

Machine Learning project – Phase 1

Dataset - Personal Key Indicators of Heart Disease

1) **Q:** What are the business needs you are trying to address

A: My business needs I'm trying to address is to understand how adult people can improve their health status

2) **Q:** What are the objectives of the Data Science project – use supervised learning (regression or classification) and define how your target, y, addresses the business need/problem

A: my y target is "HeartDisease" and I will predict the result (0/1) based on my features and when I will know the weight of each feature I will understand how adult people can Improve the chance of not getting heart disease.

3) **Q:** What kind of data is available

A: The dataset is about " Key Indicators of Heart Disease" based on " 2020 annual CDC survey data of 400k adults related to their health status" and has columns about heart disease (the target), BMI, Smoking, Alcohol Drinking, Age Category and etc....

4) **Q:** How large is your data

A: My data has 400K rows and 18 columns

5) **Q:** Main features

A: My main features are: Smoking, Stroke, AgeCategory, PhysicalHealth, DiffWalking, Diabetic, SkinCancer and KidneyDisease

6) **Q:** Can you find a new data set online that you could merge and increase your insights

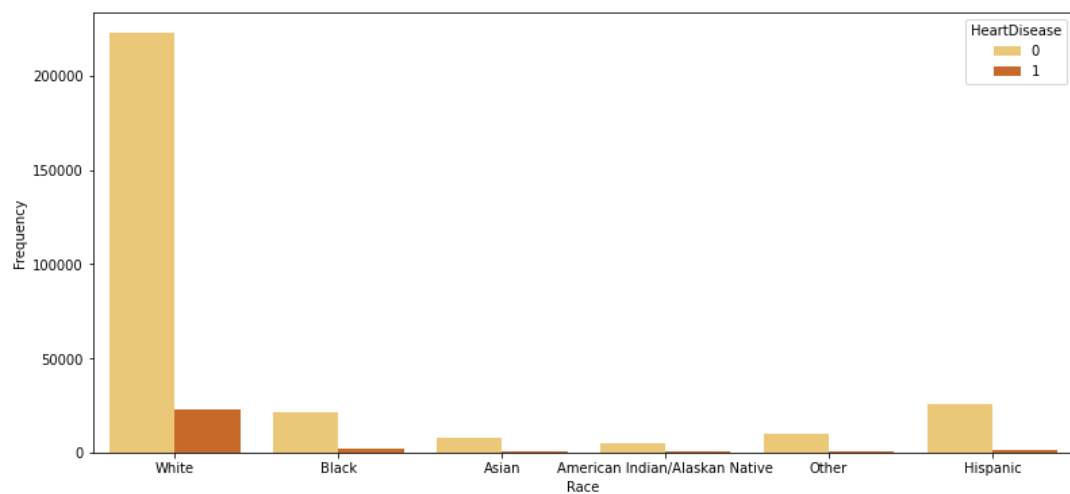
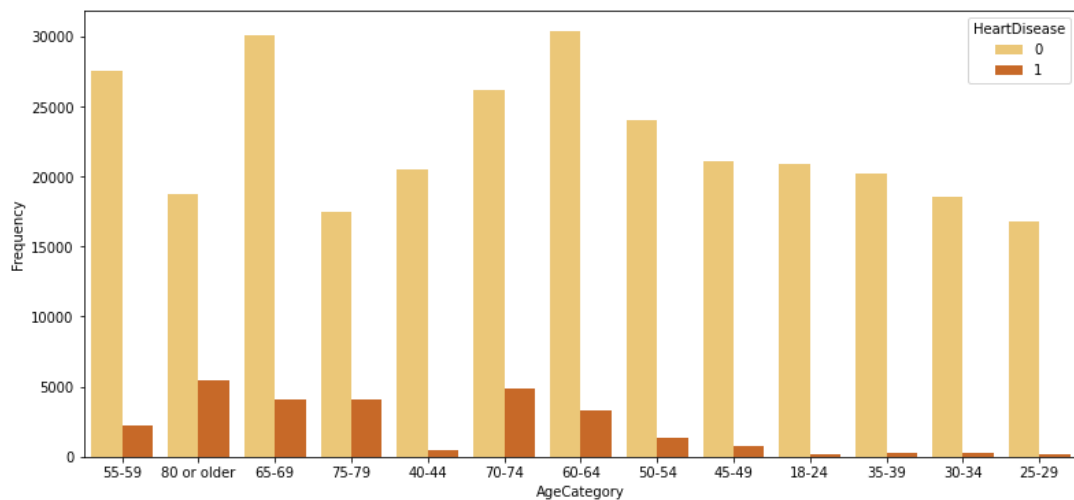
A: Yes, I took my data from Kaggle, the name of the dataset is "Personal Key Indicators of Heart Disease"

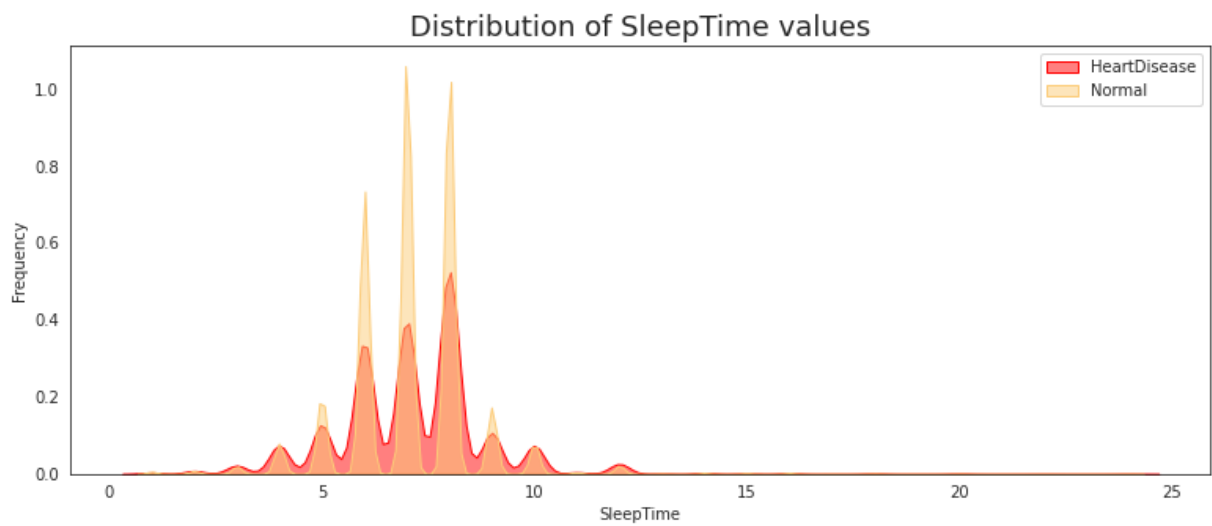
The original location came from the CDC and is a major part of the Behavioral Risk Factor Surveillance System (BRFSS), which conducts annual telephone surveys to gather data on the health status of U.S. residents.

I can go to their articles and get more data from their surveys.

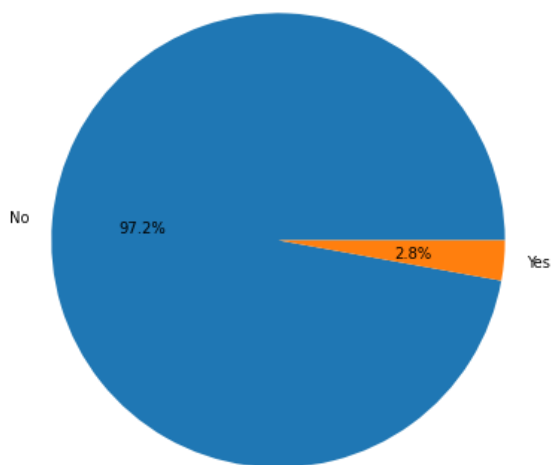
7) **Q:** Exploratory Data Analysis – 2-6 interesting visualizations. Self-explanatory

A:

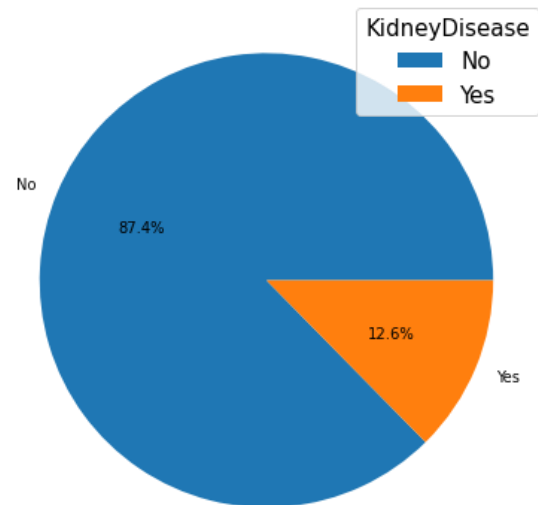




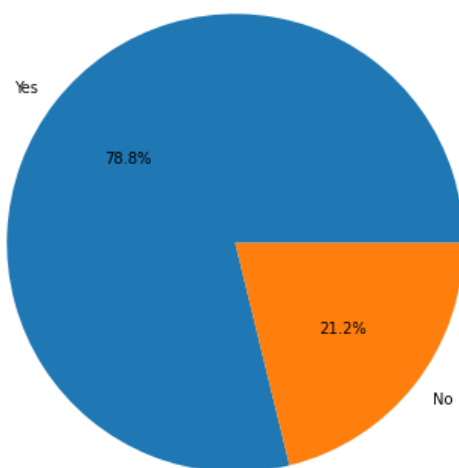
No Heart Disease



Yes Heart Disease



No Heart Disease



Yes Heart Disease

