# Heart disease detection project

**P122 - Group 4**
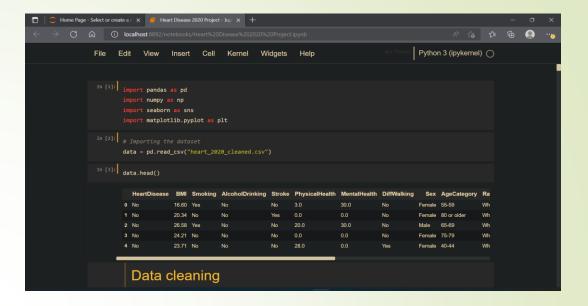
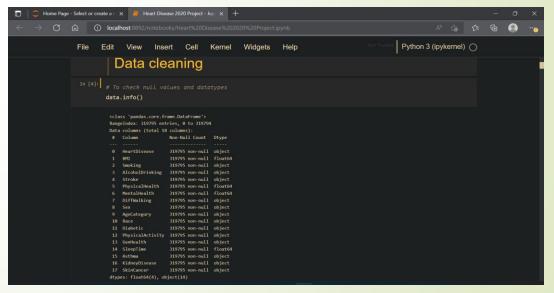Harshada More

# Names of group members:

- ➢ Ms. Harshada More
- ➢ Ms. Vaibhavi Taide
- ➢ Mr. Nagaraja M. R.
- ➢ Mr. Ganesh N. R.
- ➢ Mr. Parag Wani
- ➢ Ms. Mansi Solanki

# Importing dataset
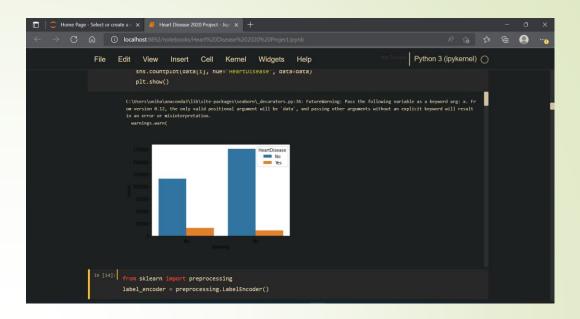
# Data cleaning

# EDA
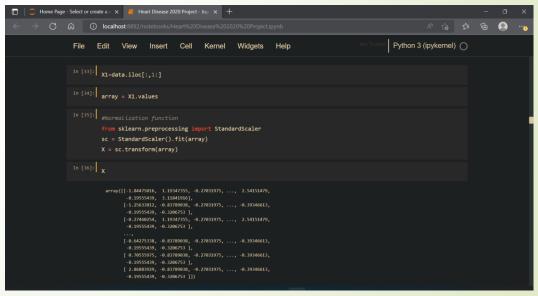




Harshada More

# Visualisation

# Label encoder

# Defining variables

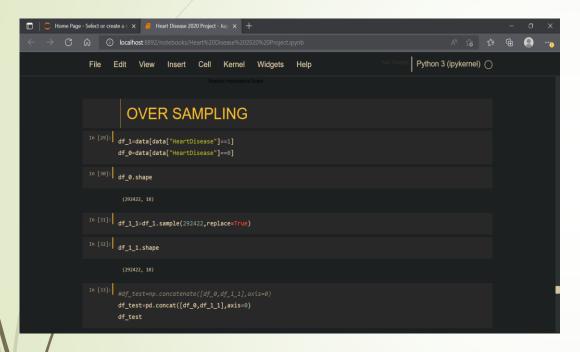# Train-Test split data





Harshada More

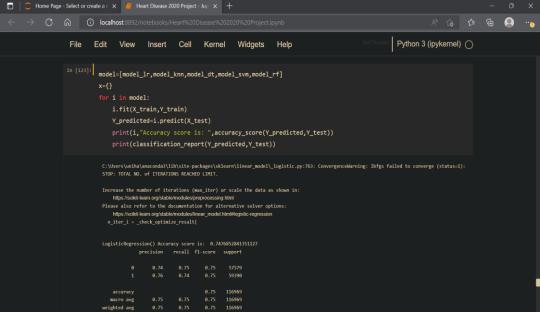# Model analysis performed for the following models:

- Logistic regression
- Decision Tree Classifier
- KNN
- Naïve Bayes
- XGBM
- LGBM
- SVM
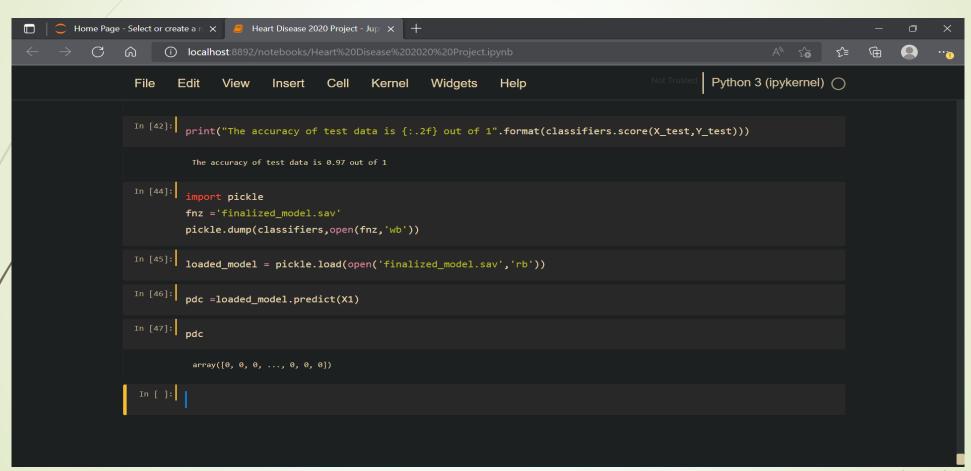- Random forest

Harshada More

# Model analysis performed:

- Confusion matrix
- Model accuracy for training data
- Model accuracy for test data
- Classification report

Harshada More

# Over sampling model analysis



Harshada More

# Saving the finalized model – Random forest by using pickle



Harshada More

# Streamlit code (Part of the code)

```python
import streamlit as st
import pandas as pd
import numpy as np
#import matplotlib as plt
#from sklearn.ensemble import RandomForestClassifier
import matplotlib.image as mp
import pickle




st.title('Hello user')

img = mp.imread("heart2.png")
st.image(img)


st.sidebar.header('User, please give your inputs for the following:')

loaded_model1 = pickle.load(open('finalized_model.sav','rb'))


def user_input_features():
    BMI = st.sidebar.number_input('Insert your BMI',0,100)
    Smoking = st.sidebar.selectbox('Do you smoke?',["Yes","No"])
    AlcoholDrinking = st.sidebar.selectbox('Do you drink alcohol',["Yes","No"])
    Stroke = st.sidebar.selectbox("Did you ever have stroke before?",["Yes","No"])
    PhysicalHealth = st.sidebar.number_input("Insert your physical health status",0,30)
```

```python
    elif i == race[2]:
        df["Race"] = 2
    elif i == race[3]:
        df["Race"] = 3
    elif i == race[4]:
        df["Race"] = 4
    else:
        df["Race"] = 5


predictions = loaded_model1.predict(df)


st.subheader('Predicted Result')




def result():
    if predictions == 0:
        results = "You do not have a heart disease."
    else:
        results = "Heart disease detested"
    return results


results = result()

st.write(results)
```

Harshada More

# Model final output



Harshada More

Harshada More