# Uncovering Patterns and Predicting Chronic Disease Risks in the U.S

*Van Le - Stanley He - James Li*

**6 IN 10**

Adults in the US have a **chronic disease**
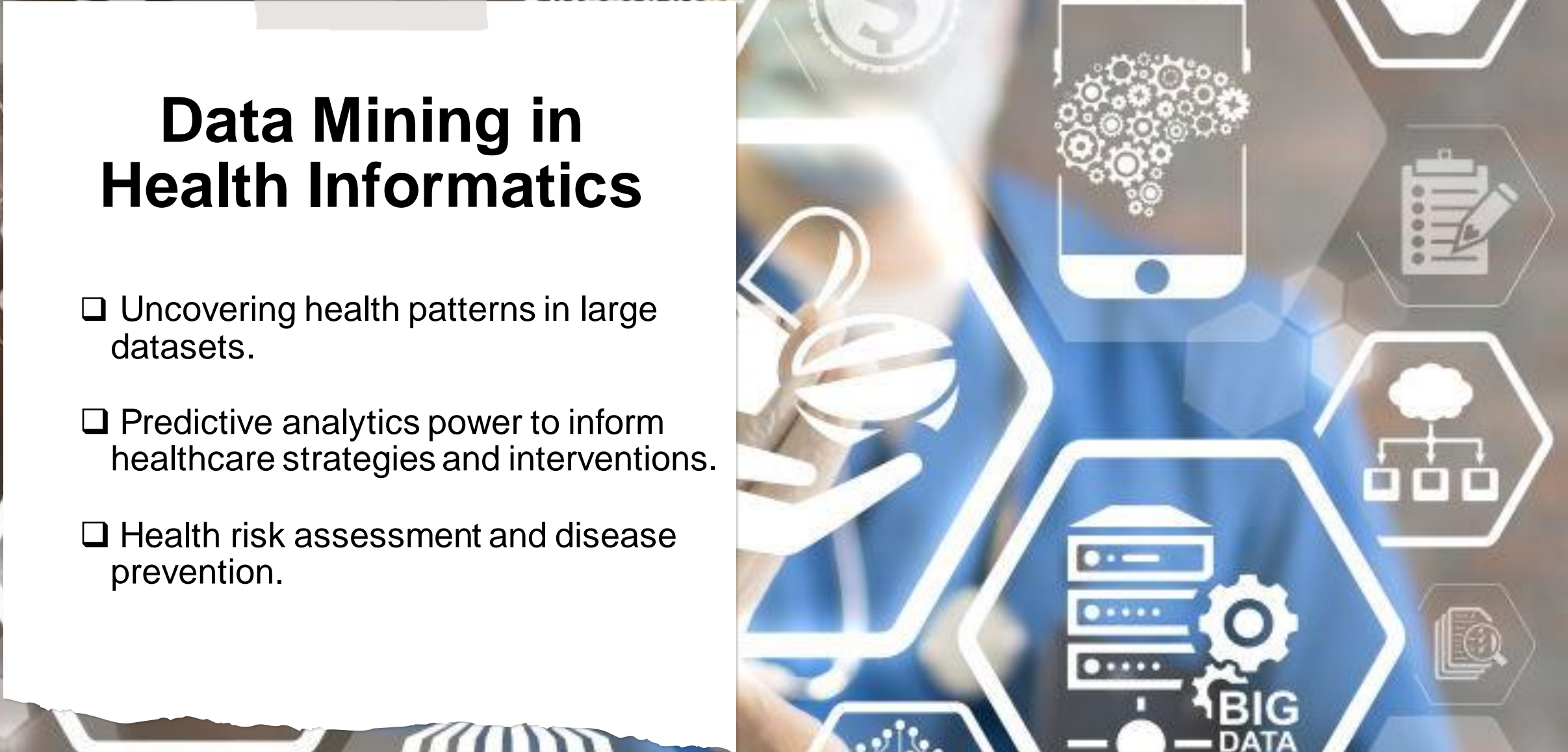
**4 IN 10**

Adults in the US have **two or more**

# The Burden of Chronic Diseases

❑ Chronic diseases are the leading causes of death and disability in the U.S.

❑ Importance to predict and manage disease outbreaks.

❑ These indicators are essential for surveillance, prioritization, and evaluation of public health interventions.

# Data Mining in Health Informatics

❑ Uncovering health patterns in large datasets.

❑ Predictive analytics power to inform healthcare strategies and interventions.

❑ Health risk assessment and disease prevention.

# Project Objectives

Classification of risk levels at the state level that are essential for the evaluation of public health interventions.

Pattern of chronic disease across different states.

Predictive models of mortality risk that important for early detection and future strategies.
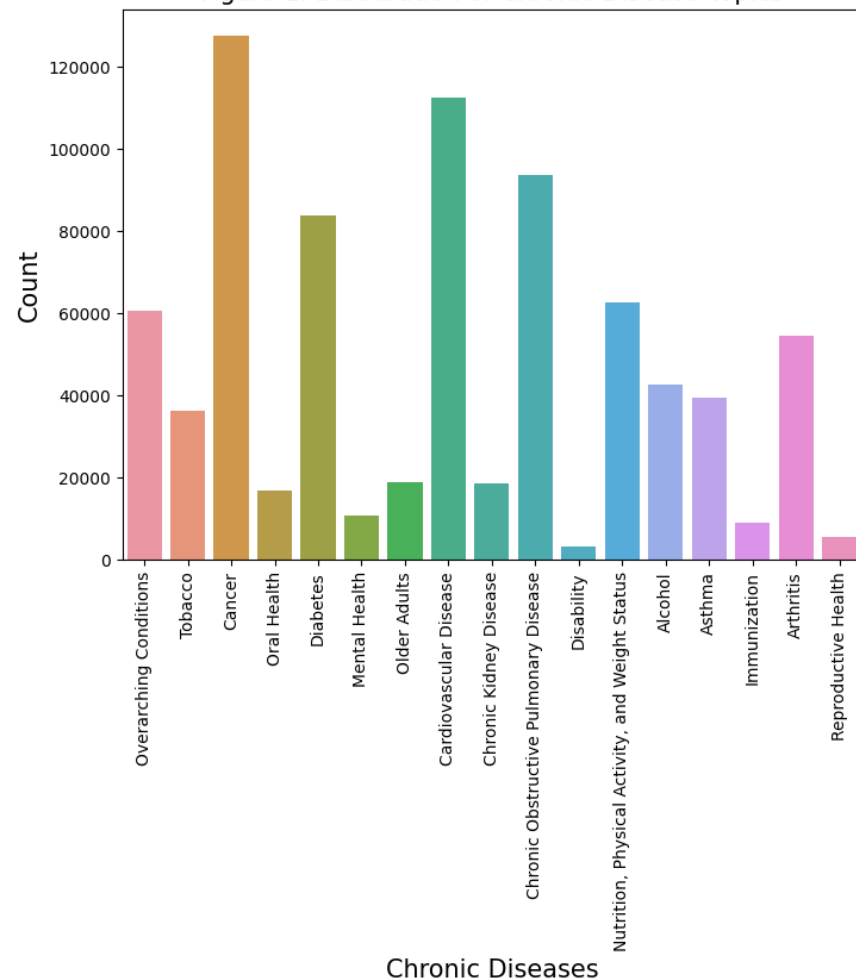
# Dataset Overview

❑ `U.S. Chronic Disease Indicators (CDI)` was developed by consensus that allows states and large metropolitan areas to uniformly define, collect, and report chronic disease data that are important to public health practice. (catalog.data.gov)

❑ The Chronic Condition Indicator was created to facilitate health services research on diagnoses using administrative data.

❑ There are more than 900K data rows in our dataset, with more than 24 indicators over 50 states that include 17 topic groups.
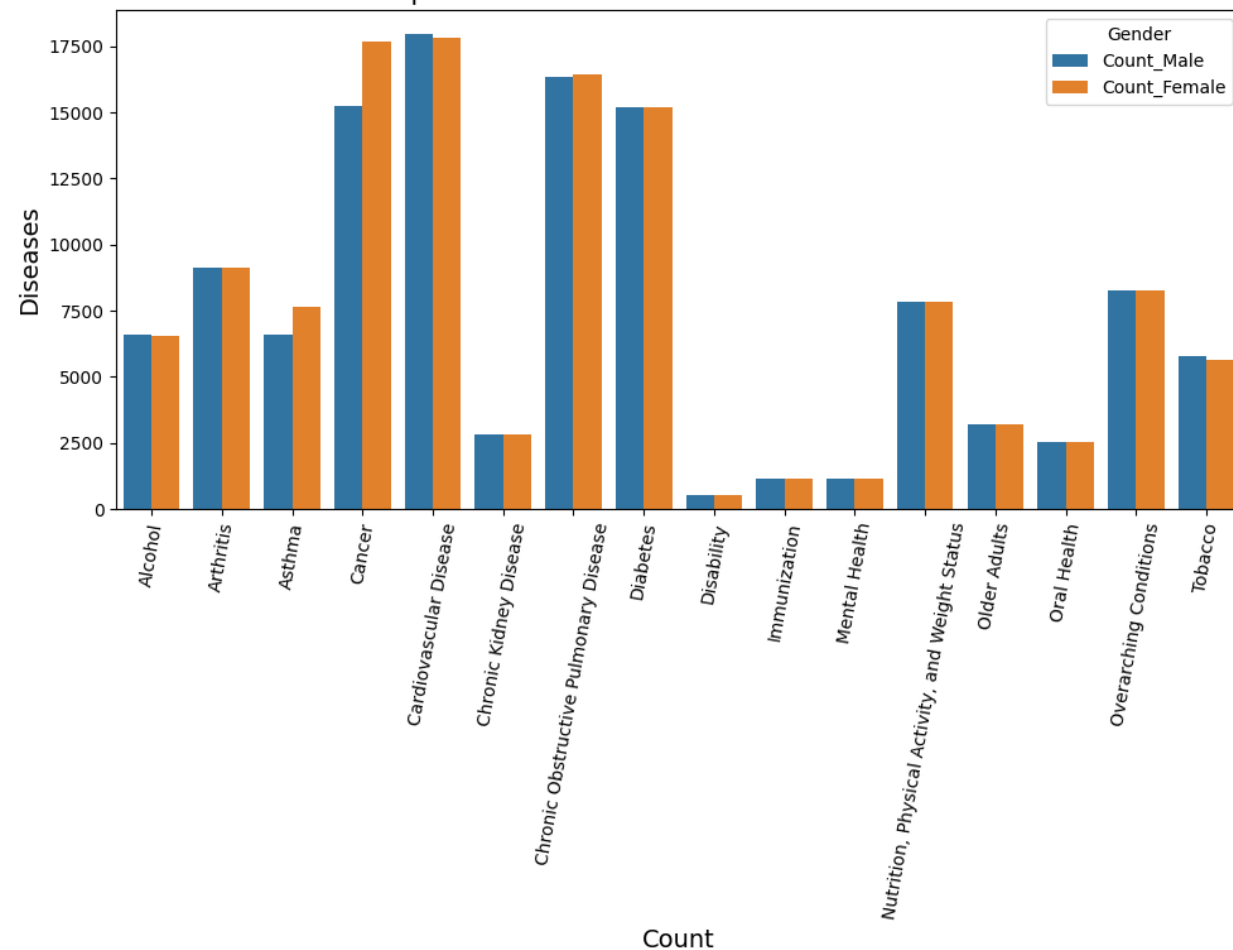
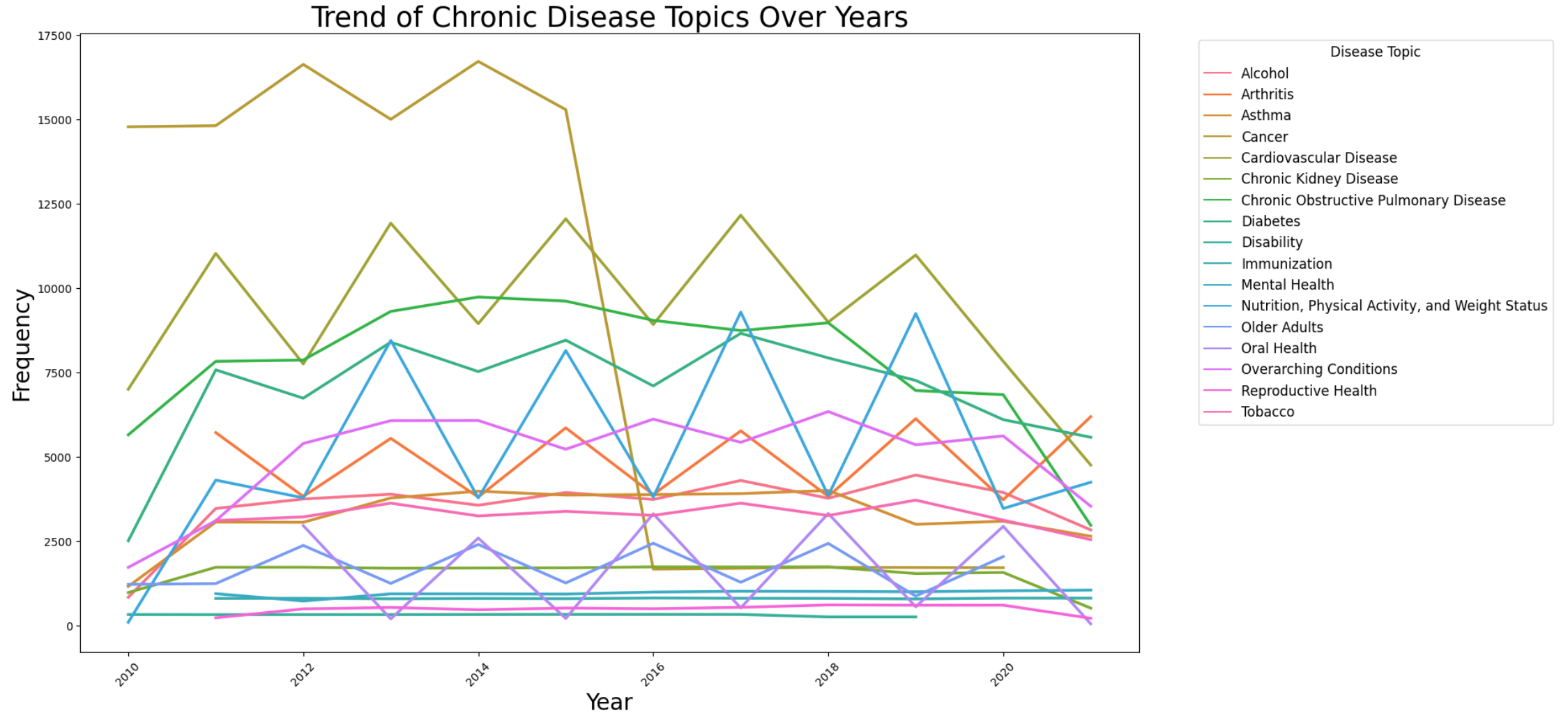| | YearStart | YearEnd | LocationAbbr | Topic | Question | DataValueUnit | DataValueType | DataValue | DataValueAlt | StratificationCategory1 | Stratification1 | LocationID | TopicID | QuestionID | DataValueTypeID | StratificationCategoryID1 | StratificationID1 | Longitude | Latitude | DeadStatus |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2001 | 2001 | NJ | Overarching Conditions | Life expectancy at birth | Years | Number | 77.6 | 77.6 | Overall | Overall | 34 | OVC | OVC4_1 | NMBR | OVERALL | OVR | -74.273691 | 40.13057 | 1 |
| 1 | 2001 | 2001 | OH | Overarching Conditions | Life expectancy at birth | Years | Number | 76.5 | 76.5 | Overall | Overall | 39 | OVC | OVC4_1 | NMBR | OVERALL | OVR | -82.404260 | 40.06021 | 1 |
| 2 | 2001 | 2001 | NJ | Overarching Conditions | Life expectancy at age 65 years | Years | Number | 18.1 | 18.1 | Overall | Overall | 34 | OVC | OVC4_2 | NMBR | OVERALL | OVR | -74.273691 | 40.13057 | 1 |
| 3 | 2001 | 2001 | MA | Overarching Conditions | Life expectancy at age 65 years | Years | Number | 18.7 | 18.7 | Overall | Overall | 25 | OVC | OVC4_2 | NMBR | OVERALL | OVR | -72.082691 | 42.27687 | 1 |
| 4 | 2001 | 2001 | KS | Overarching Conditions | Life expectancy at birth | Years | Number | 77.8 | 77.8 | Overall | Overall | 20 | OVC | OVC4_1 | NMBR | OVERALL | OVR | -98.200781 | 38.34774 | 1 |

# Data Summary & Statistics



Figure 1. Distribution of Chronic Disease Topics



Comparison of Disease between Male and Female
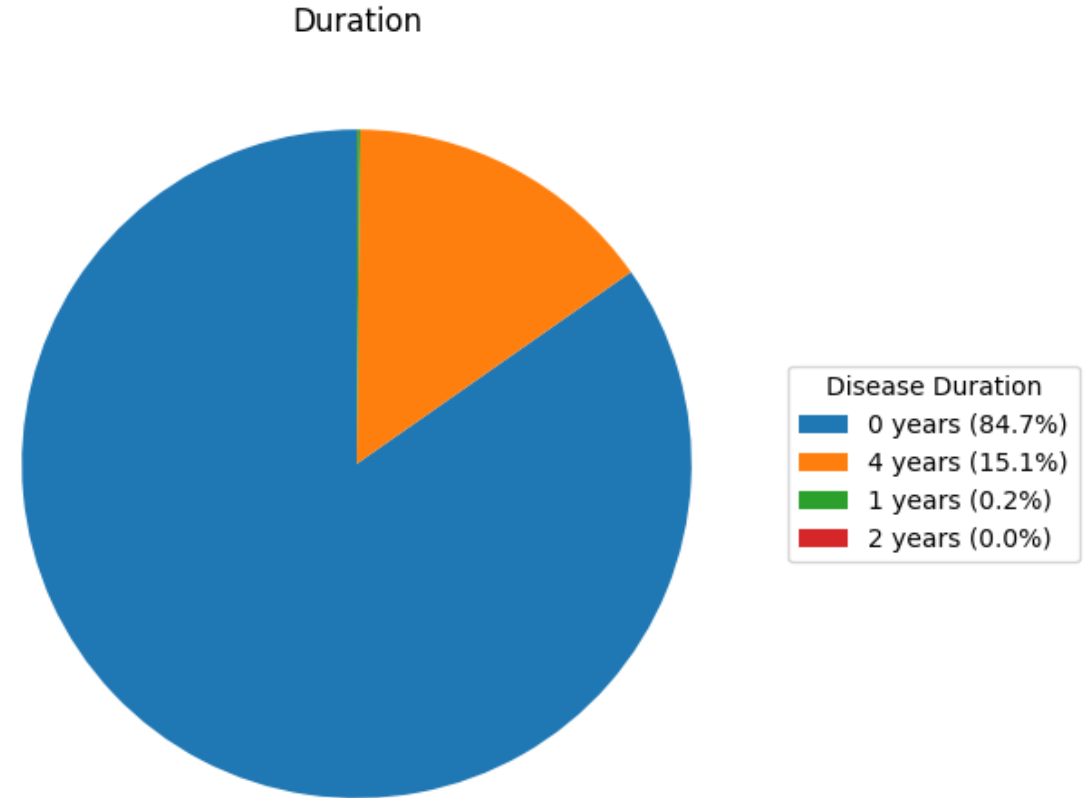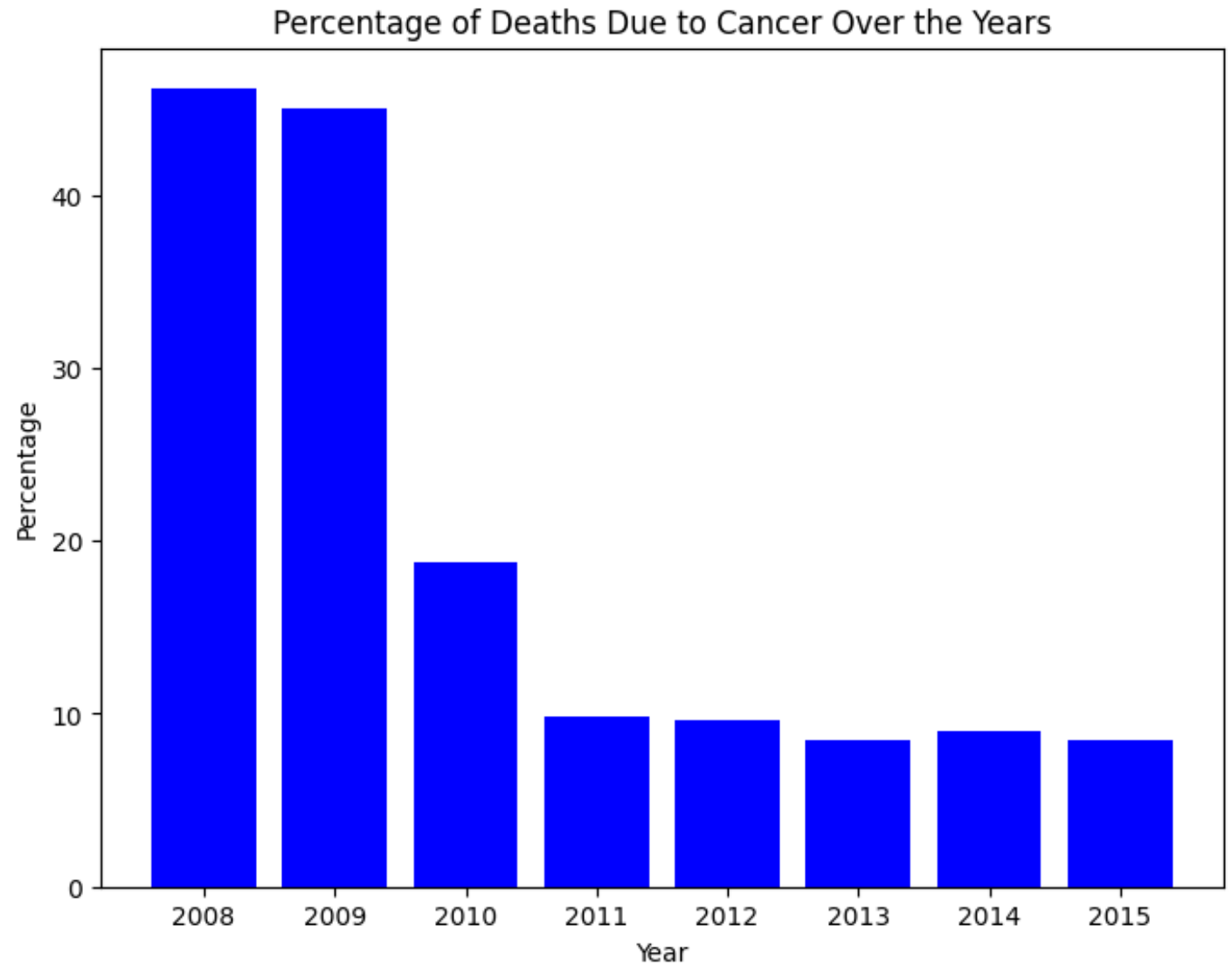
# Trends of Topic Over Years



Trend of Chronic Disease Topics Over Years

# Recovery Time Analysis

❑ The majority recovered swiftly during the initial year.

❑ A significant portion required 4 years for recovery.
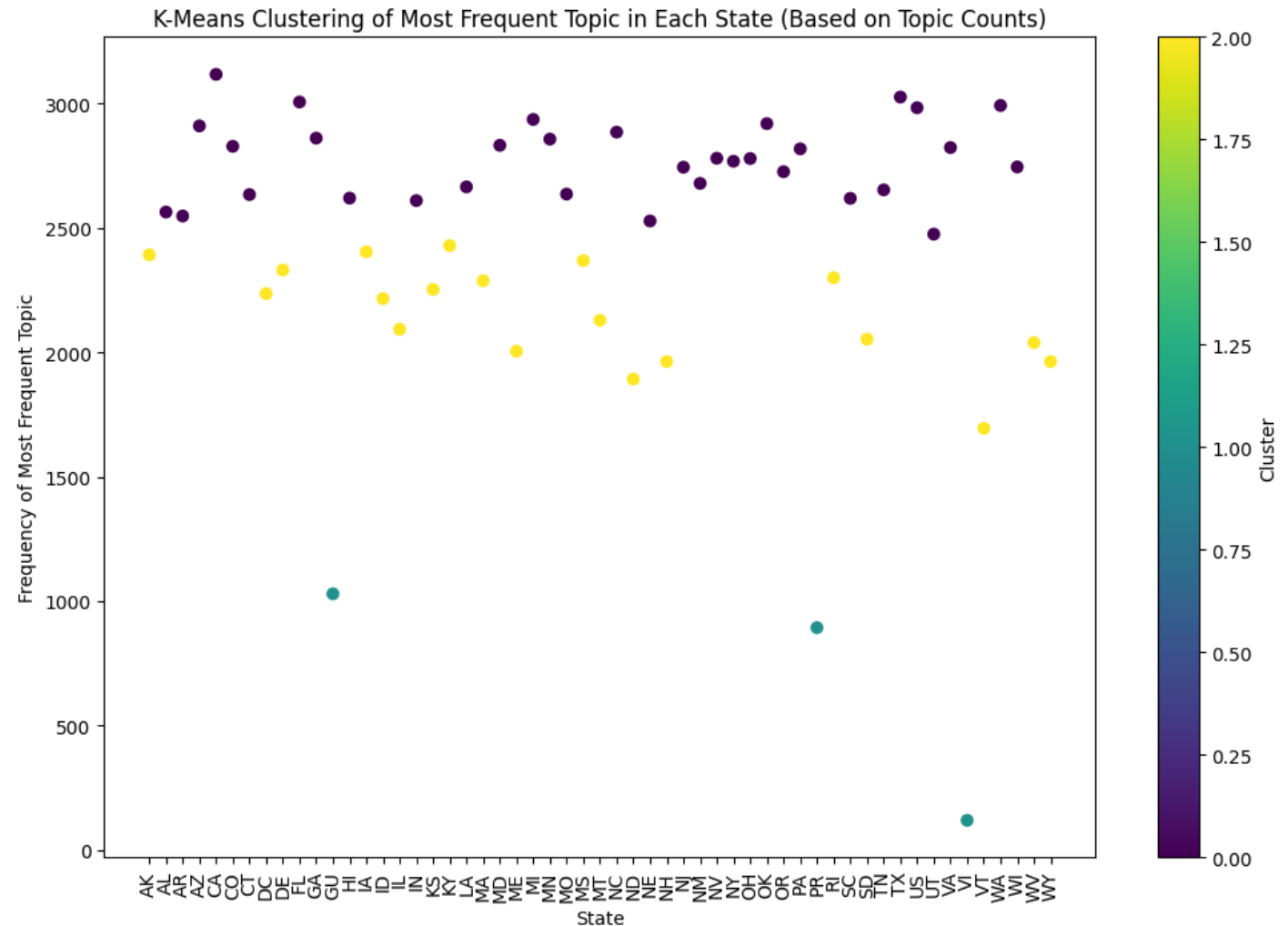
❑ A notable subset recovered after 1 or 2 years.



Duration

Disease Duration
- 0 years (84.7%)
- 4 years (15.1%)
- 1 years (0.2%)
- 2 years (0.0%)

# K-means Clustering Analysis

❑ Purpose: Determine patterns in disease prevalence across states.

❑ Input data consists of the health issues or the chronic diseases that are the most frequent in each state in the US.

❑ The K-means cluster analysis was chosen because it is relatively easy to understand and implement.
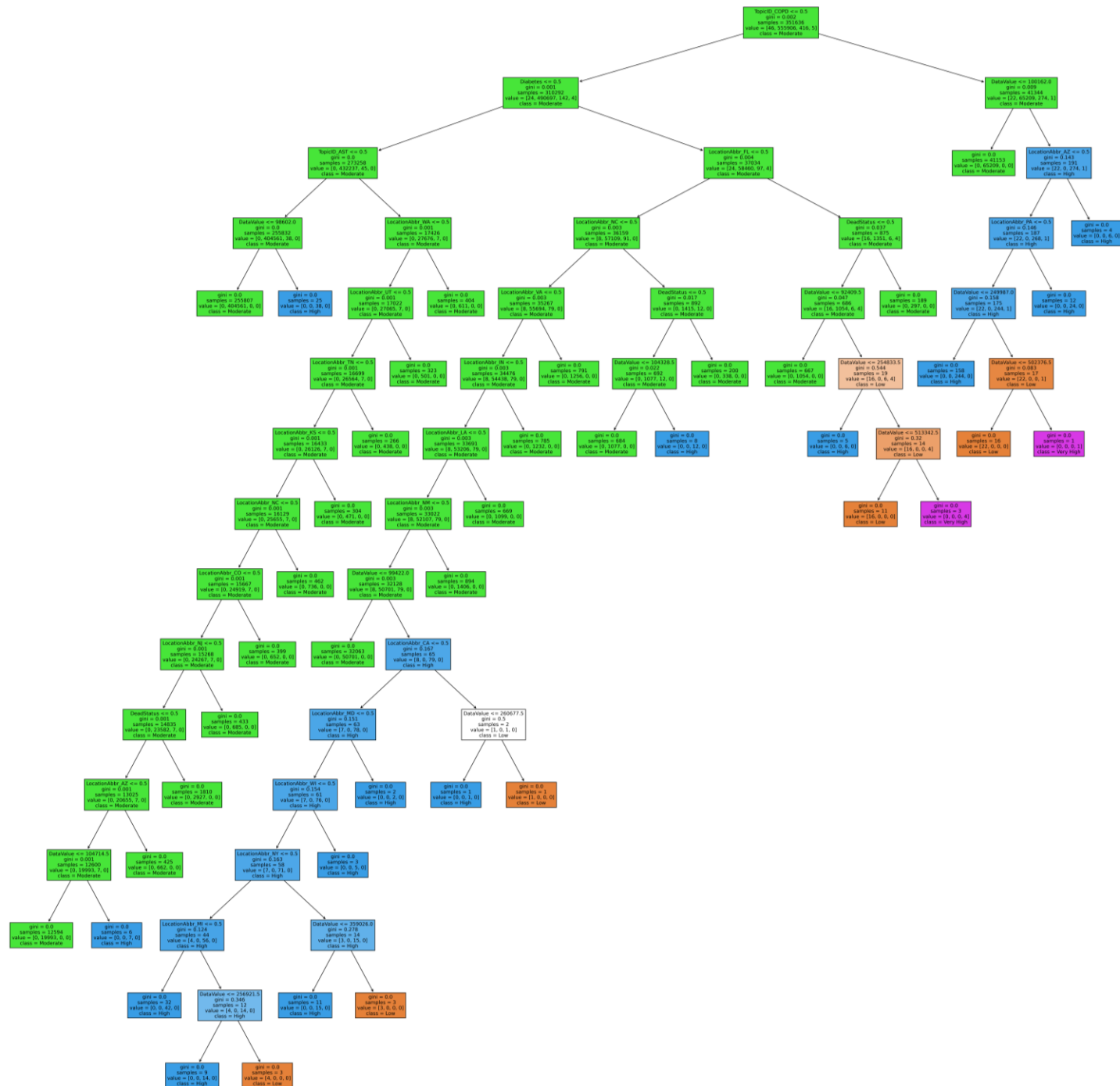
❑ Silhouette score: 0.63622 which is relatively high.



K-Means Clustering of Most Frequent Topic in Each State (Based on Topic Counts)

# Risk Level Assessment

❑ Random Forest Classifier

❑ Attributes: `LocationAbbr`, `YearStart`, `DataValue`, `DeadStatus`, `TopicID`

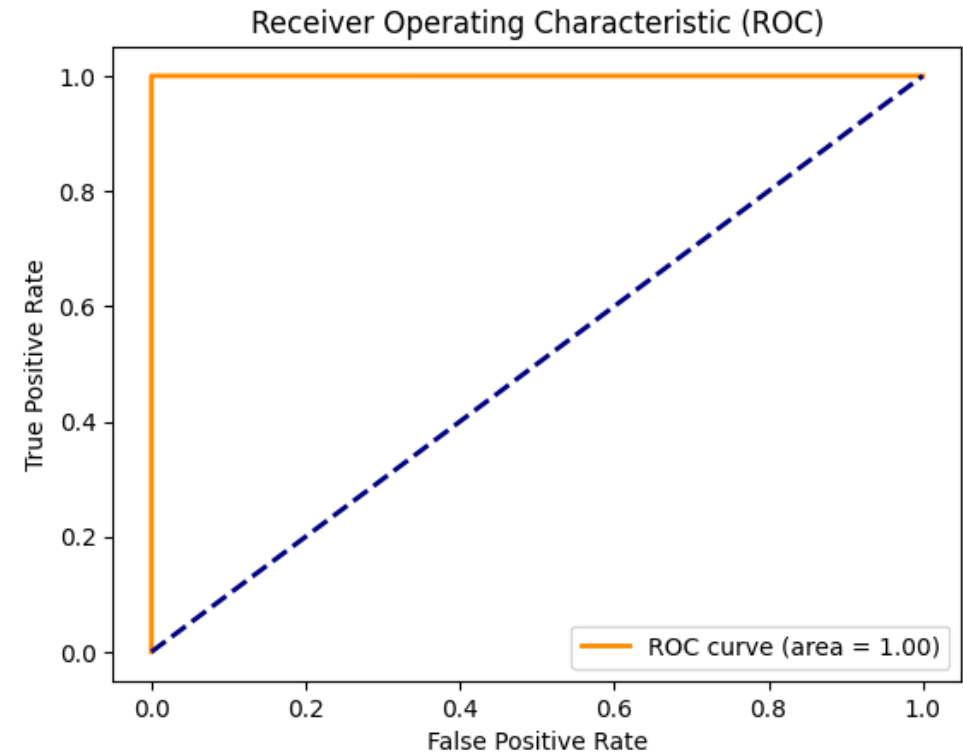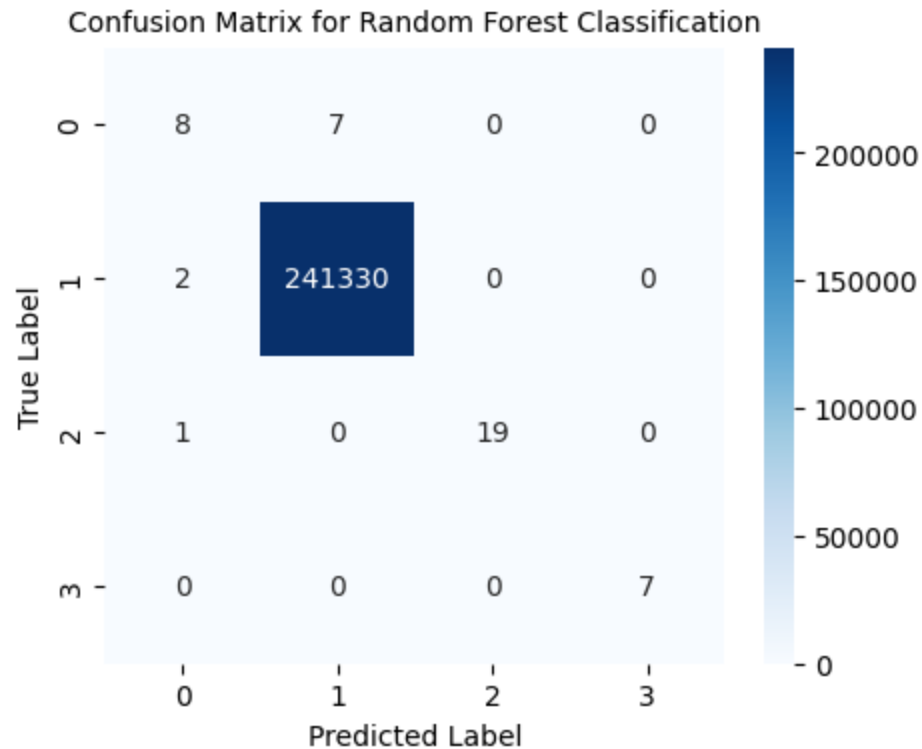❑ Target: `RiskLevel` that is defined based on `DataValue` and `DataValueAlt`.

| Risk Level | DataValue |
|------------|-----------|
| Low | [-1, 500000) |
| Medium | [500000 – 1000000) |
| High | [1000000 – 2000000) |
| Very High | >= 2000000 |

**Structure of Random Forest**

# RandomForest Classifier Evaluation

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| High | 0.93 | 0.93 | 0.93 | 27 |
| Low | 1.00 | 1.00 | 1.00 | 238236 |
| Moderate | 0.98 | 1.00 | 0.99 | 180 |
| Very High | 1.00 | 0.33 | 0.50 | 3 |
| | | | | |
| accuracy | | | 1.00 | 238446 |
| macro avg | 0.98 | 0.81 | 0.85 | 238446 |
| weighted avg | 1.00 | 1.00 | 1.00 | 238446 |



Confusion Matrix for Random Forest Classification



Receiver Operating Characteristic (ROC)

# Chronic Disease Mortality Prediction

❑ Logistic Regression

   ▪ Binary classification tasks. (1 = Deceased, 0 = Alive)

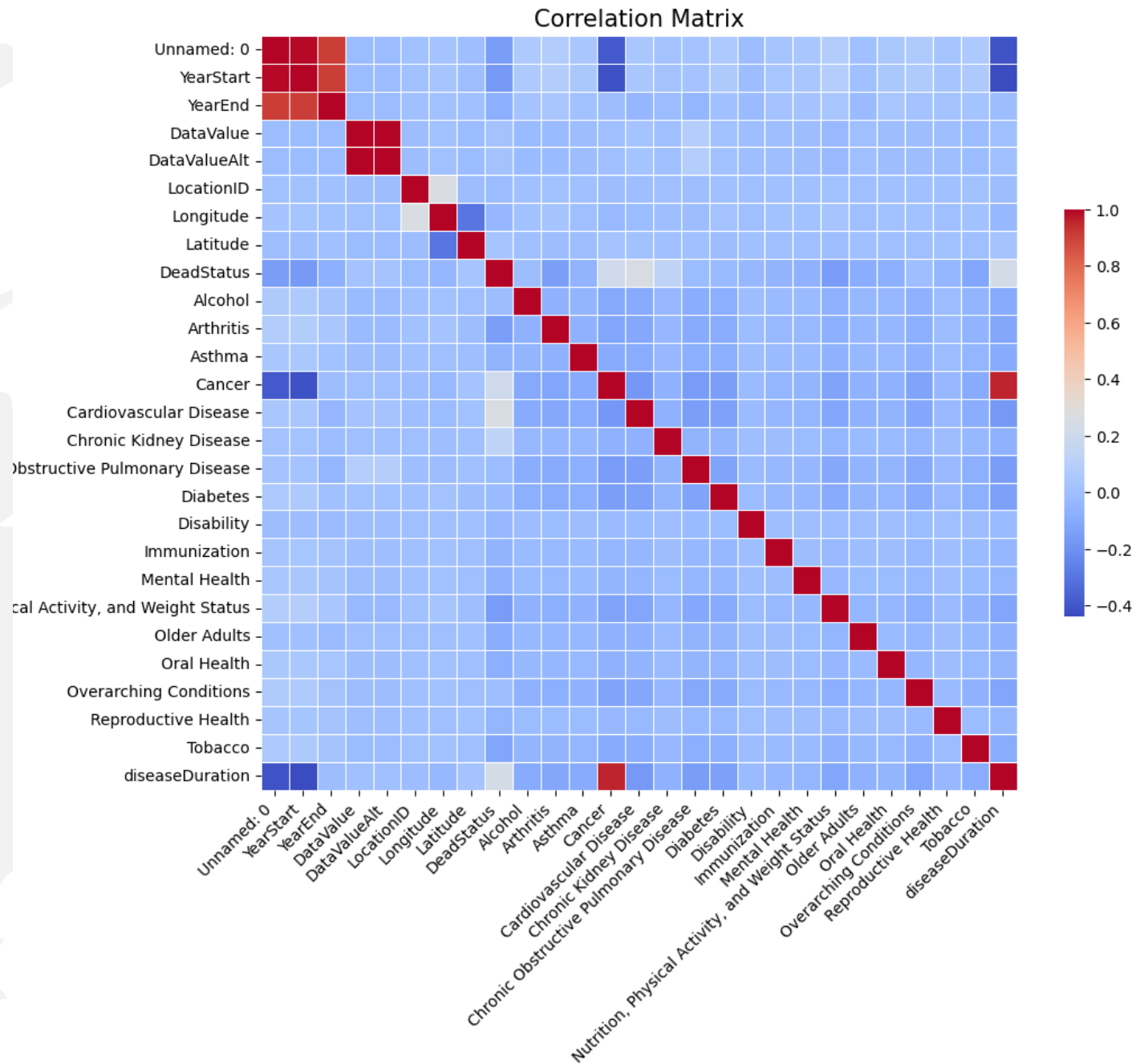❑ Naïve Bayes Classification

   ▪ Able to handle many features

❑ Attributes used for model training

   ▪ `YearStart`, `Arthritis`, `Cancer`, `Cardiovascular Disease`, `Chronic Kidney Disease`, `Nutrition, Physical Activity, and Weight Status`, `Tobacco`, `diseaseDuration`
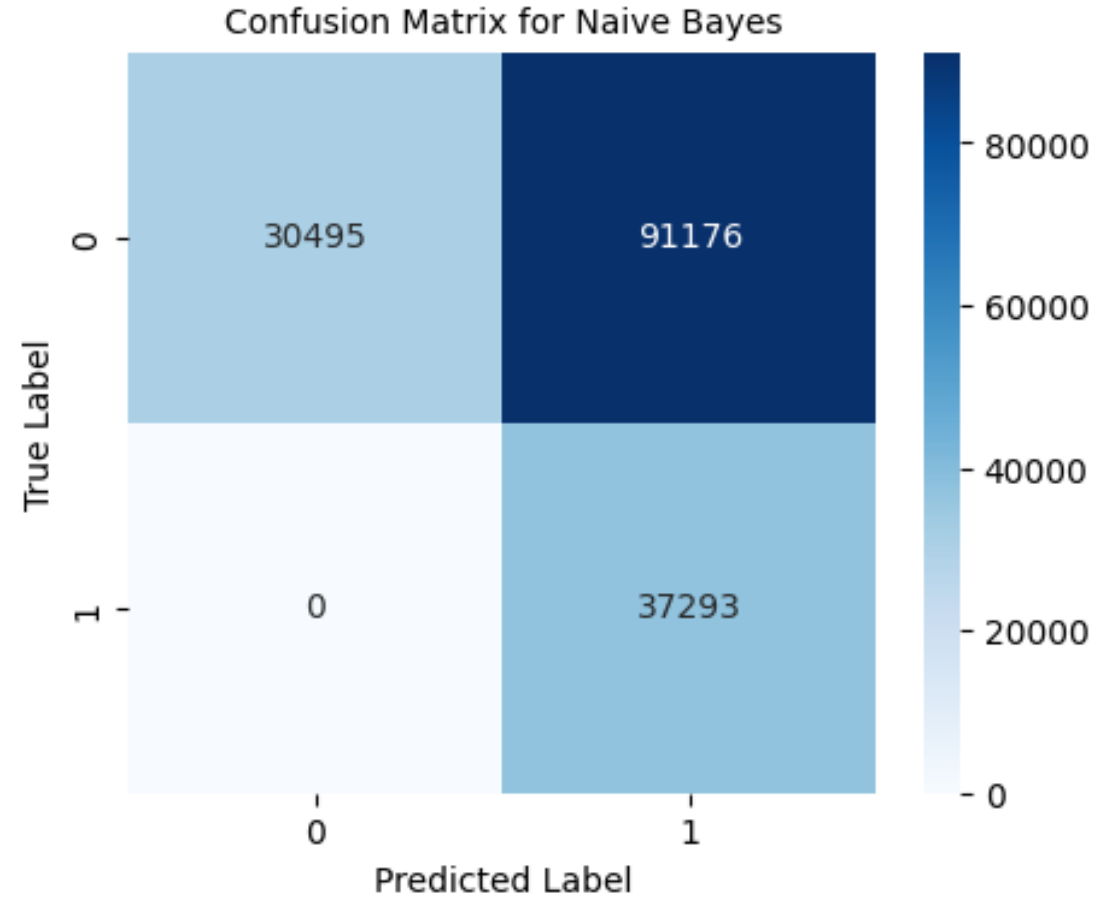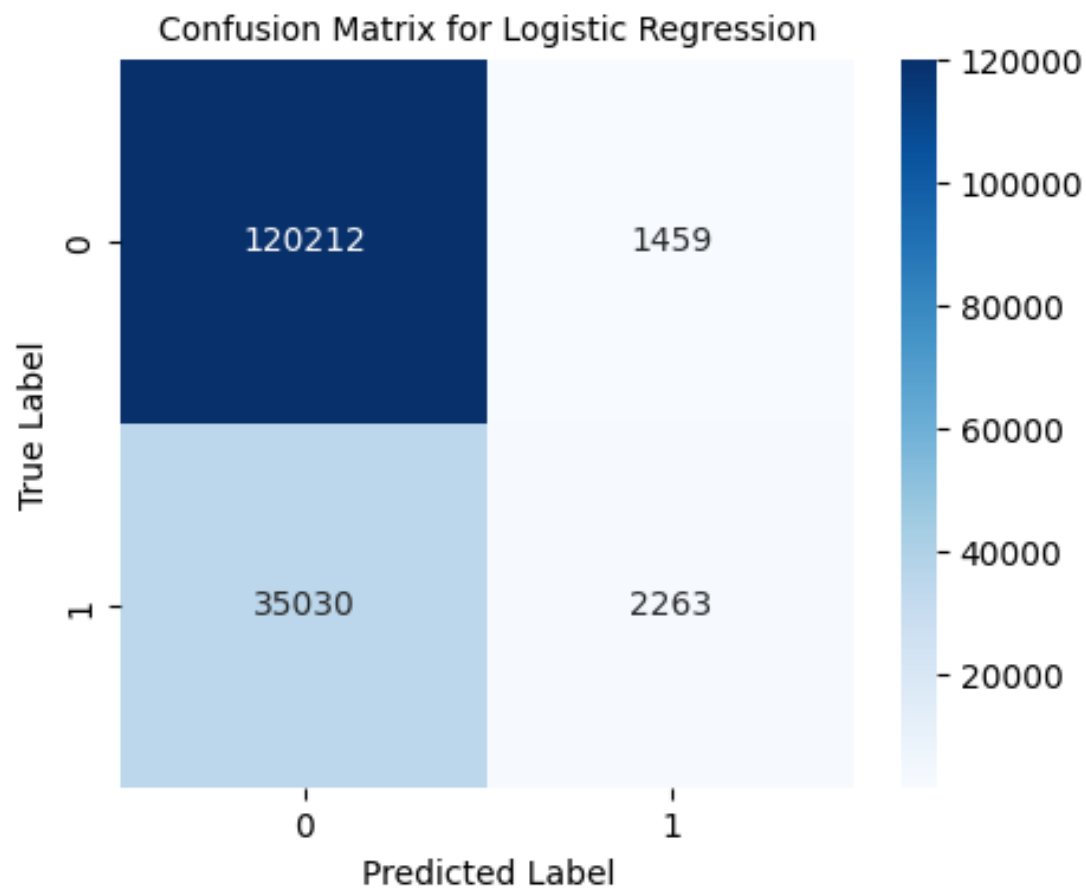
# Correlation Matrix

❑ Selected features based on correlation threshold (> 0.1)

*'YearStart', 'DeadStatus', 'Arthritis', 'Cancer', 'Cardiovascular Disease', 'Chronic Kidney Disease', 'Nutrition, Physical Activity, and Weight Status', 'Tobacco', 'diseaseDuration'*
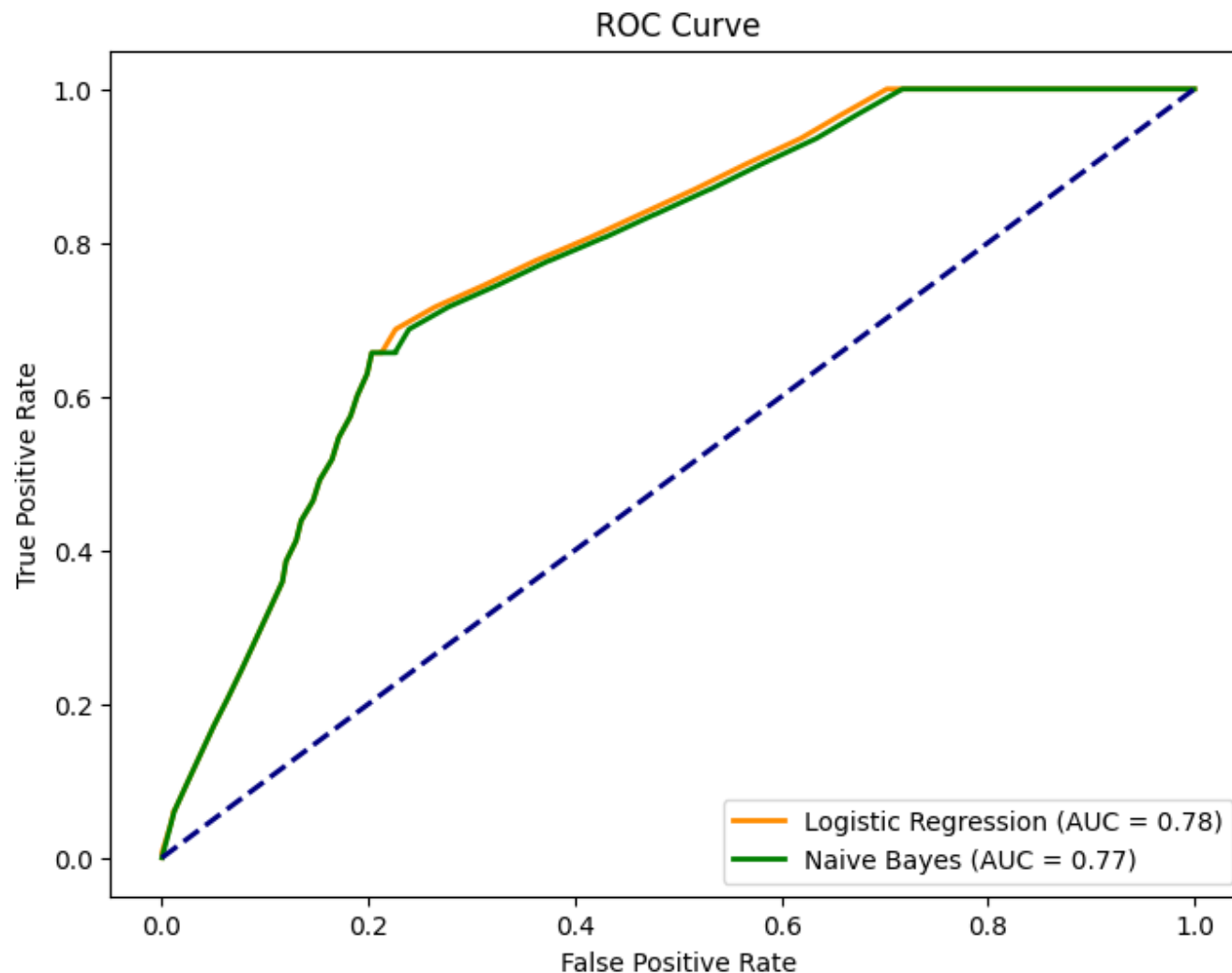
# Logistic Regression vs. Naïve Bayes



Predicted 1, Actual 1 = TP    Predicted 0, Actual 1 = FN

Predicted 1, Actual 0 = FP    Predicted 0, Actual 0 = TN

# Logistic Regression vs. Naïve Bayes Result



Logistic Regression Classification Report

```
Accuracy: 0.77
Classification Report :
              precision    recall  f1-score   support

          0       0.77      0.99      0.87    121671
          1       0.61      0.06      0.11     37293

   accuracy                           0.77    158964
  macro avg       0.69      0.52      0.49    158964
weighted avg       0.74      0.77      0.69    158964
```

Naïve Bayes  Classification Report

```
Accuracy: 0.43
Classification Report:
              precision    recall  f1-score   support

          0       1.00      0.25      0.40    121671
          1       0.29      1.00      0.45     37293

   accuracy                           0.43    158964
  macro avg       0.65      0.63      0.43    158964
weighted avg       0.83      0.43      0.41    158964
```

# Challenges

❑ Unbalanced dataset

❑ Data types are diverse: hard to compile into model training, so we need intensive preprocessing data.

❑ Data was not available for all `Topic` for all years.

# Questions