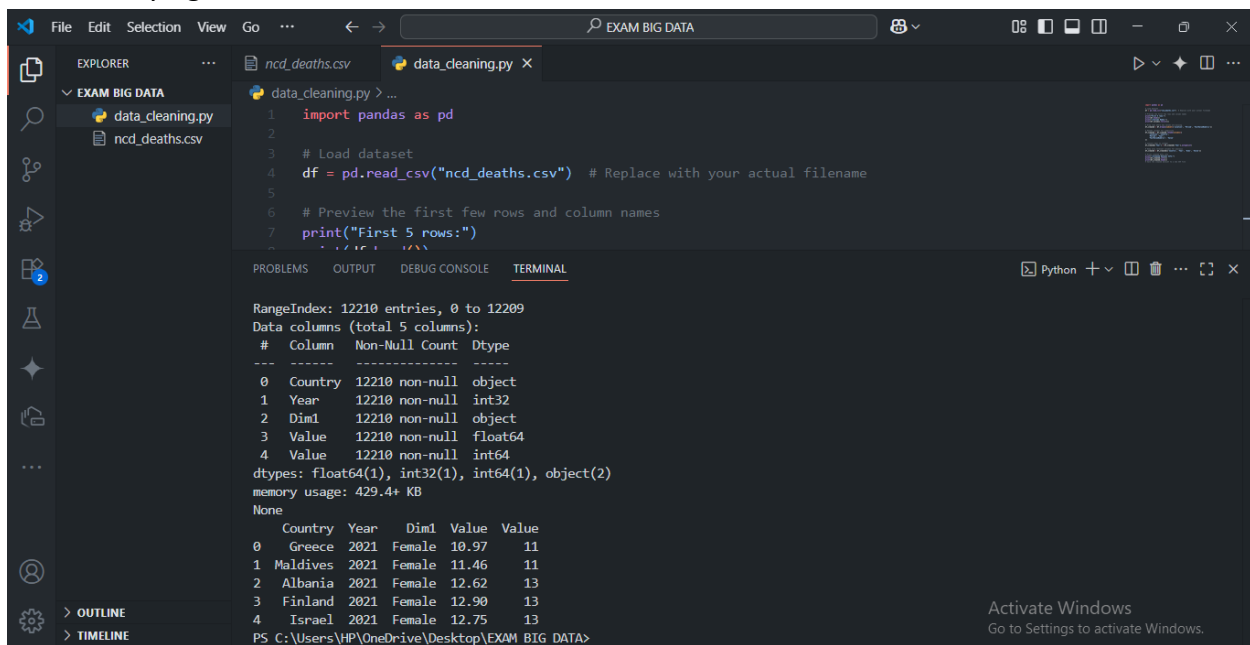Names: MUHIRE Samuel

Id: 26092

INTRODUCTION TO  BIG DATA EXAM SCREEN SHOOT TAKEN


1. Clean the Dataset ▪ Handle missing values, inconsistent formats, and outliers ▪ Apply necessary data transformations (e.g., encoding, scaling)
Data Cleaning

We start by loading the WHO NCD dataset and handling missing values, standardizing formats, and identifying outliers.

Conduct Exploratory Data Analysis (EDA) ▪ Generate descriptive statistics ▪ Visualize distributions and relationships among variables
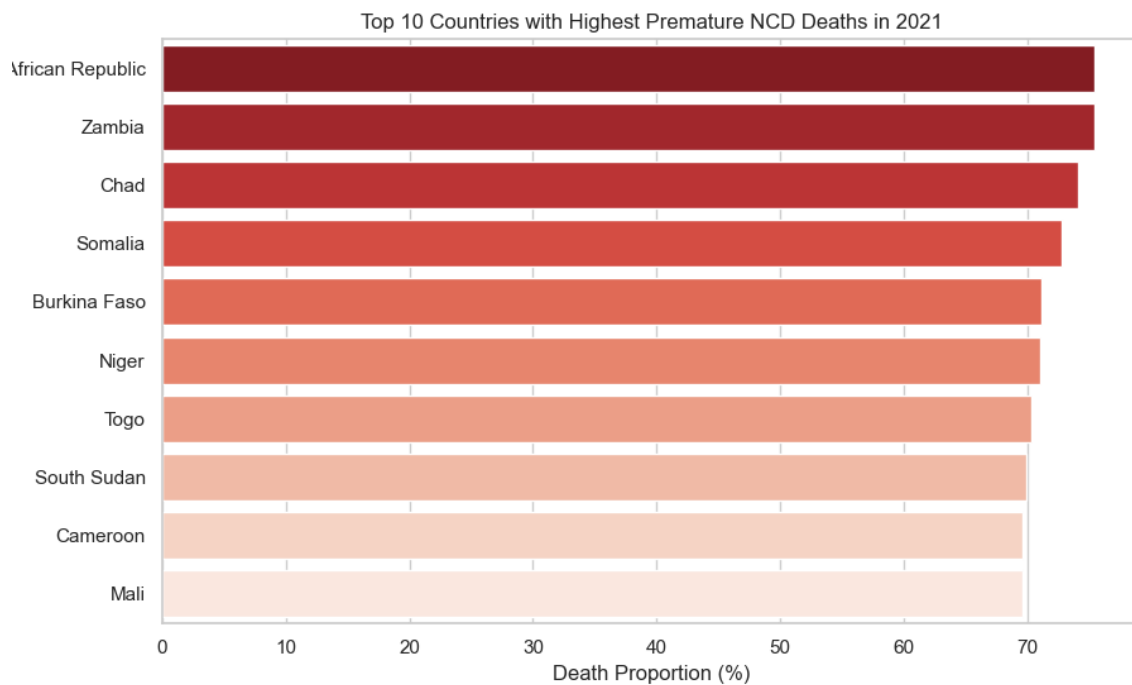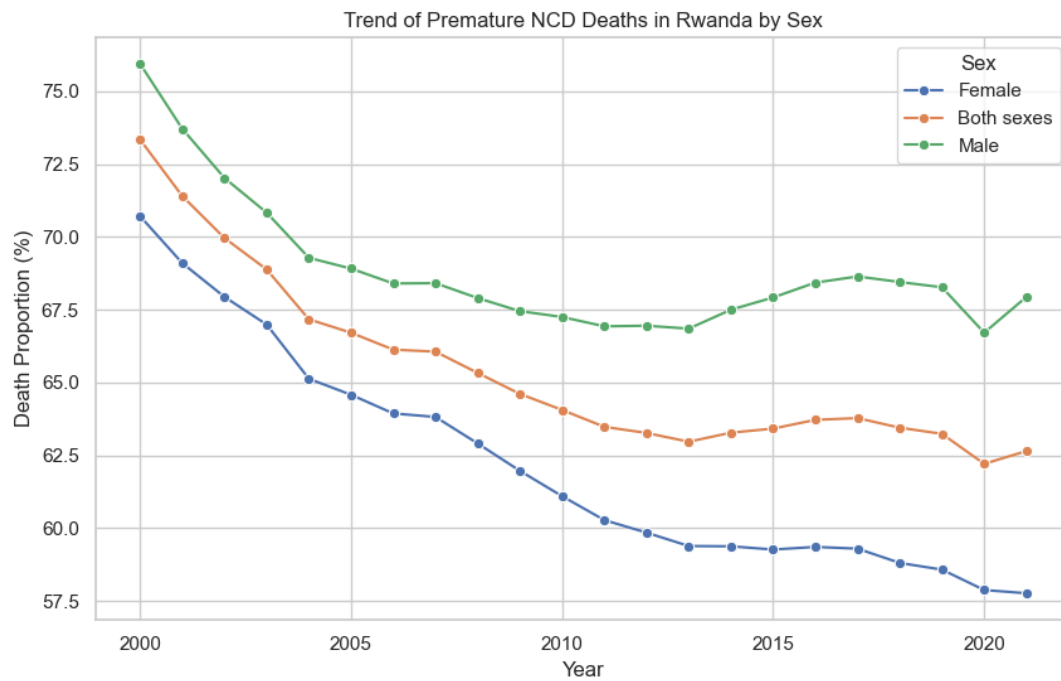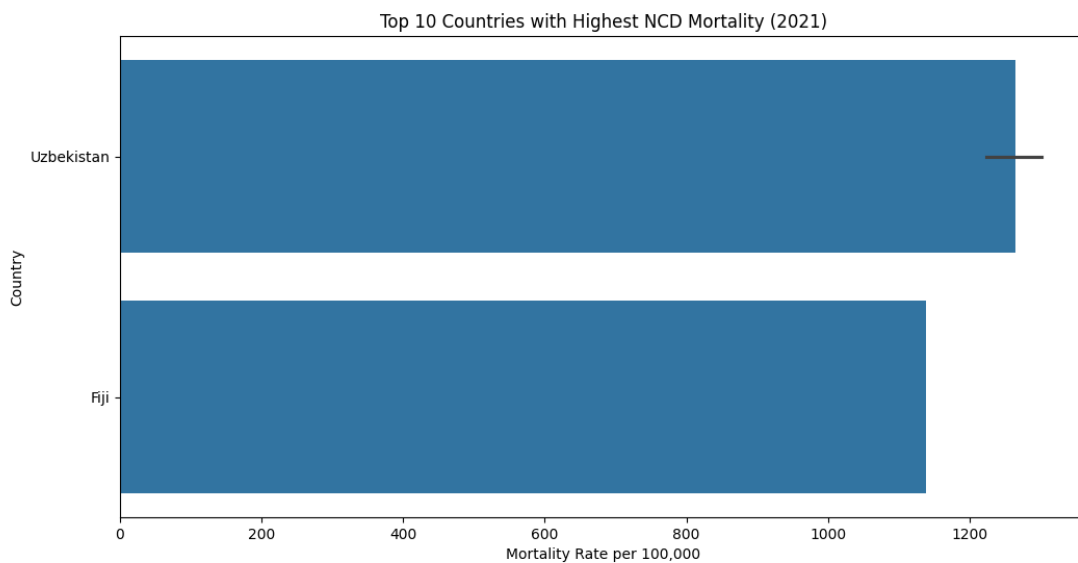
Figure 1

Distribution of Premature NCD Deaths by Sex



Figure 1

Top 10 Countries with Highest Premature NCD Deaths in 2021

Trend of Premature NCD Deaths in Rwanda by Sex

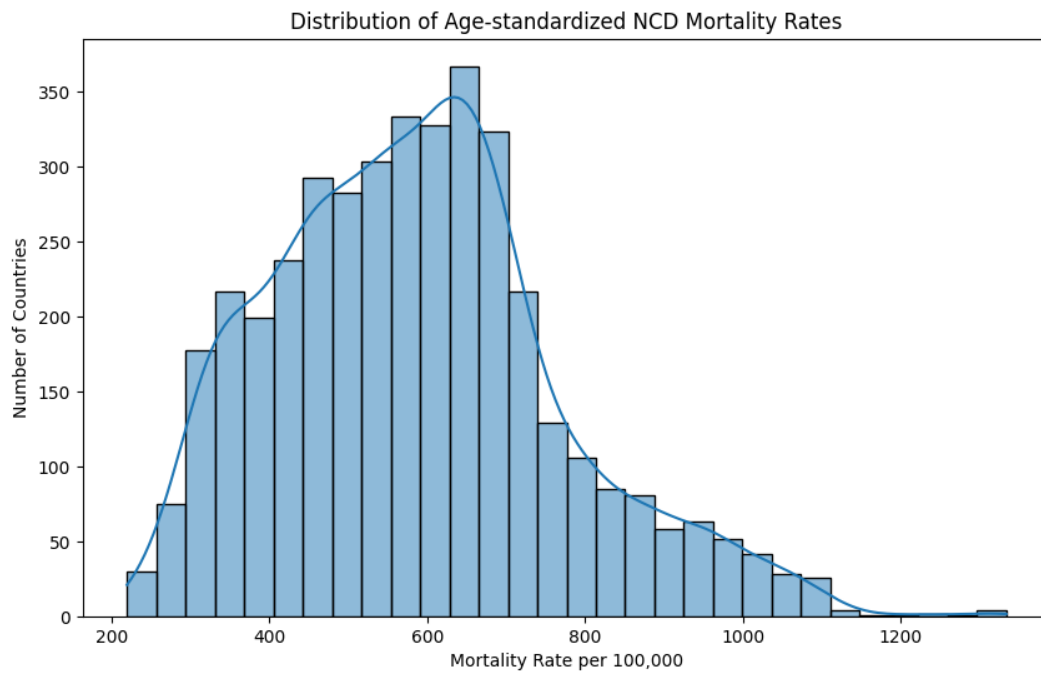**Exploratory Data Analysis (EDA)**

**Objectives:**

- Understand distributions

- Compare countries/regions

- Highlight Rwanda

Distribution of Age-standardized NCD Mortality Rates



Top 10 Countries with Highest NCD Mortality (2021)

### 3. Apply a Machine Learning or Clustering Model

Because this is not a predictive task, we can use **Clustering (e.g., KMeans)** to group countries by mortality levels.

## Evaluate the Model

For clustering, we use **silhouette score**.



## 5. Structure Your Code with Markdown & Functions

Example of reusable function:

6
1. Custom Function: Highlight Countries Exceeding a Threshold



## DBSCAN Clustering for High-Risk Region Detection

Python

DBSCAN Clustering of Countries by Mortality Rate

Time Series Forecasting with ARIMA for Rwanda

Forecast of NCD Mortality in Rwanda