small, and that is because the data isn't balanced.

According to CDC, people with HD are also vulnerable to having strokes because it's an important risk factor that affects cardiovascular diseases due to unhealthy lifestyles and physical inactivity.

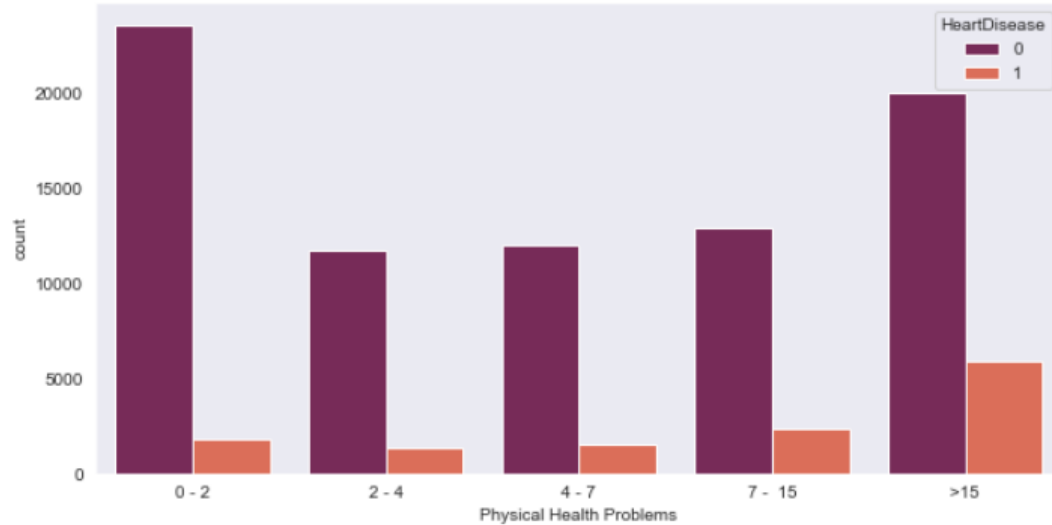
(<https://www.cdc.gov/chronicdisease/resources/publications/factsheets/heart-disease-stroke.htm#:~:text=Making%20blood%20sticky%20and%20more,and%20narrowing%20of%20blood%20vessels>)

Physical Health Problems:



Symptoms of Heart disease include:

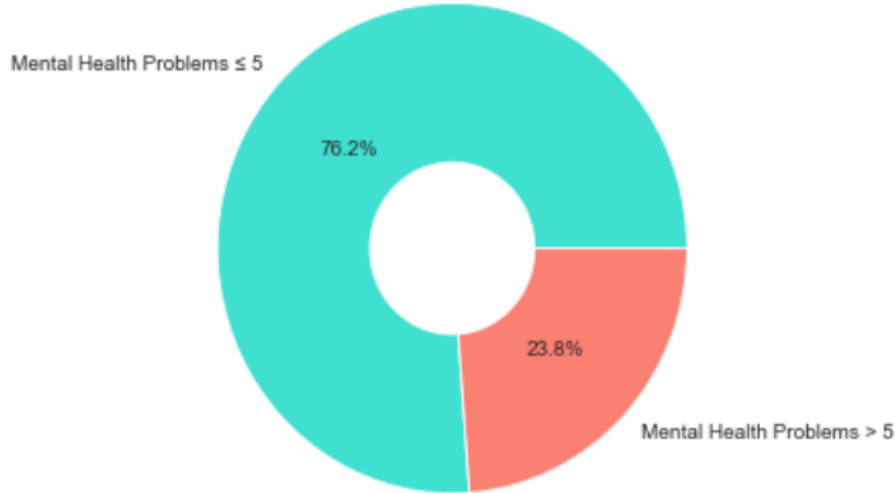
- Chest pain, chest tightness, and Angina
- Shortness of breath
- Pain and coldness in the legs or arms.



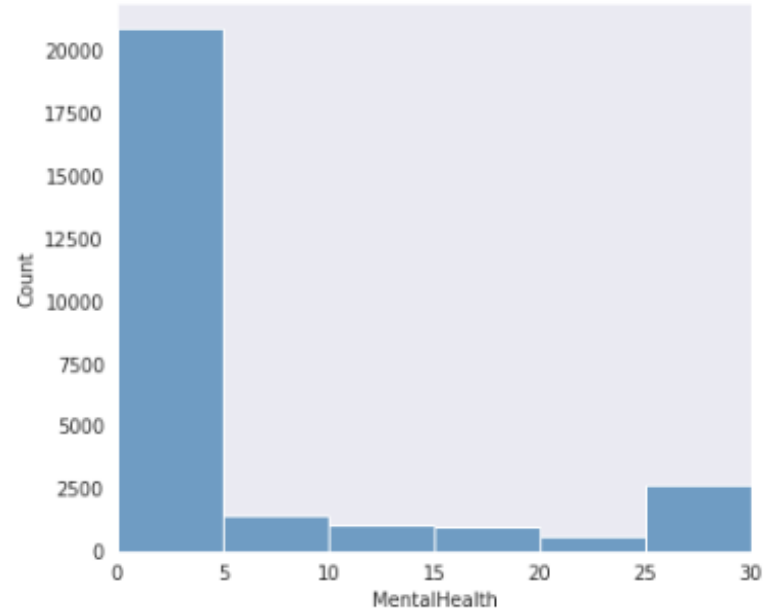
Mental Health Problems and HD:



% Mental Health Problems of HD



Mental Health Problems of People with HD



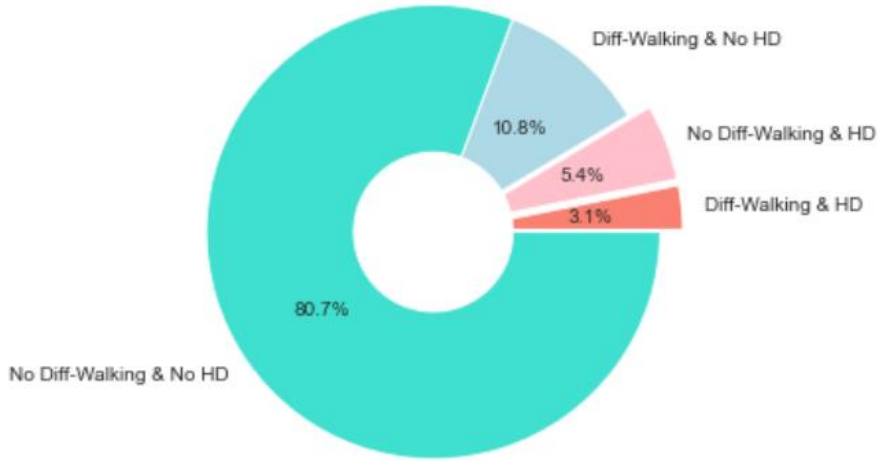
According to CDC, there is a strong relation between mental health problems and heart disease.

(<https://www.cdc.gov/heartdisease/mentalhealth.htm>)

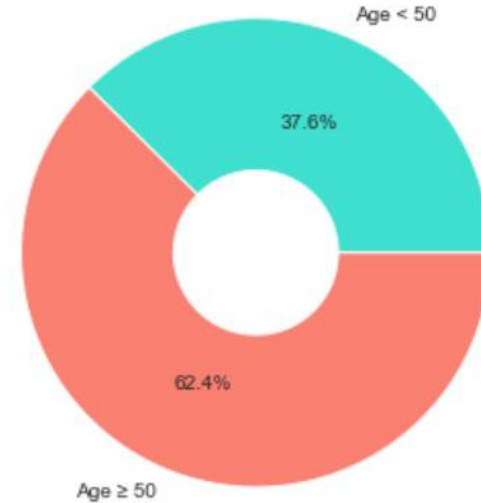
Difficulty walking and Age Category



% Diff-Walking & HD



% Age of Less Than 50 and the Complement of People with HD

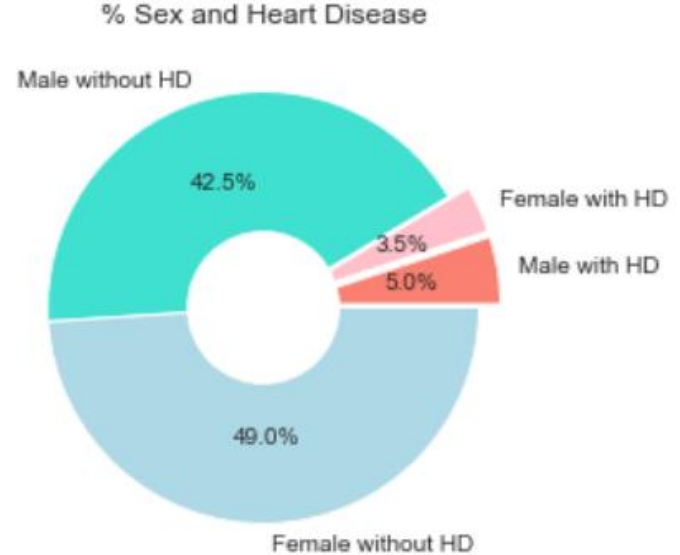


It is clear that old people are exposed to such a disease easily more than young people and sometimes they are even according to CDC (https://www.cdc.gov/heartdisease/any_age.htm)

Gender and Heart Disease:

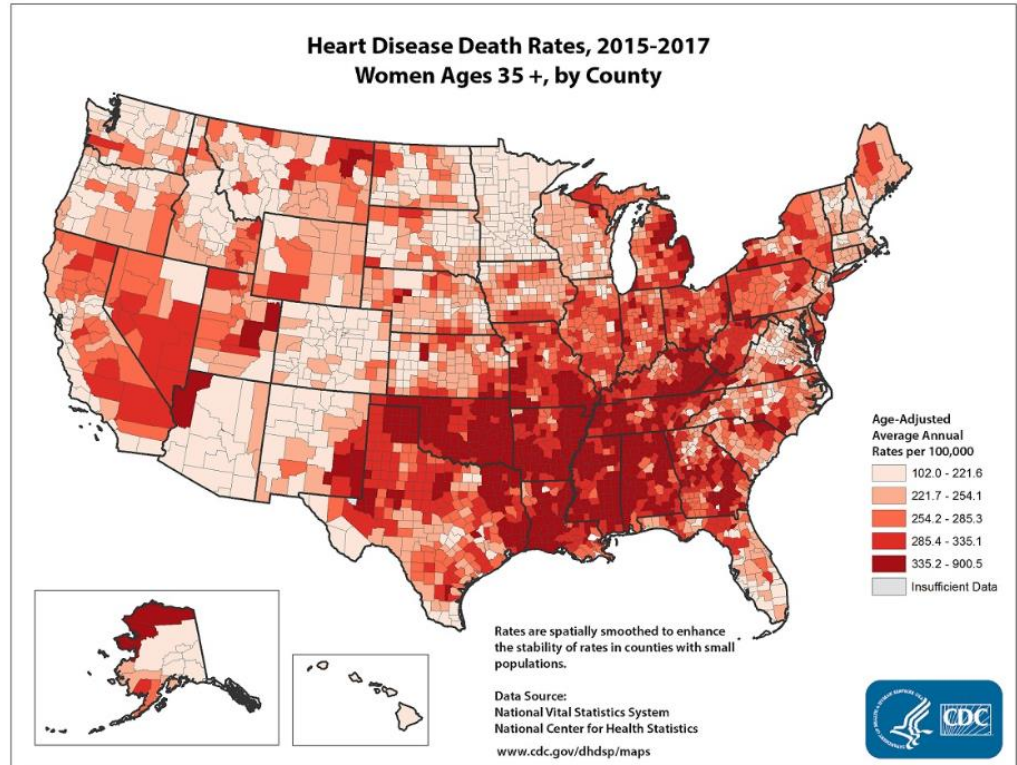


The CDC stated that despite the increases in awareness over the past decades, only about half 56% of women recognize that heart disease is their number 1 killer as HD develops 7 to 10 years later in women than men. In addition, women were more likely than men to be older and have a more complicated medical history at the time of their heart attacks.



Gender and Heart Disease:

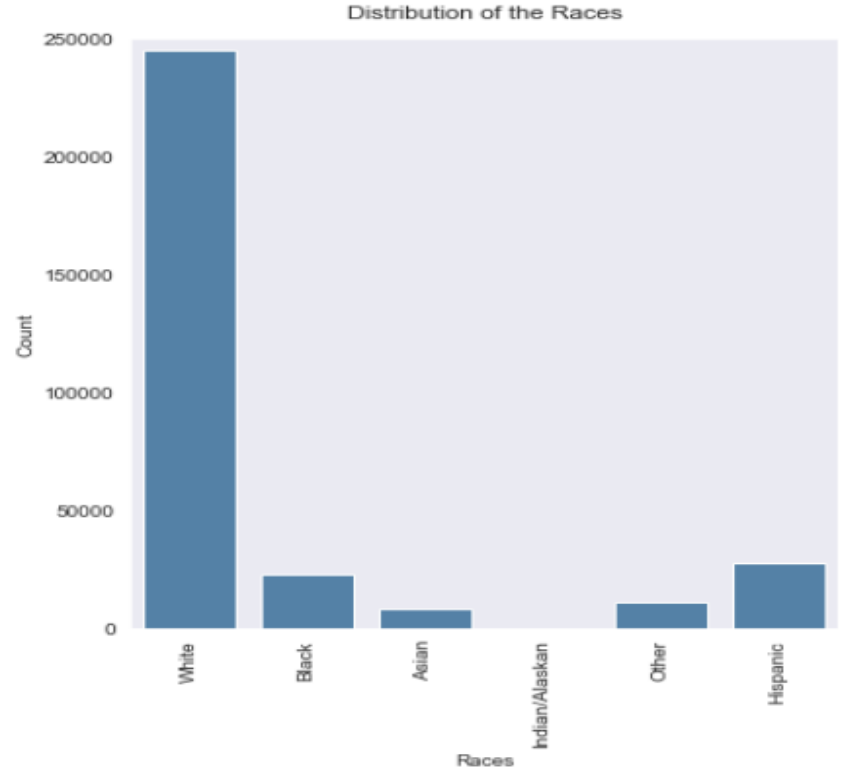
According to CDC, This map shows death rates from heart disease in women in the United States. The darker red indicates a higher death rate. (<https://www.cdc.gov/heartdisease/women.htm>)



Race and Heart Disease:



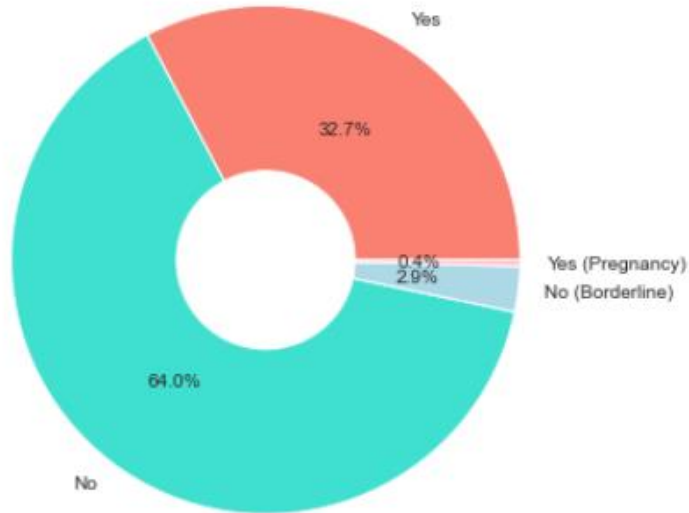
As the Cleveland Clinic website, black people are more vulnerable to heart disease than white people according to some diagnoses, due to the health disparities of race and ethnicity. ([Heart Disease Risk: How Race and Ethnicity Play a Role](https://clevelandclinic.org) (clevelandclinic.org))



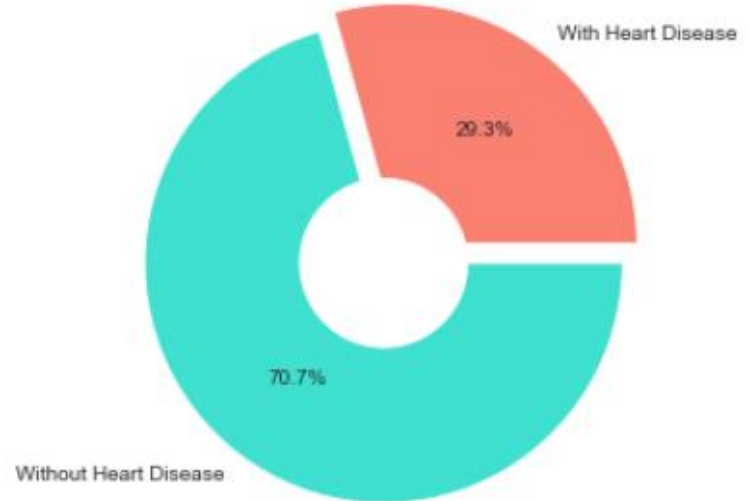
Diabetes and Kidney Disease conditions



% Diabetic Conditions among People with HD



% HD Among People with Kidney Disease

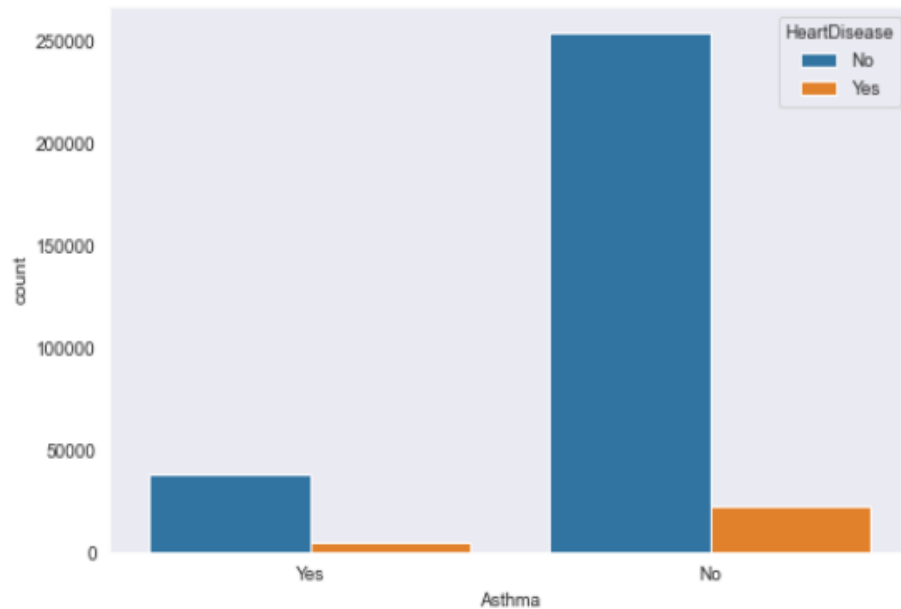


According to CDC, People with diabetes are more likely to have other conditions that raise the risk for heart disease (<https://www.cdc.gov/diabetes/library/features/diabetes-and-heart.html#:~:text=Over%20time%2C%20high%20blood%20sugar,and%20can%20damage%20artery%20walls>)

Asthma Heart Disease:



Previous studies have associated asthma with an increased risk of cardiovascular conditions, including Heart attack and cardiovascular disease.

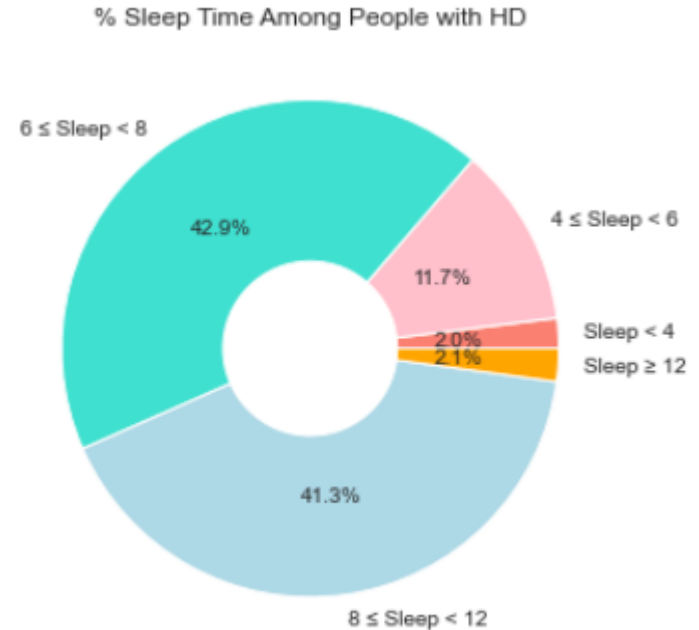


Sleep time and Heart Disease:

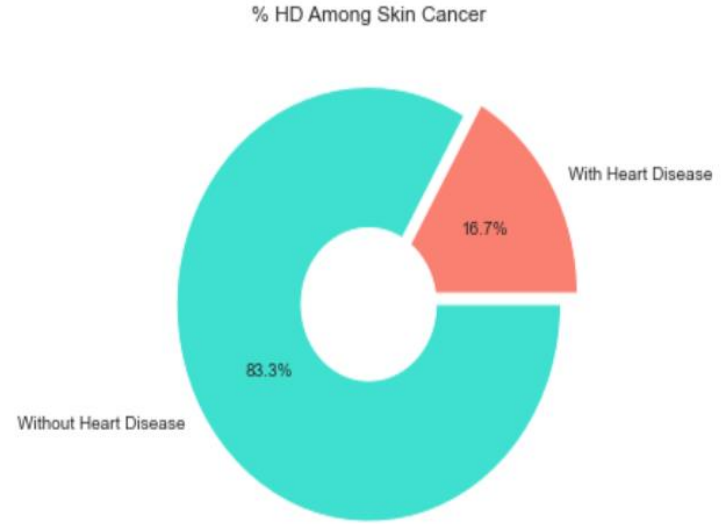
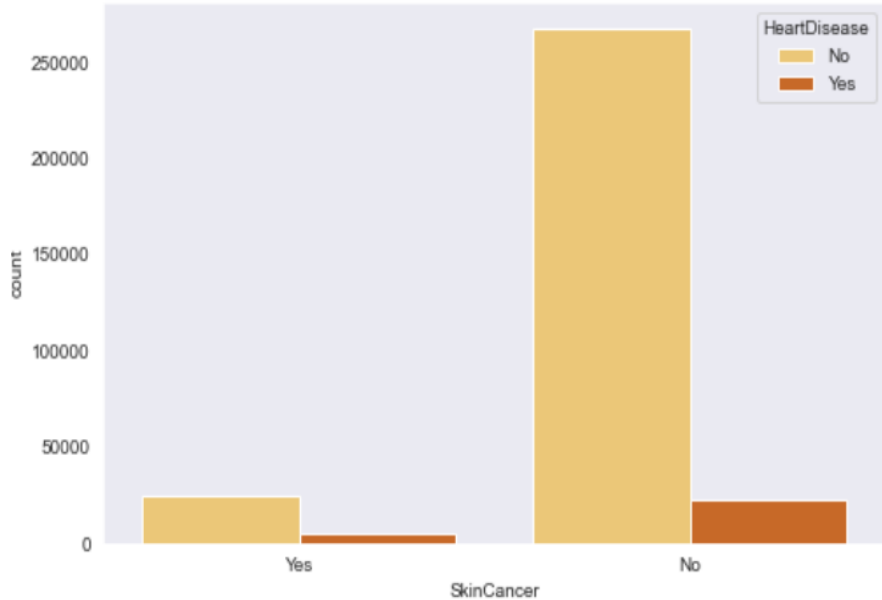


Going to sleep between 10:00 and 11:00 pm is associated with a lower risk of developing heart and circulatory disease compared to earlier or later bedtimes. According to CDC, Adults who sleep less than 7 hours each night are more likely have health problems, including heart attack, asthma, and depression.

([How Does Sleep Affect Your Heart Health? | cdc.gov](https://www.cdc.gov/heartdisease/sleep/index.html))([Too Much Sleep May Bring Heart Disease, Death Risk \(webmd.com\)](https://www.webmd.com/heart-disease/heart-disease/death-risk-too-much-sleep))



Skin Cancer and Heart Disease



It is clear from the graphs above that there is a very weak relationship between skin cancer and heart disease.

3

Data Preprocessing

What's should be done for data to fit into the models?



Preprocessing stages:

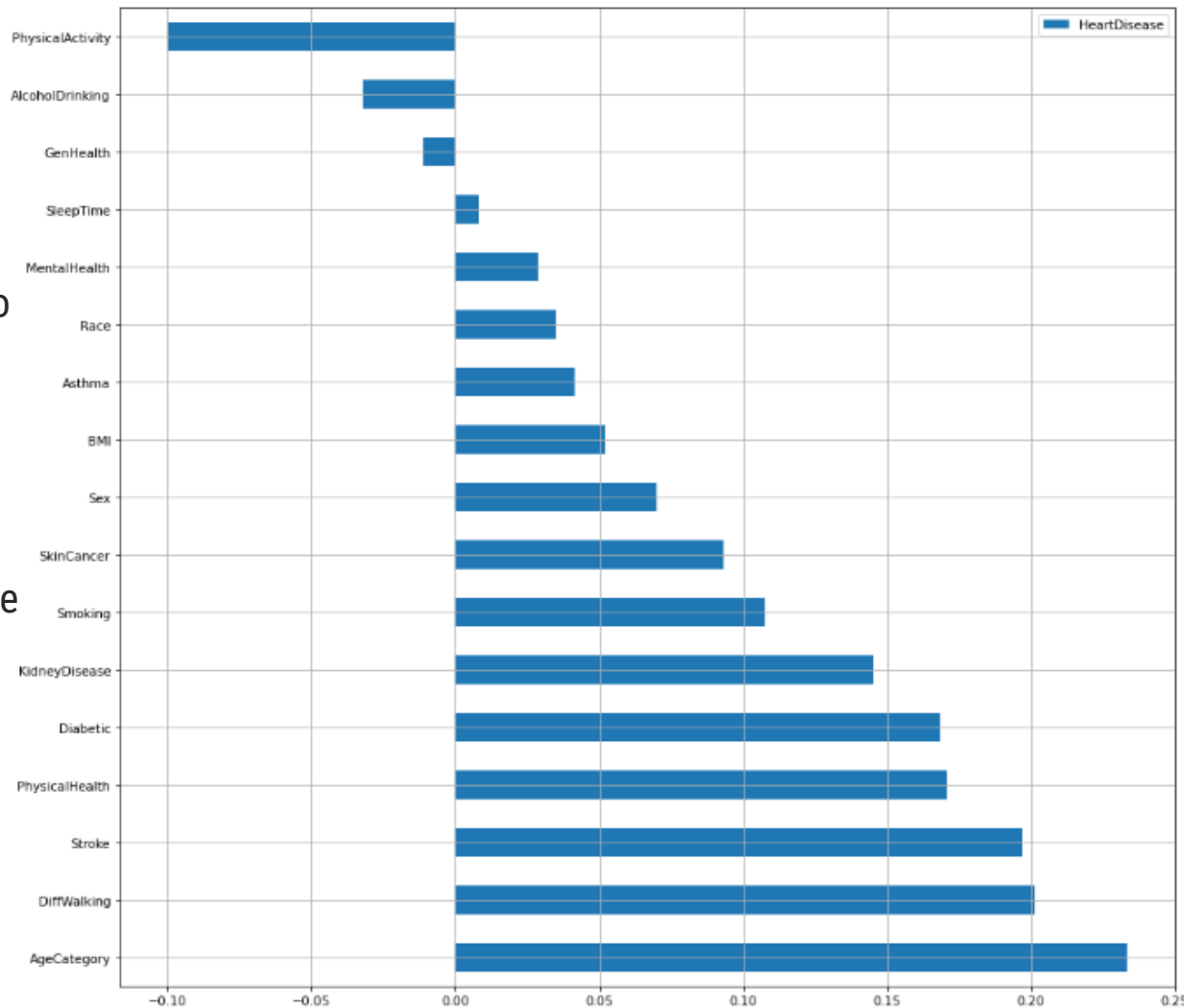
- ❖ Checking null values.

- ❖ Encoding categorical data.

Which var + visualization -> which algo

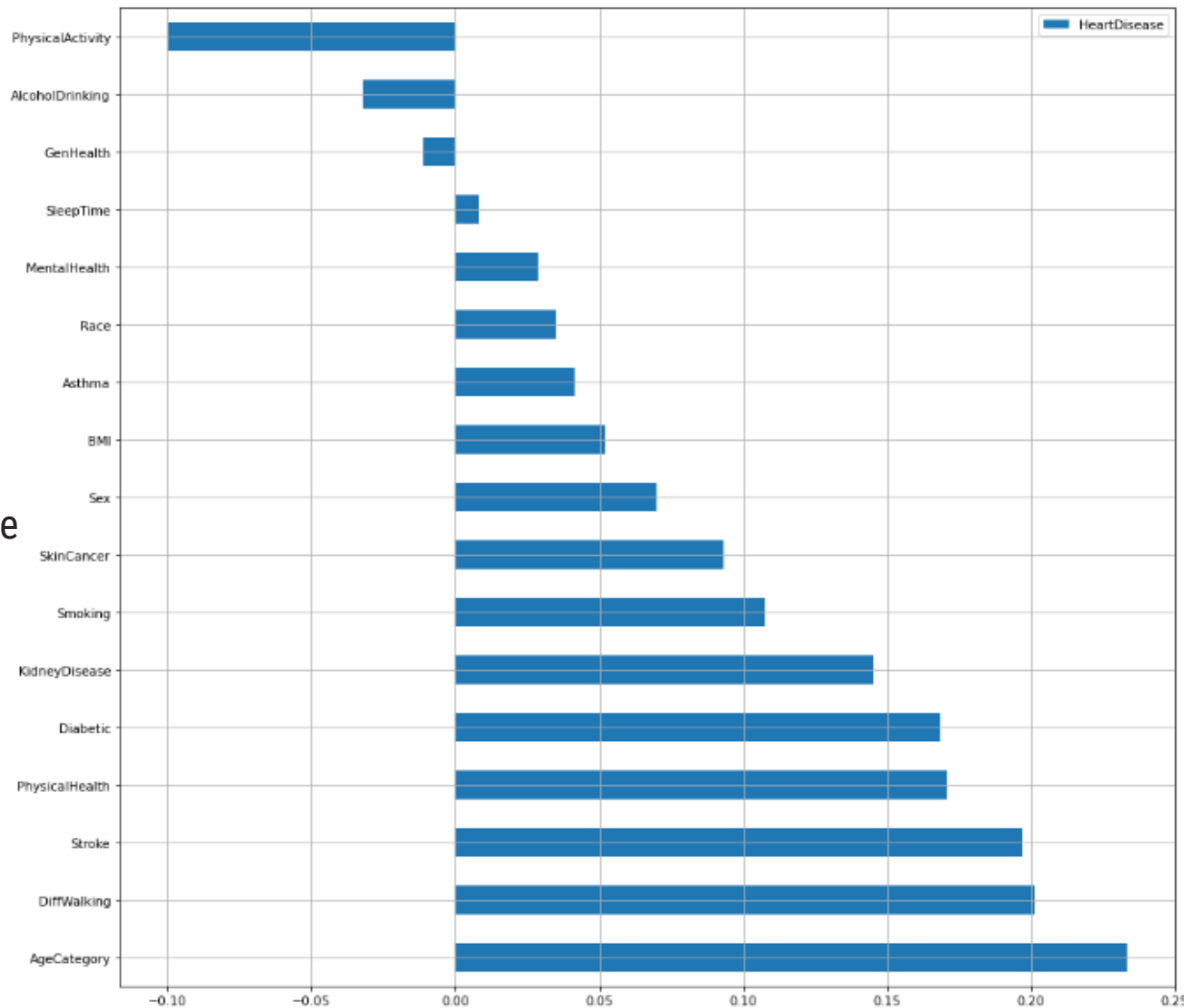
- ❖ Presenting correlation between features. We found some strong relationships like DiffWalking & PhysicalHealth with 0.43. There are also weak relationships like heart disease & Stroke with 0.2.

- ❖ Handling outliers.



Preprocessing stages:

- ❖ Checking null values.
- ❖ Encoding categorical data.
- ❖ Presenting correlation between features. We found some strong relationships like DiffWalking & PhysicalHealth with 0.43. There are also weak relationships like heart disease & Stroke with 0.2.
- ❖ Handling outliers.



Modeling and Prediction



Logistic Regression Classifier:

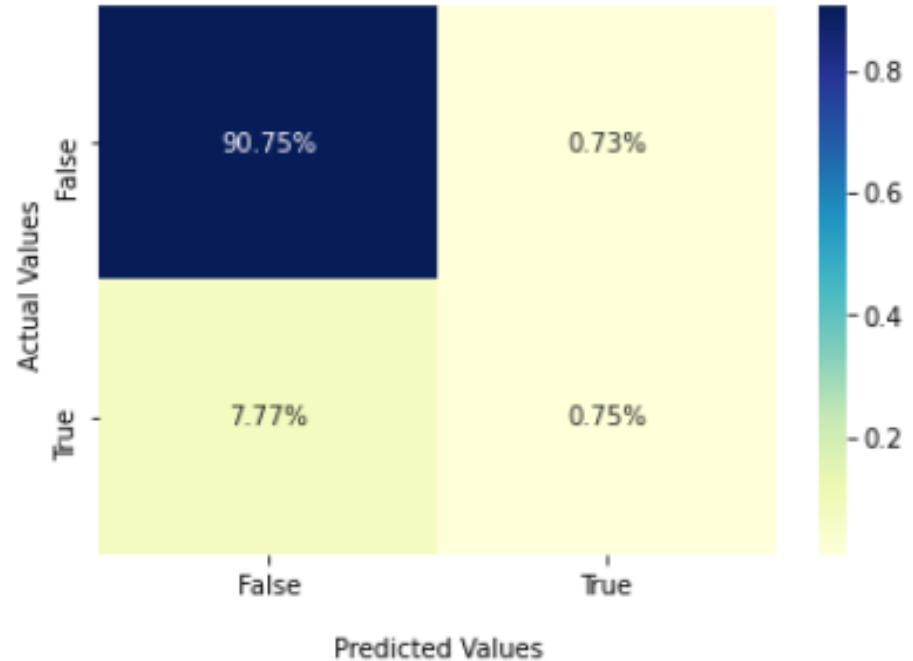
The Classification Report for LR Classifier:

	precision	recall	f1-score	support
0	0.92	0.99	0.96	58512
1	0.51	0.09	0.15	5447
accuracy			0.92	63959
macro avg	0.71	0.54	0.55	63959
weighted avg	0.89	0.92	0.89	63959

Without using Sampling:

❖ Accuracy: 91.5%

This model guess the TN more than TP



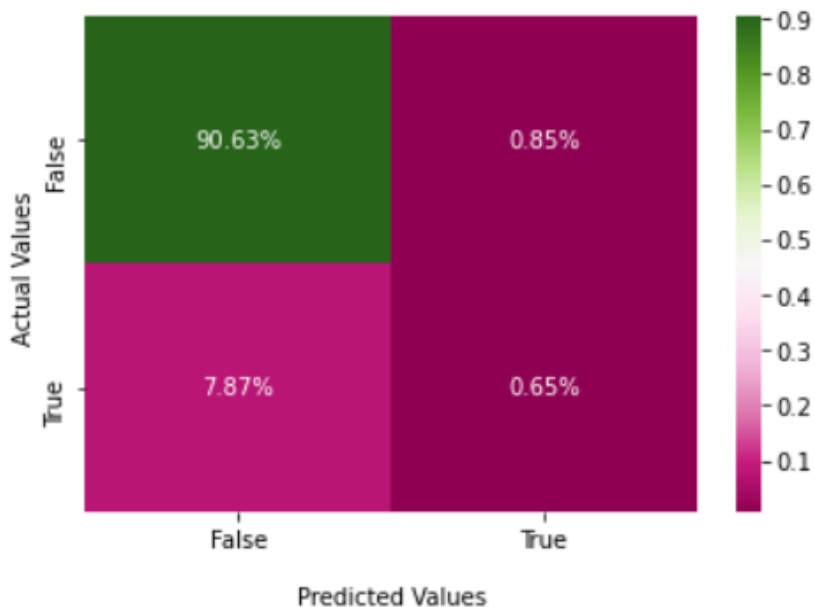
KNieghbours Classifier:



This model also doesn't need Sampling as it is not influenced in any way by the size of the class:

- ❖ Accuracy: 91.28%
- ❖ Precision: 82%
- ❖ Recall: 96%
- ❖ F1-Score: 89%

It is a good indicator for TN not for TP

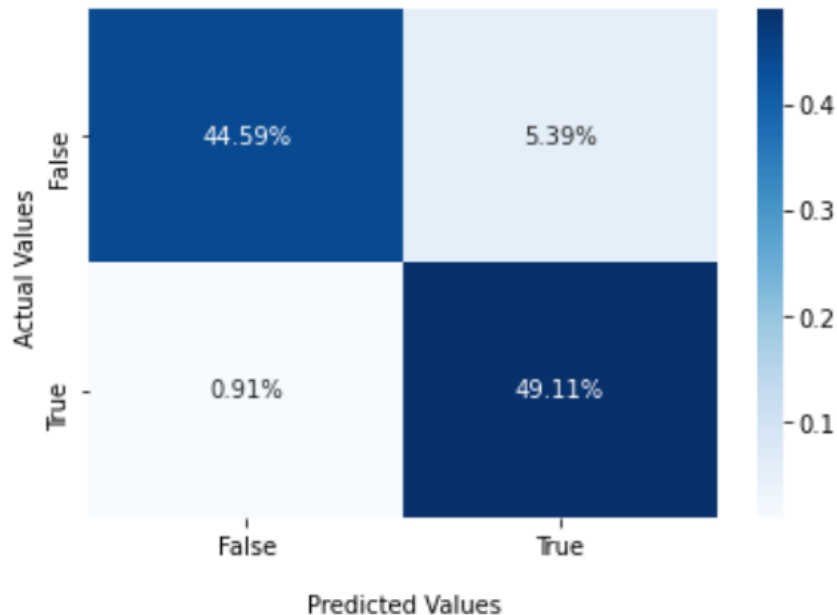


Random Forest Classifier:



The Classification Report for RF Classifier:

	precision	recall	f1-score	support
0	0.98	0.89	0.93	58461
1	0.90	0.98	0.94	58508
accuracy			0.94	116969
macro avg	0.94	0.94	0.94	116969
weighted avg	0.94	0.94	0.94	116969

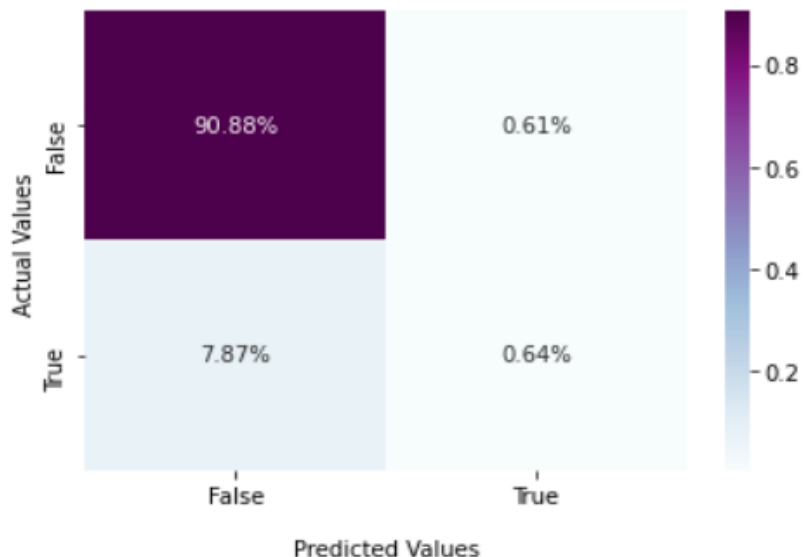


Using Random OverSampling:

❖ Accuracy: 93.72%

This model Guesses the TP more than TN

XGBoost Classifier:



This model doesn't need sampling cause it is know with the great ability to handle overfitting and imbalanced data:

- ❖ Accuracy: 91.97%
- ❖ Precision: 52%
- ❖ Recall: 8%
- ❖ F1-Score: 13%

This model Guesses the TN way more than TP

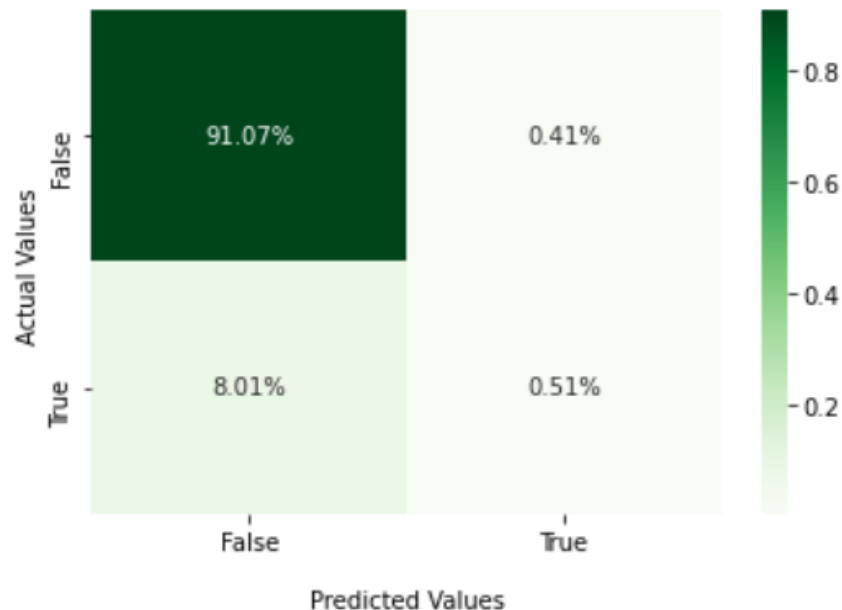
Votting Classifier:



Without using Sampling:

- ❖ Accuracy: 91.72%
- ❖ Precision: 56%
- ❖ Recall: 6%
- ❖ F1-Score: 11%

As the rest of the models, this one predicts the TN better.



The Best Model for our data is:



Random Forest Classifier:

It produces the best accuracy, which is about 94% under the condition of oversampling, besides that, it predicts the true positives and true negatives in nearly equal ranges unlike other models that only guess the true negatives with high ranges.



SAMSUNG

THANKS!

Together for Tomorrow!

Enabling People

Education for Future Generations

©2021 SAMSUNG. All rights reserved.

Samsung Electronics Corporate Citizenship Office holds the copyright of book.

This book is a literary property protected by copyright law so reprint and reproduction without permission are prohibited.

To use this book other than the curriculum of Samsung innovation Campus or to use the entire or part of this book, you must receive written consent from copyright holder.