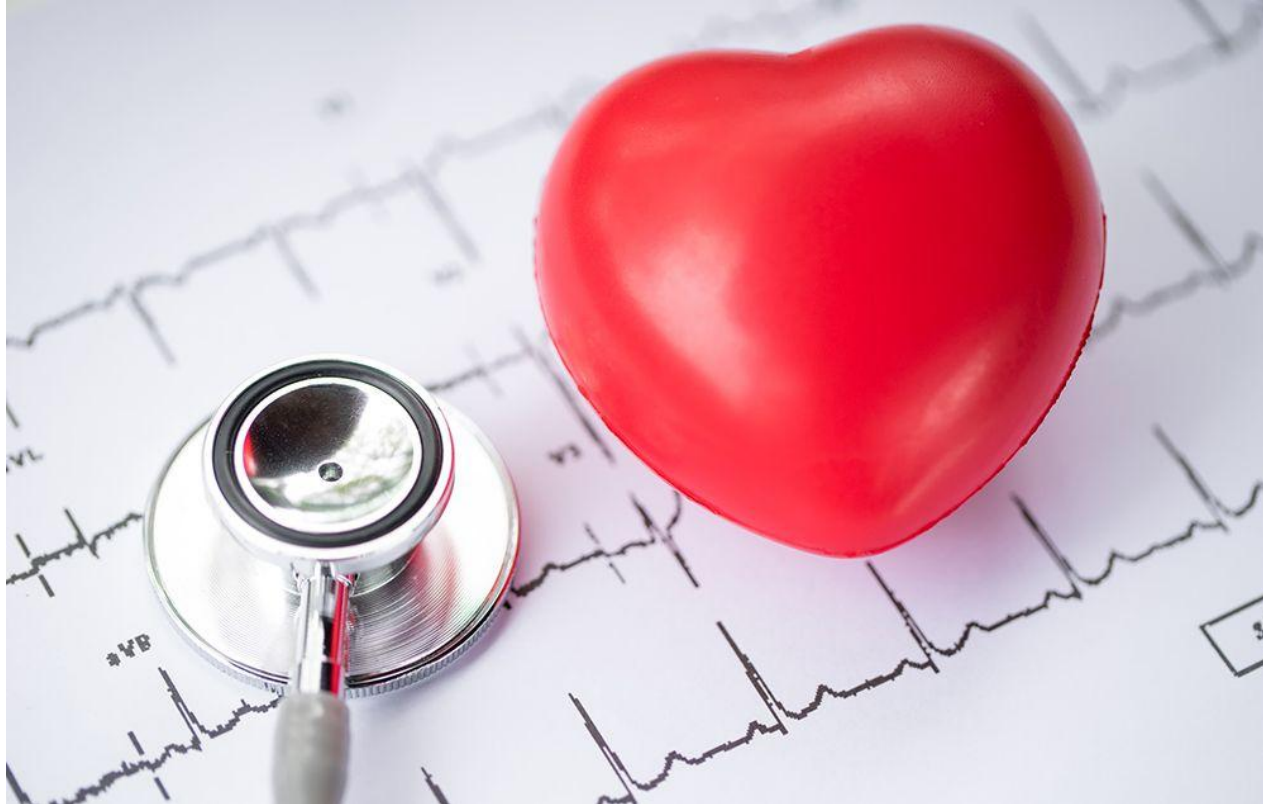


Heart Disease Classification

Group 7: Lucas Ogilvie
Mark Biernacki
Ran Ma



Background





Data Overview

1. Source: Behavioral Risk Factor Surveillance System (BRFSS) 2021

Developed by the Centers for Disease Control and Prevention (CDC)

Annual health-related telephone surveys to gather data across the U.S.

2. Size:

438,693 rows

303 columns

3. Content of the data:

Demographic detail: age, sex, race

Behavioral factors: smoking, alcohol consumption, physical activity levels

Chronic health conditions

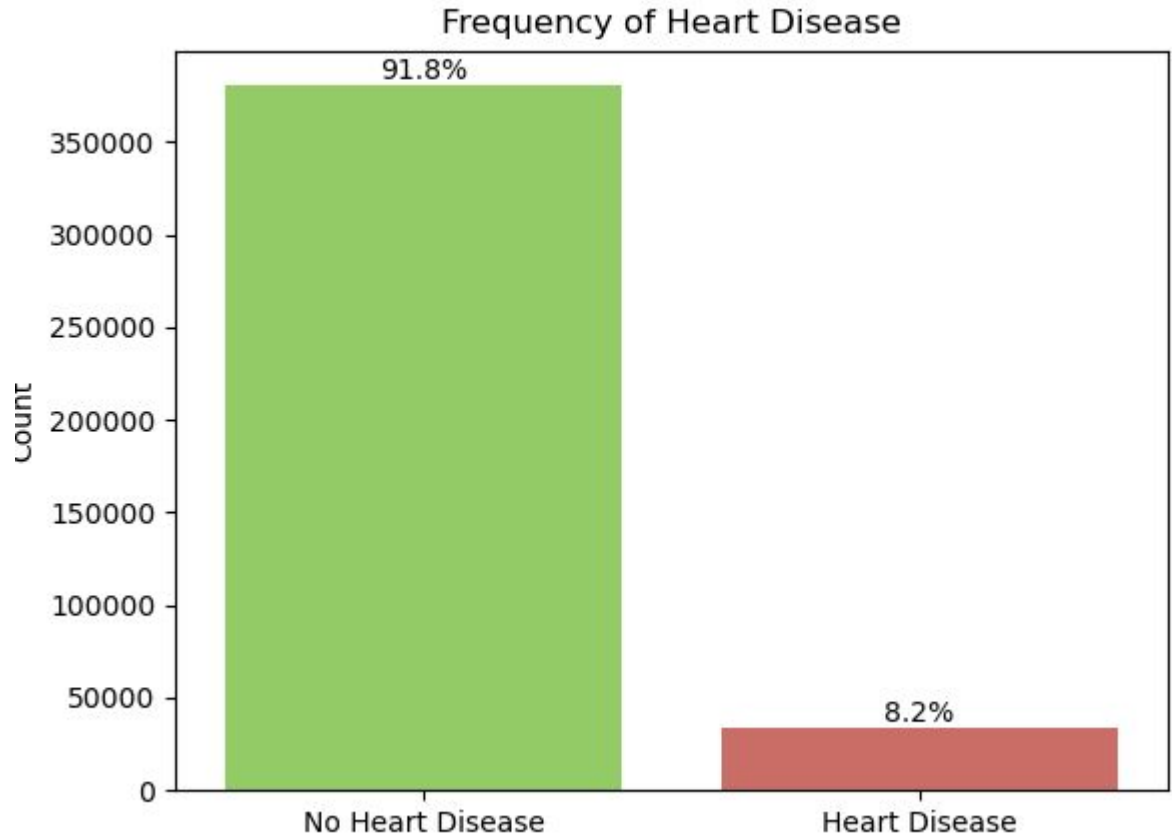
Preventive health measures



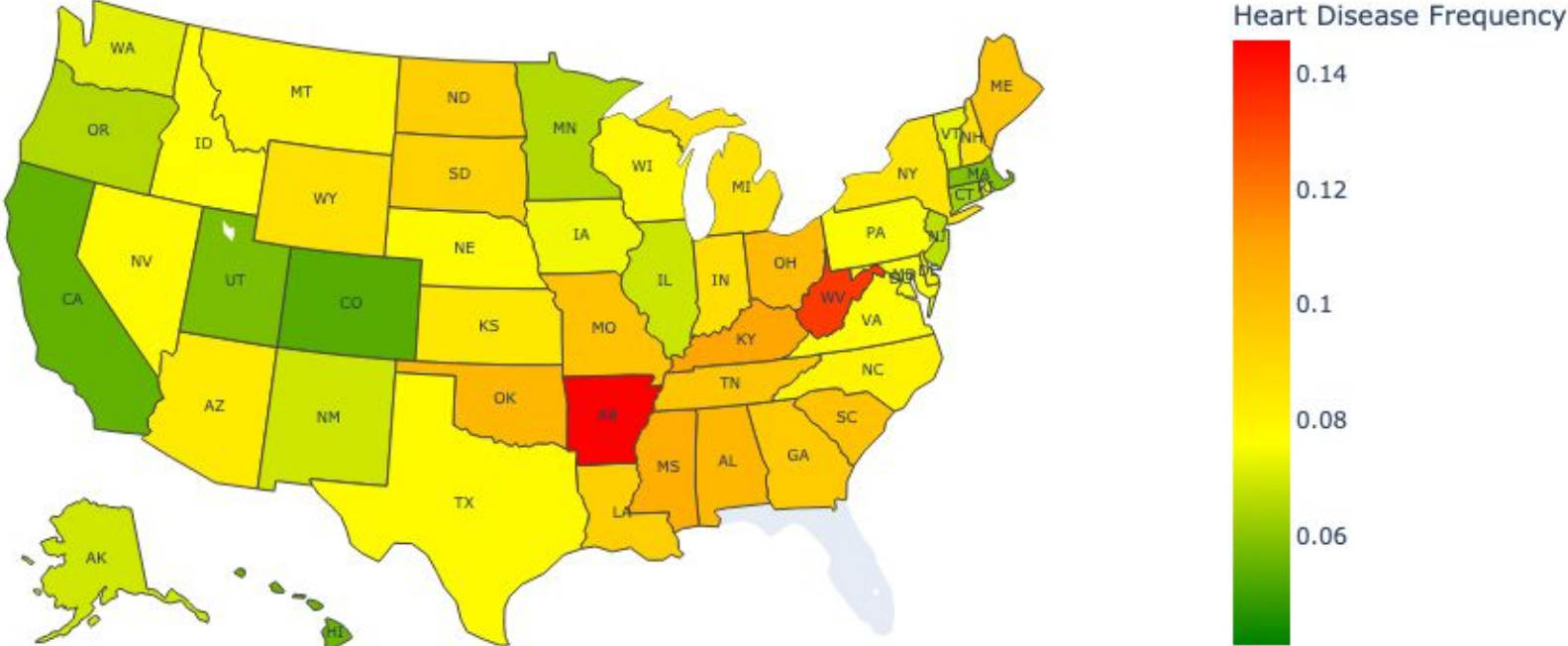
Data Cleaning

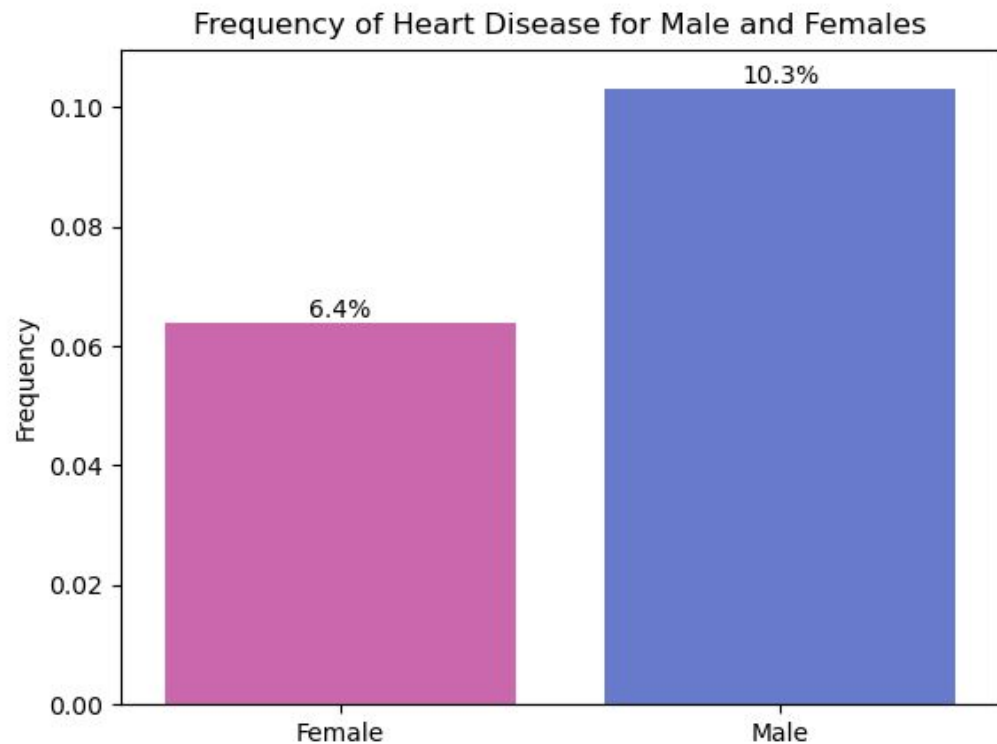
- Identify and remove irrelevant columns
- Remove columns with excessive missing values
- Missing values replacement
- Numeric to category conversion
- Missing Value imputation

Imbalanced Data



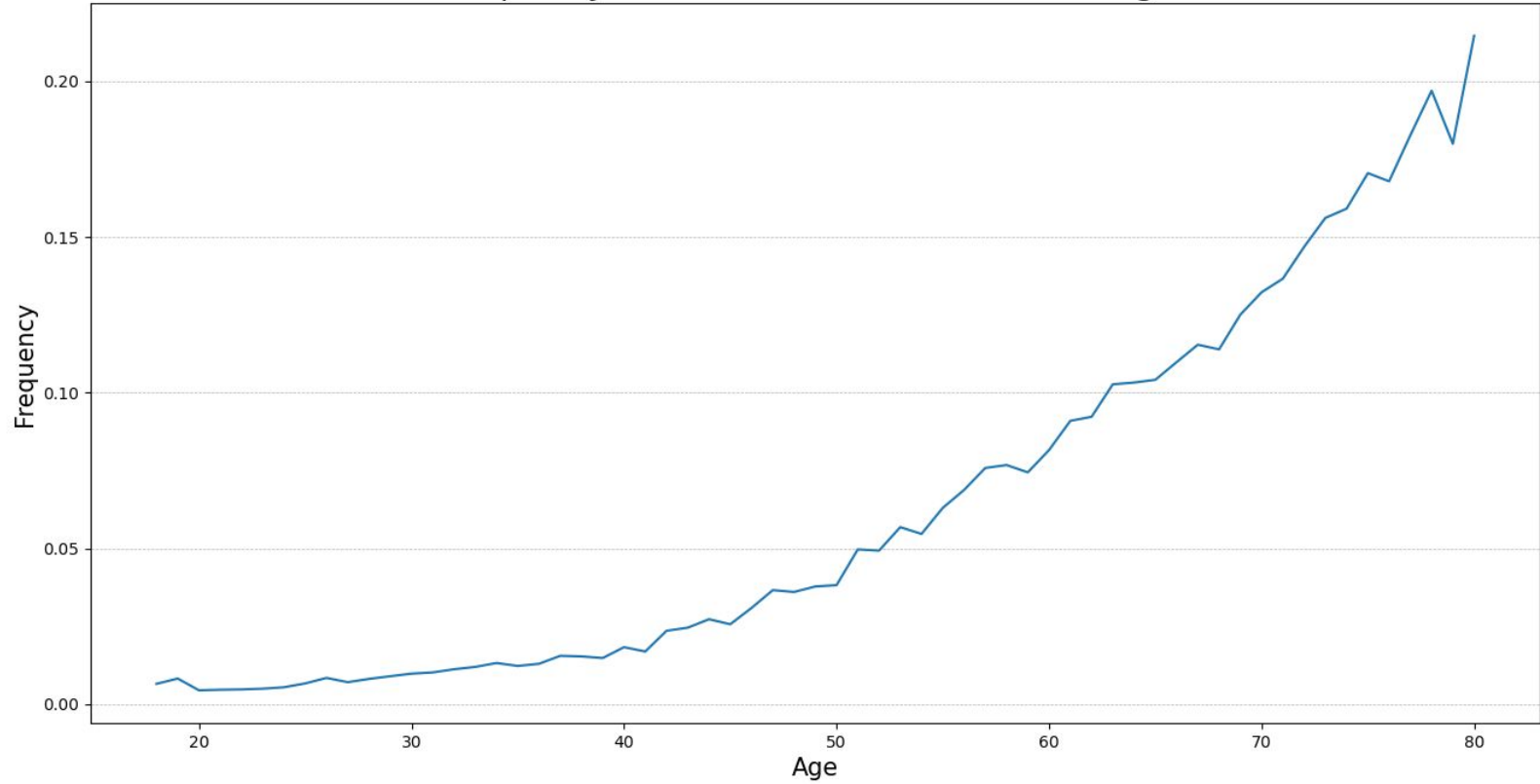
Heart Disease Frequency by State



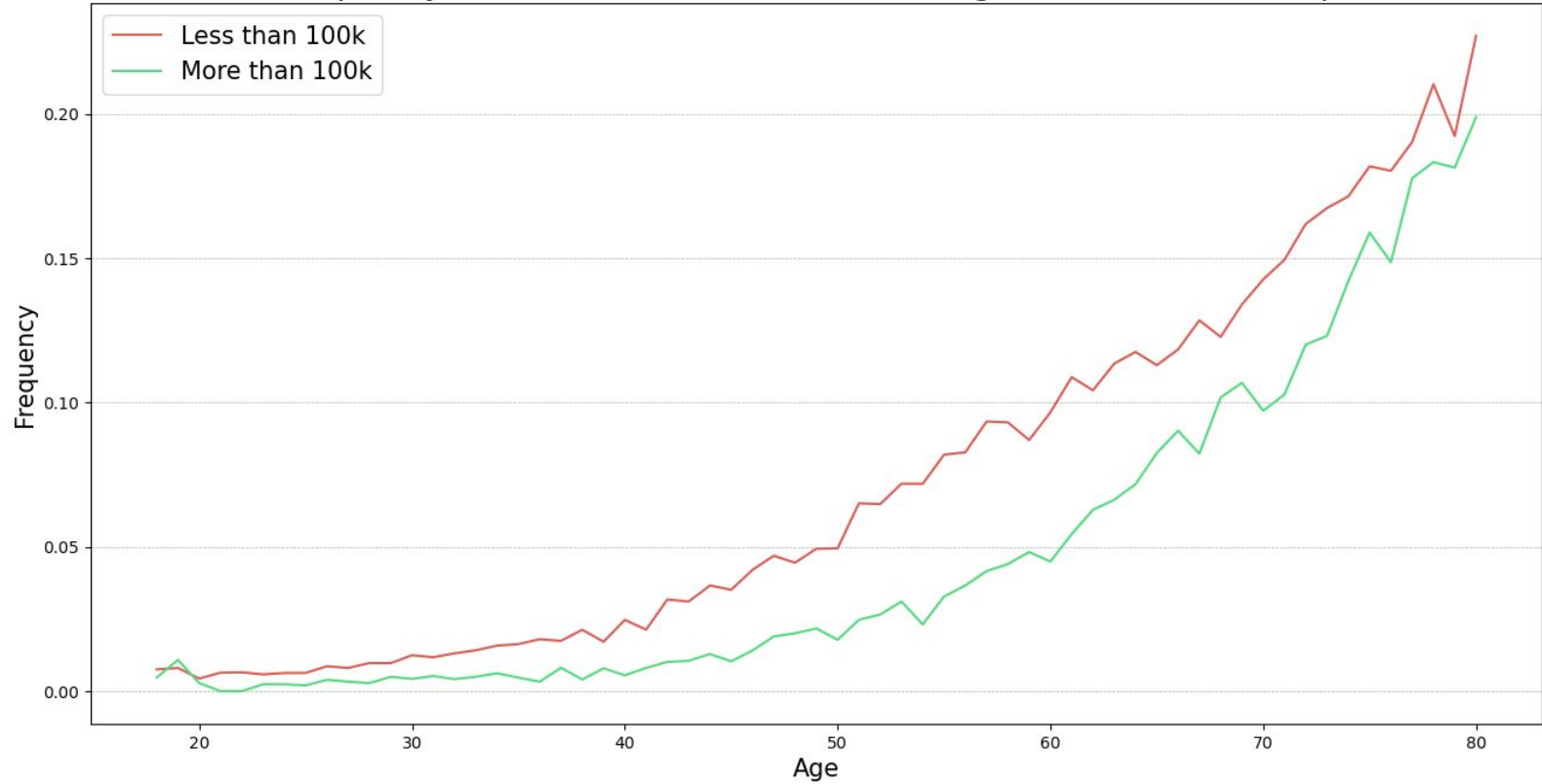


Women are less likely to have heart disease

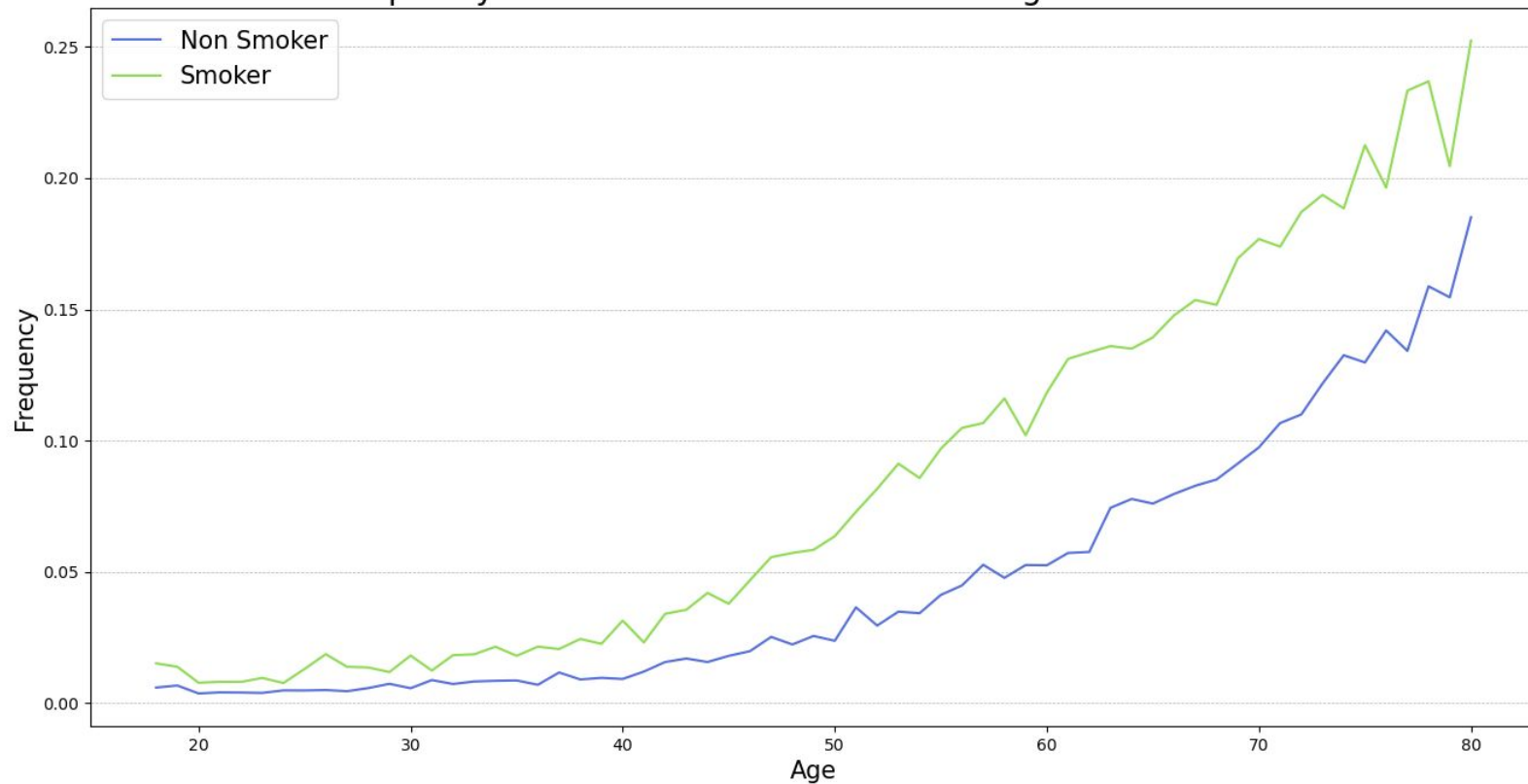
Frequency of Heart Disease For Different Ages




Frequency of Heart Disease For Different Ages and Income Groups

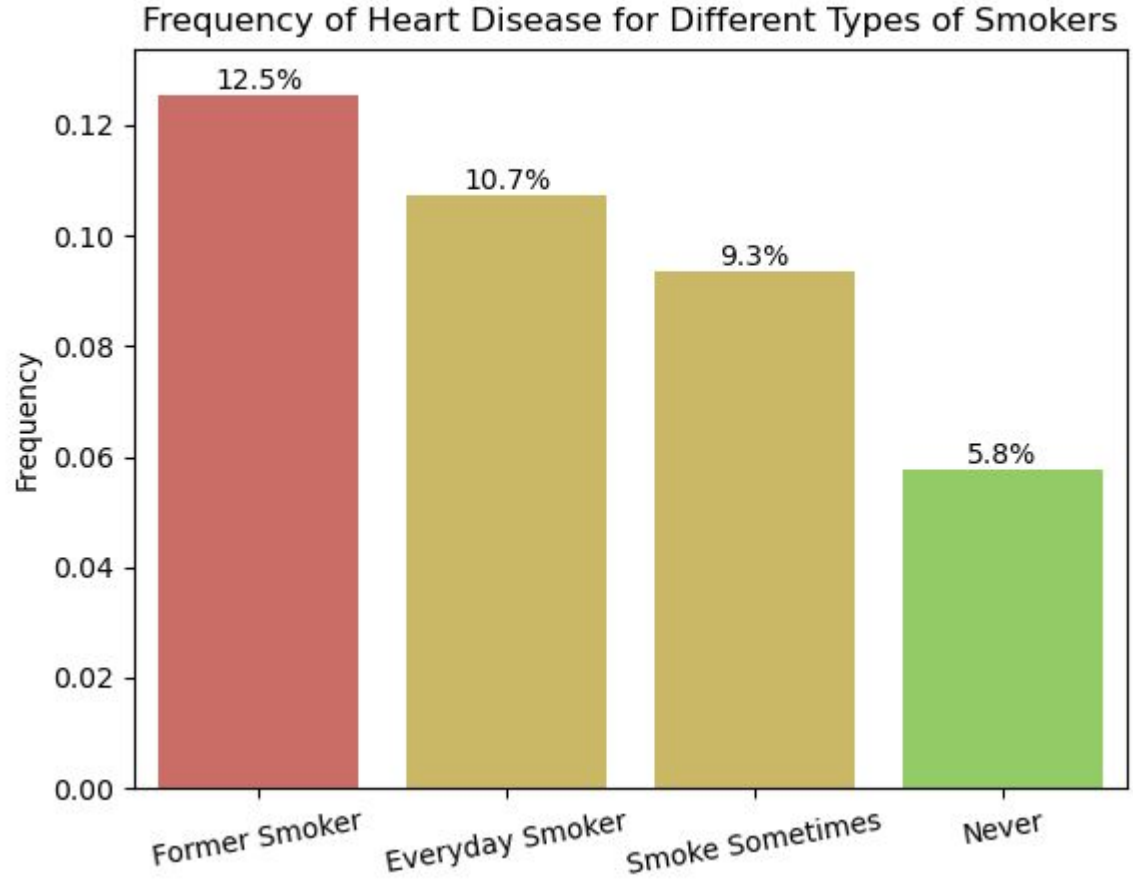


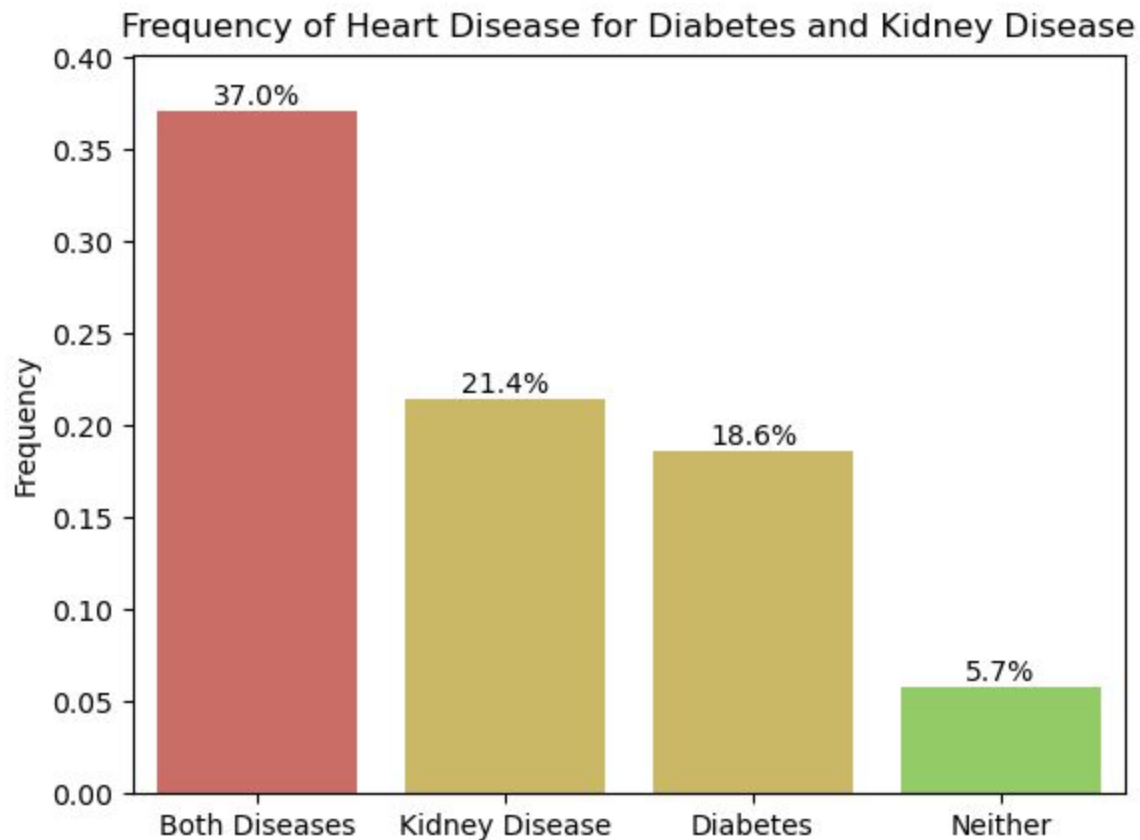
Frequency of Heart Disease For Different Ages and Smokers





**Smokers are
more likely to
have Heart
Disease**





“When it rains, it pours”

You are more likely to have heart disease if you have other underlying diseases



Modeling Techniques

- Models: Multinomial Bayes, Logistic, XGBoost, and Random Forest classification models.
- Handling imbalance classes: Oversampling and Overweight minority class, Resampling with SMOTE and ADASYN.



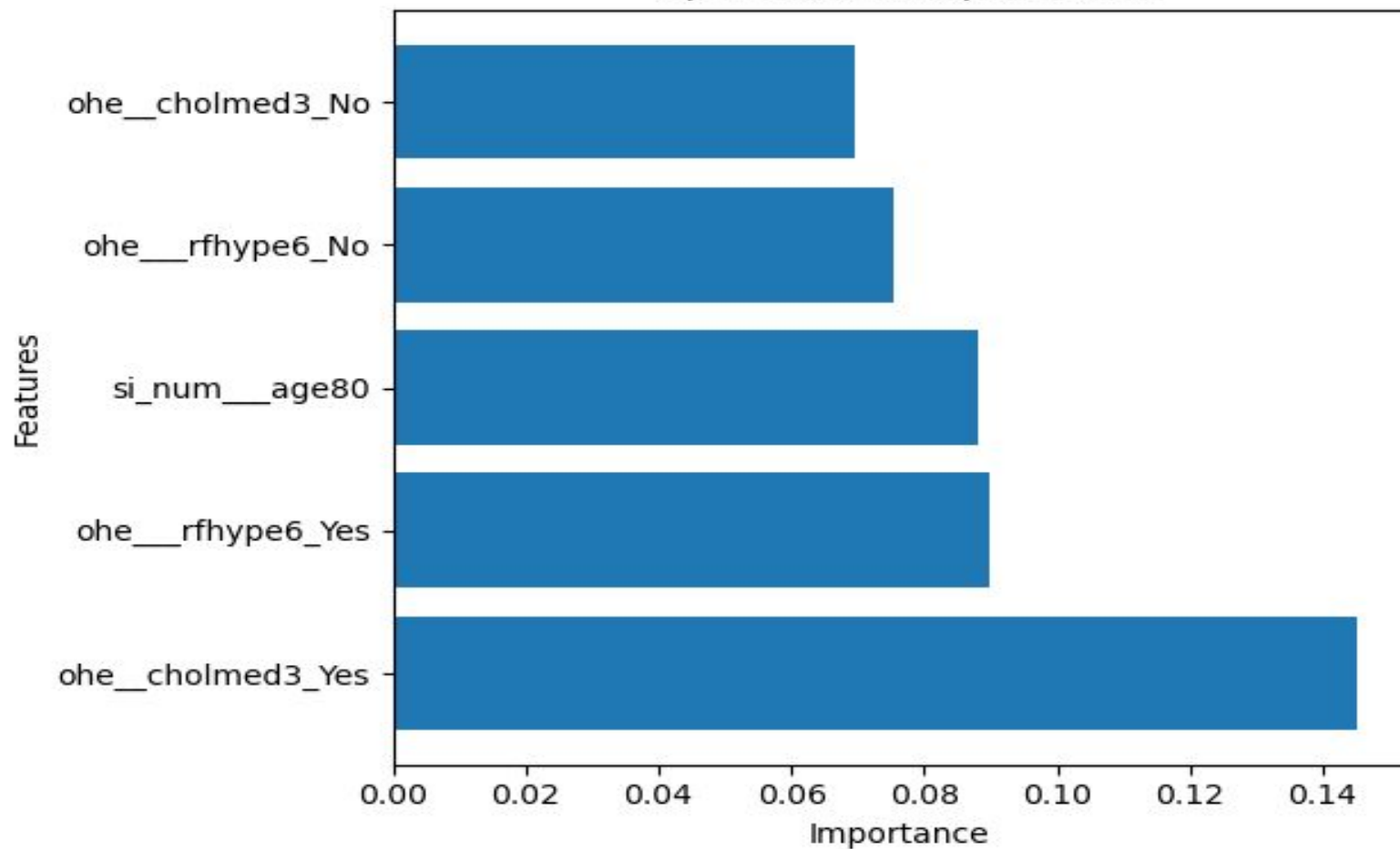
Modeling Evaluation

Recall = True Positive Rate

	balanced_accuracy	recall
Type of Model		
Best NB Model	0.751111	0.754287
Best Logistic Model	0.784727	0.811017
XGBoost Model	0.786722	0.819591
Best Random Forest Model	0.772020	0.826056

Random Forest Model = Best Model

Top 5 Feature Importances





Conclusion

- Model Application: Predicting the probability that a patient has heart disease will be useful as a screening process for healthcare professionals.
- Our model had a high balanced accuracy and an even higher recall score.
- Since the model was chosen to be sensitive towards positives, this can be used as a tool to identify people who are likely to have heart disease.