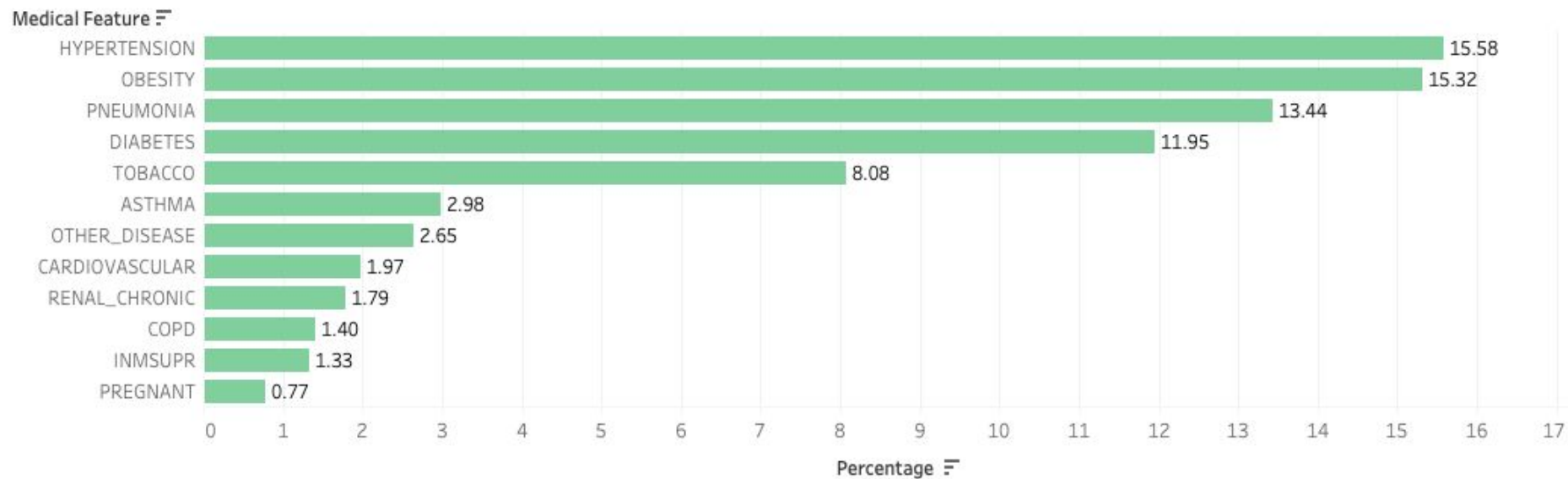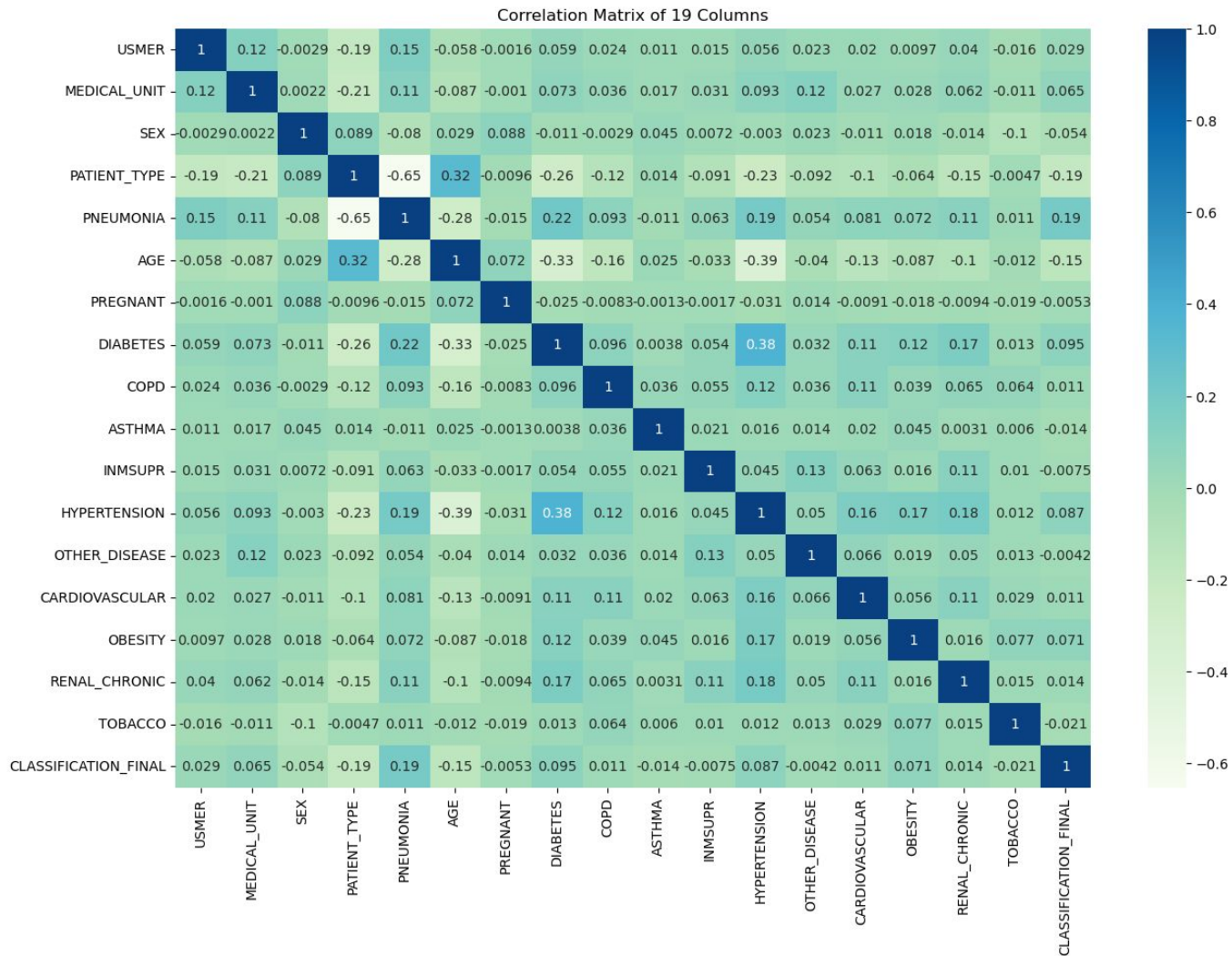# COVID-19

# Overview

- Exploration of different medical features that affect survival rate of COVID including:
    - Hypertension
    - Obesity
    - Pneumonia
    - Diabetes
- Model prediction of survival rate

# Medical Features



Medical Feature ≡

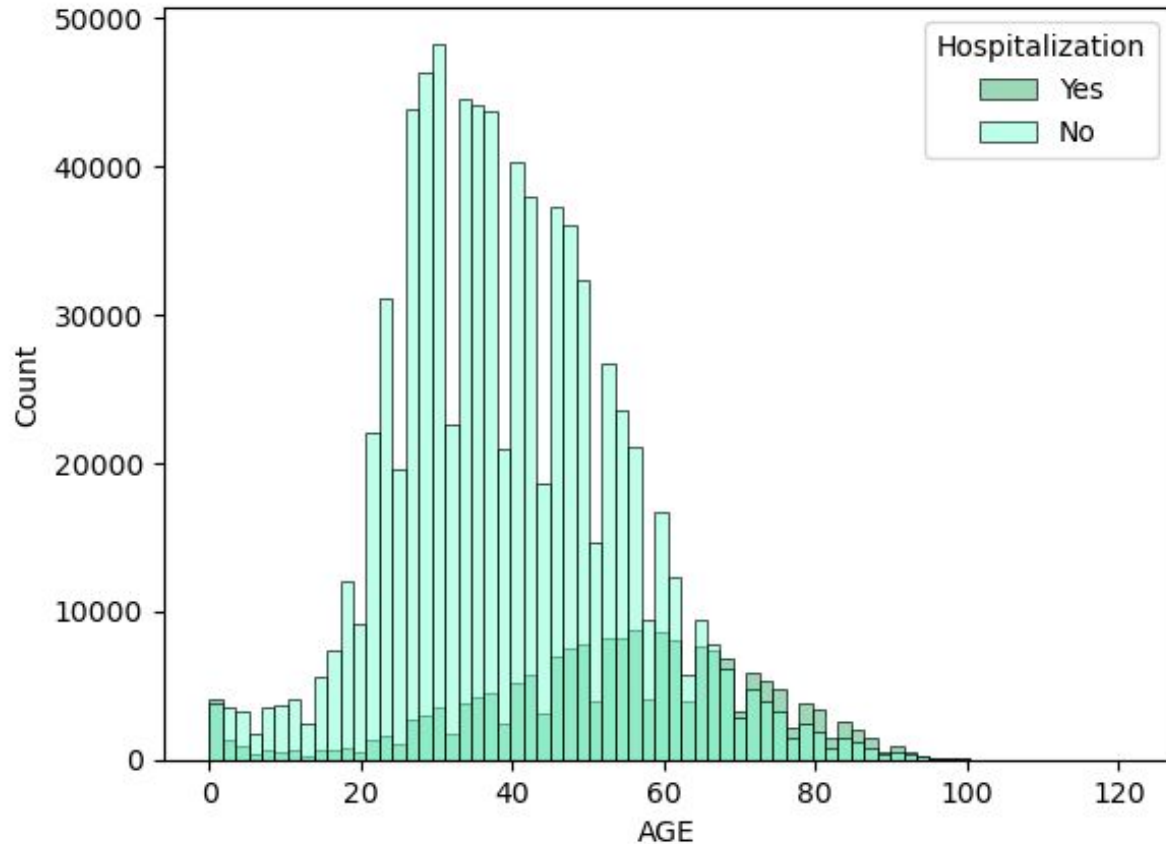| Medical Feature | Percentage |
|---|---|
| HYPERTENSION | 15.58 |
| OBESITY | 15.32 |
| PNEUMONIA | 13.44 |
| DIABETES | 11.95 |
| TOBACCO | 8.08 |
| ASTHMA | 2.98 |
| OTHER_DISEASE | 2.65 |
| CARDIOVASCULAR | 1.97 |
| RENAL_CHRONIC | 1.79 |
| COPD | 1.40 |
| INMSUPR | 1.33 |
| PREGNANT | 0.77 |

Percentage ≡

# Correlation of features

- Pneumonia and patient type are negatively correlated at -0.65, followed by hypertension and age at -0.39 then hypertension and diabetes at 0.38.
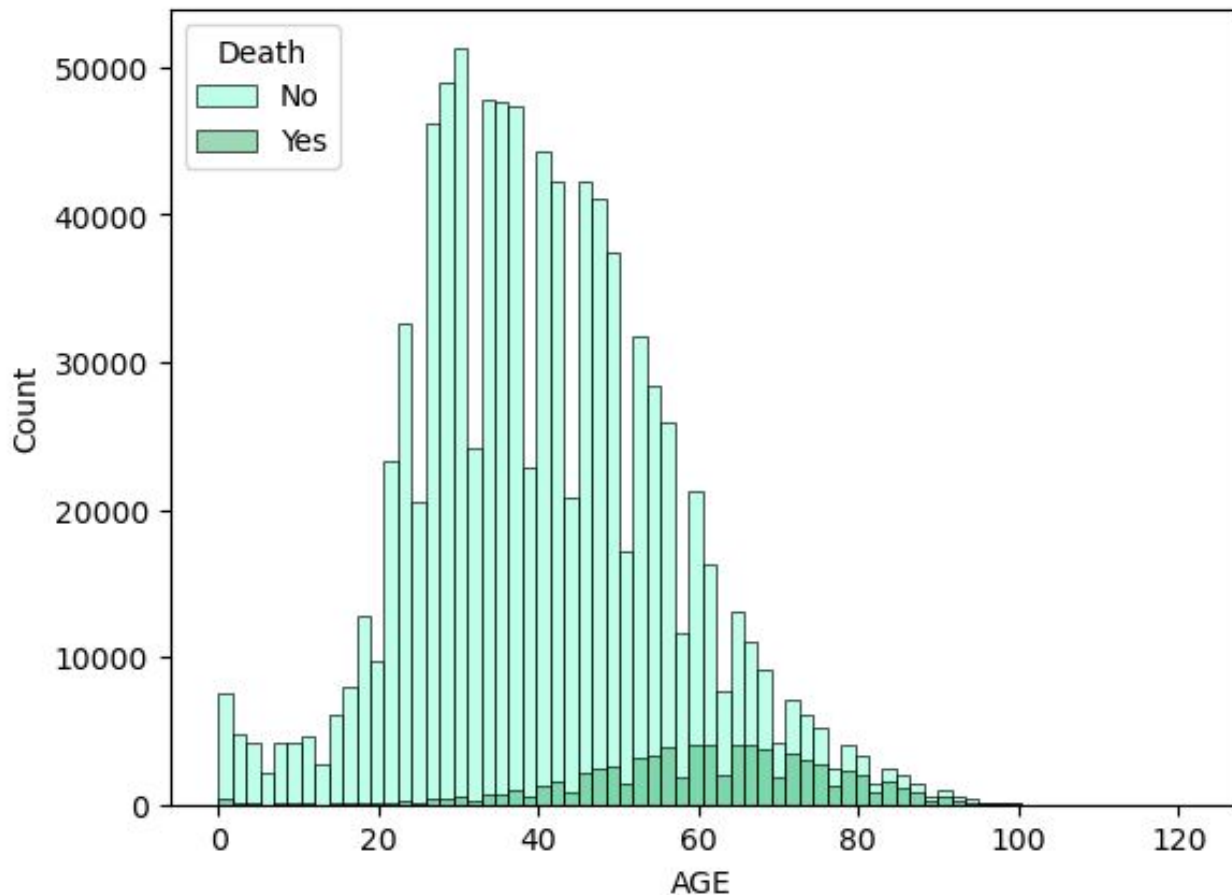


Correlation Matrix of 19 Columns

# Age vs Hospitalization

- Hospitalization occurs more frequently between aged 50-80 and also in young infants.
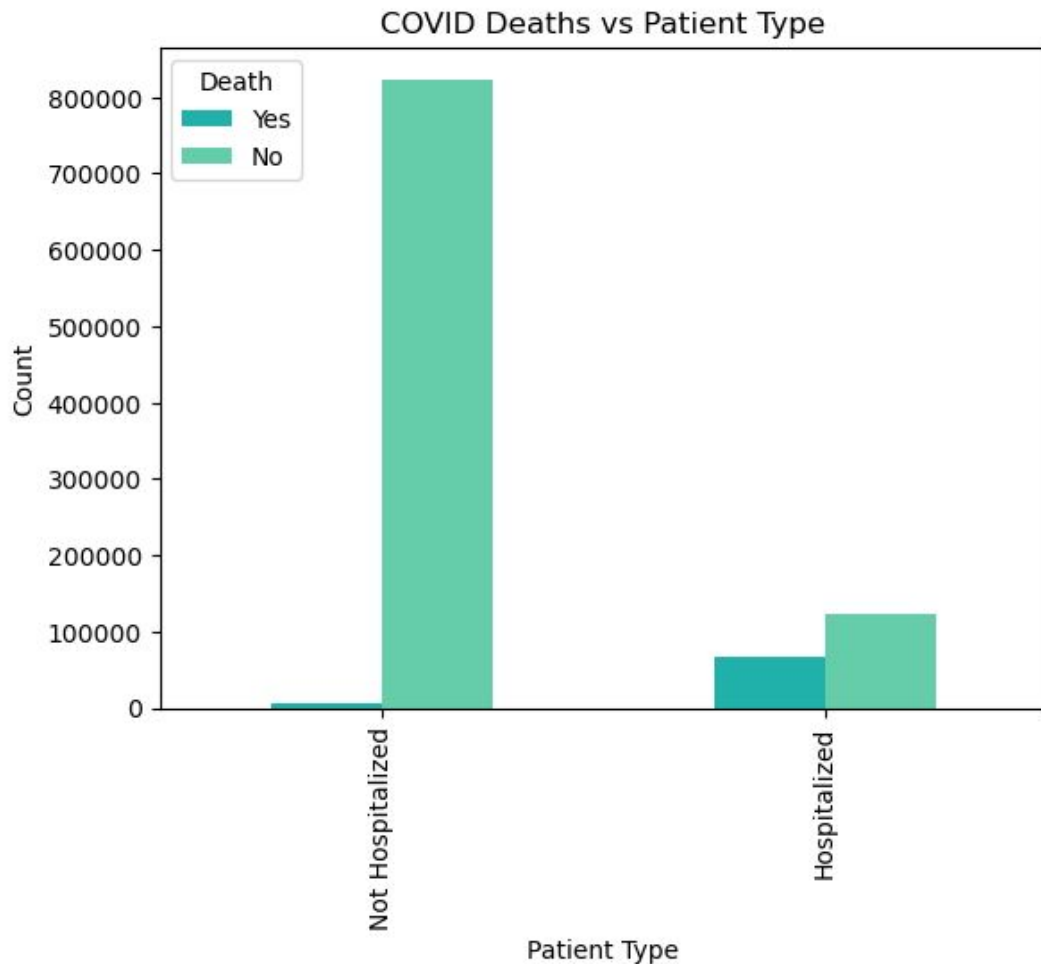
# Age vs Covid death

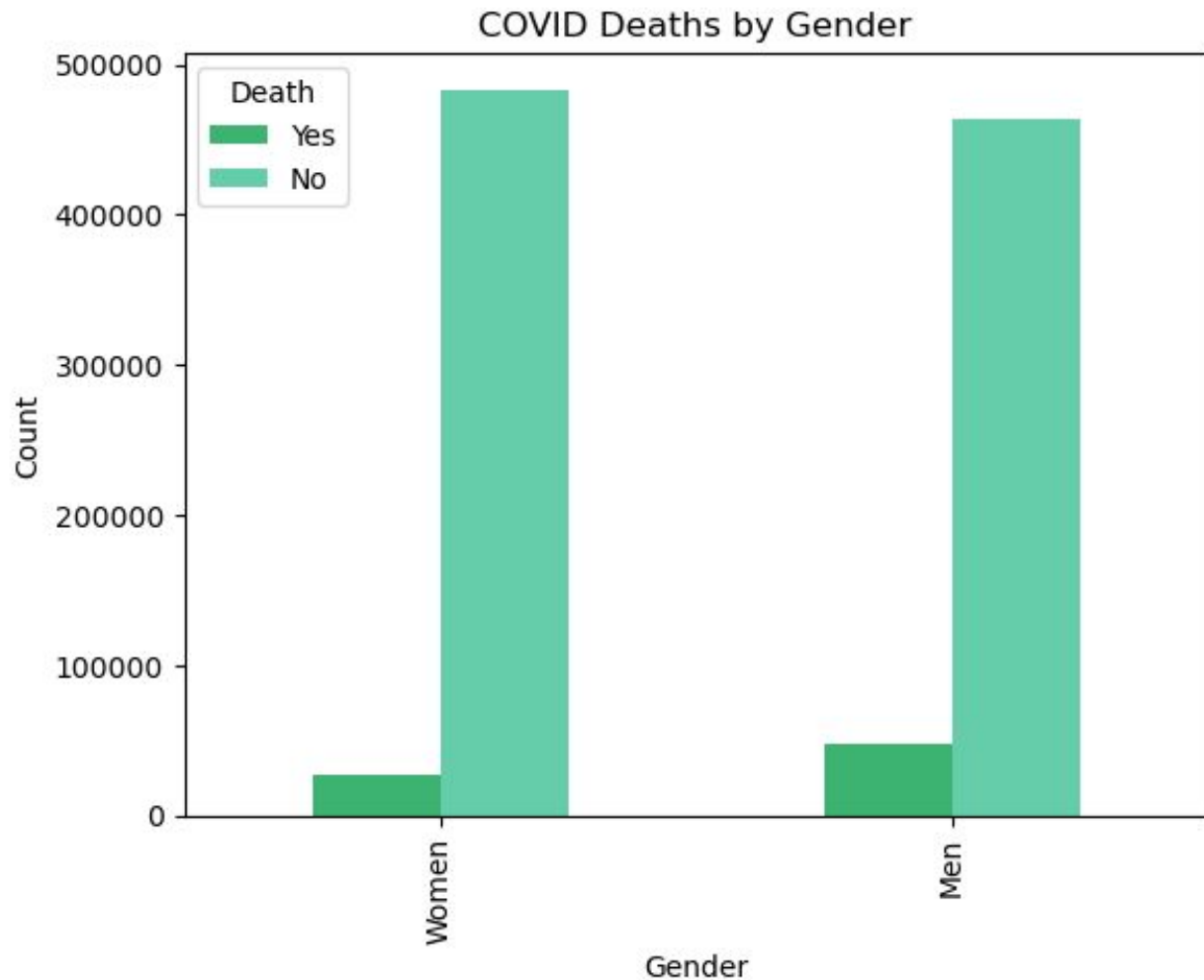- Death occurs most frequently between 50-80 inline with hospitalization.

# Hospitalization vs Covid Death

- Majority of patients are not hospitalized and of those only 0.8% have died.
- Of those who have been hospitalized death rate is much higher at 35%.



COVID Deaths vs Patient Type
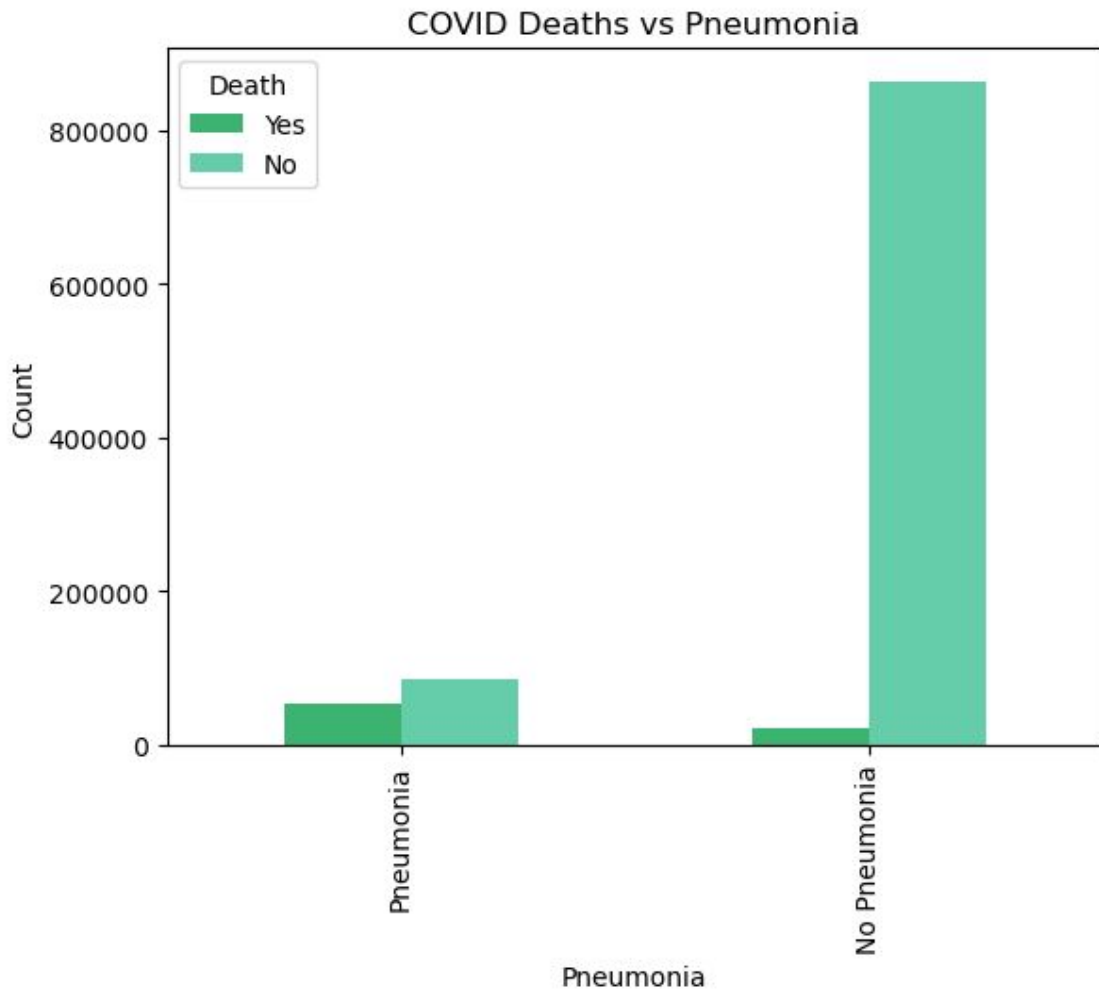
# Gender vs Covid Death

- Percentage of death in Female patients: 5.21 %
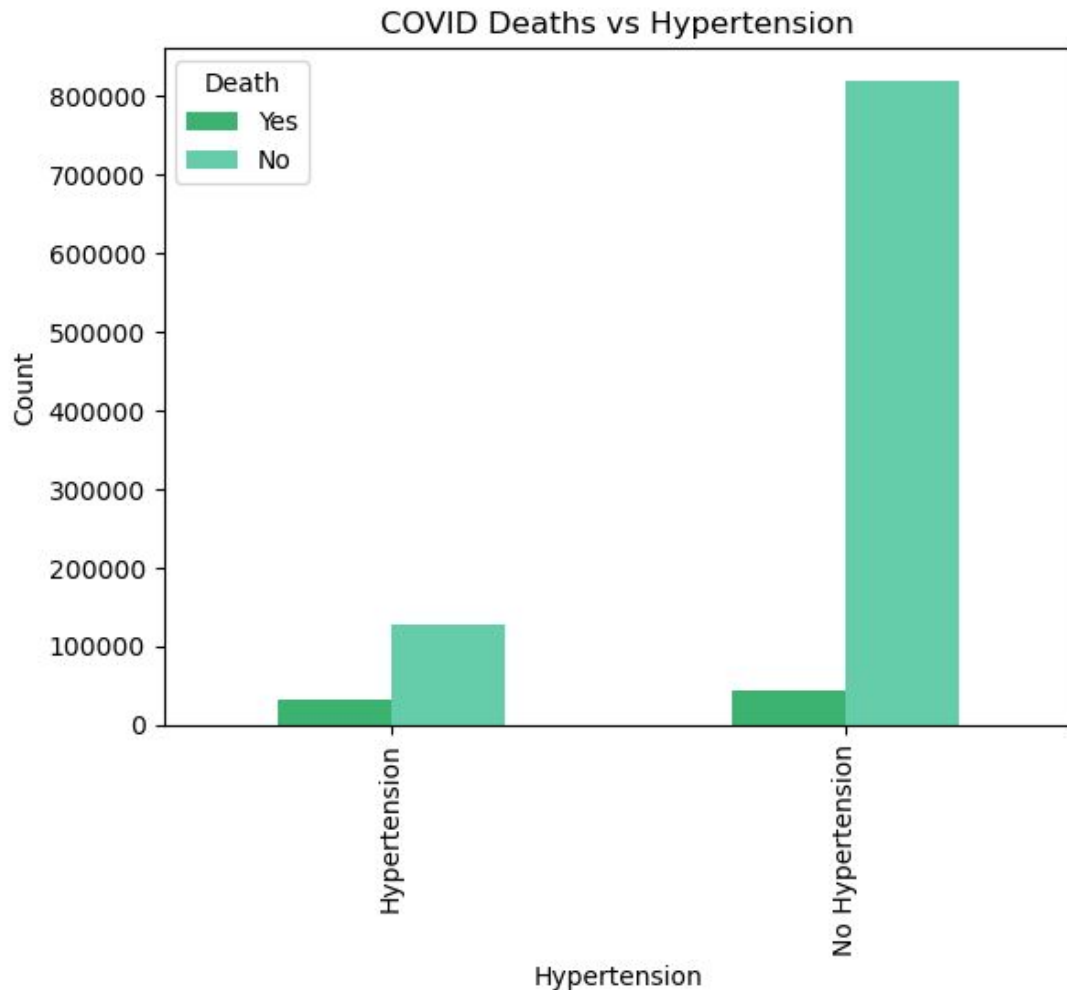- Percentage of death in Male patients: 9.39 %



COVID Deaths by Gender

# Pneumonia vs Covid Death

- Percentage of death in pneumonia patients: 38.4 %
- Percentage of death in non-pneumonia patients: 2.48 %

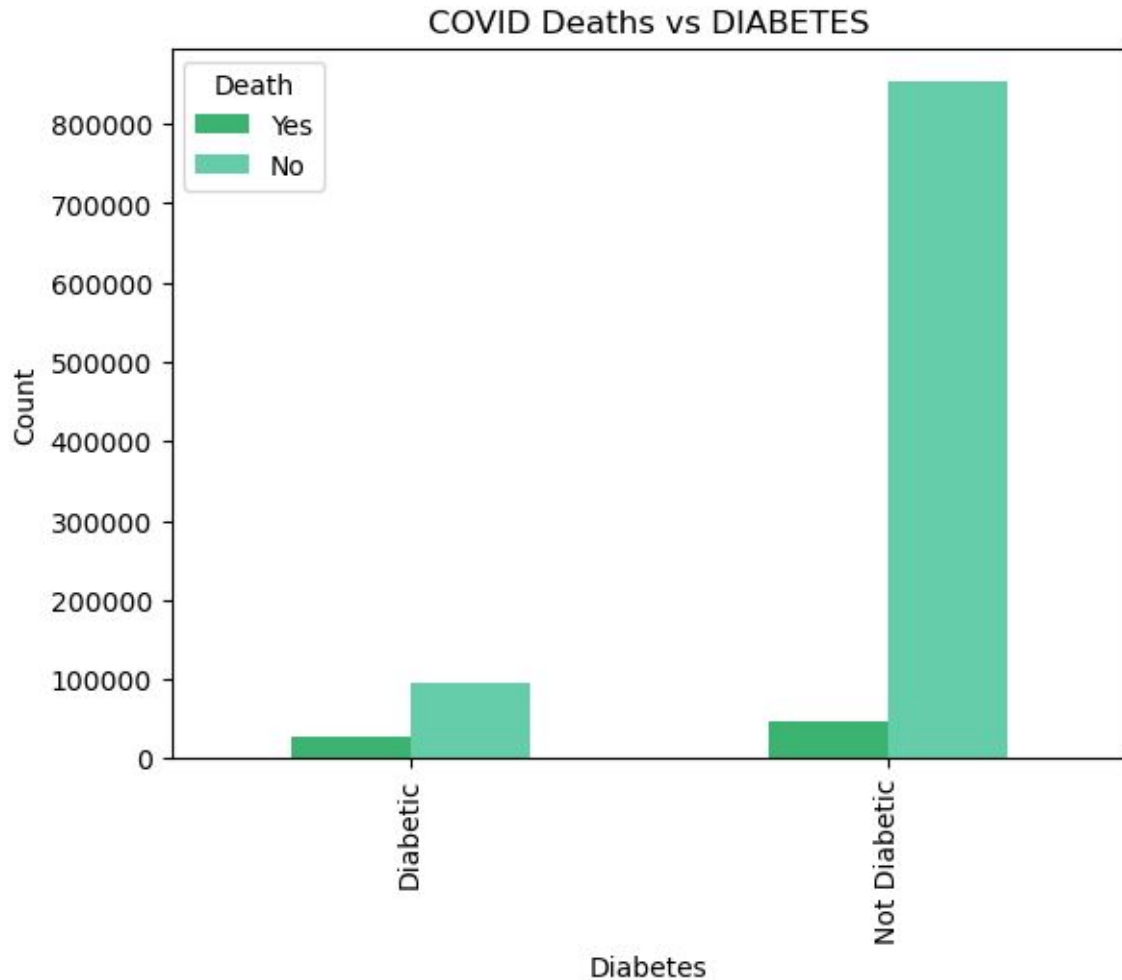

COVID Deaths vs Pneumonia

# Hypertension vs Covid Death

- Percentage of death in patients with Hypertension:  19.73 %
- Percentage of death in patients without Hypertension:  5.01 %



COVID Deaths vs Hypertension
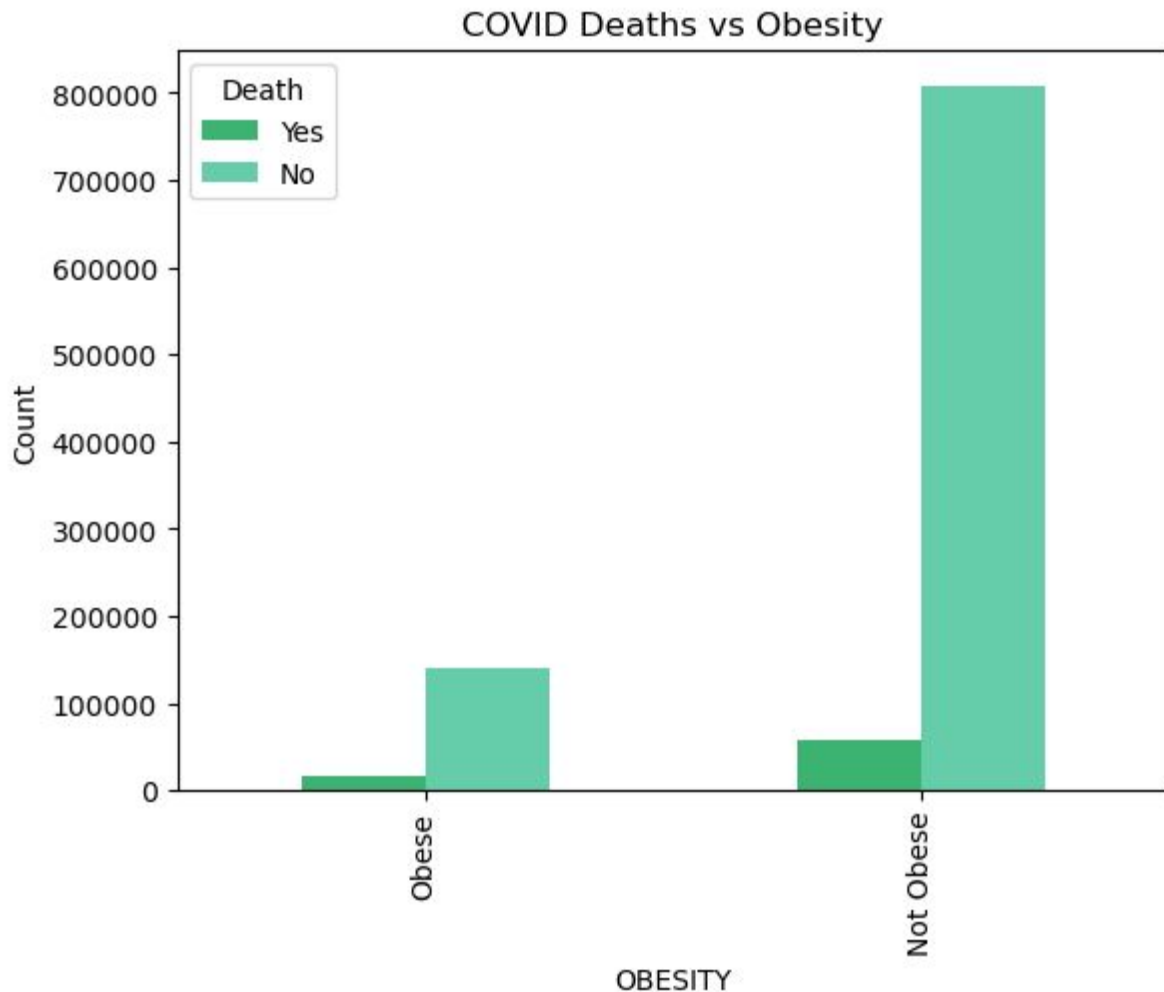
# Diabetes vs Covid Death

- Percentage of death in diabetic patients: 22.64 %
- Percentage of death in non-diabetic patients: 5.22 %

## COVID Deaths vs DIABETES

Death
- Yes
- No

# Obesity vs Covid Death

- Percentage of death in obese patients: 10.76 %
- Percentage of death in non-obese patients: 6.68 %



COVID Deaths vs Obesity

# Modelling the data

- The aim of the model is to predict the survival of a patient based on certain medical features they may have.
- Dataset is heavily imbalanced in the target variable therefore different sampling techniques were tested in the models and evaluated.
- Models that were tested: Logistic regression, random forests, random forest with gridsearch cross validation, over/under sampling, SMOTE, variance threshold and RFE.

# Model Evaluations

| Model | Class Imbalance | Under Sampling | Over Sampling | SMOTE | Variance Threshold | RFE |
|---|---|---|---|---|---|---|
| **Logistic Regression** | 0.47 | 0.51 | 0.51 | 0.51 | 0.41 | 0.37 |
| **Random Forest** | 0.32 | 0.49 | 0.49 | 0.48 | - | - |
| **Random Forest GS** | 0.47 | 0.83 | **0.88** | - | - | - |

- Random forest with grid search cross validation and oversampling is the best model with a Kappa score of 0.88.
- This indicates a high level of agreement between the model predictions and the actual outcomes with high accuracy.

| Final Model | Kappa | Precision | Recall | F1-Score |
|---|---|---|---|---|
| **Random Forest GS** | 0.88 | 0.70 | 0.89 | 0.75 |