# PROJECT 3

Fate is inevitable

Author: Rick Olsen, Cole Geiger, Ashby Maryo

# INTRODUCTION

The primary goal of this project is to conduct a comprehensive exploration, analysis, and visualization of essential demographic and economic indicators.

By leveraging a diverse set of datasets, we aim to uncover meaningful insights and trends that contribute to a holistic understanding of the socio-economic landscape.

Our project aims to analyze mortality data from mortality.org.

## Area of Focus:
- Birth Rates
- Death Rates
- Income
- Population Statistics

# PLATFORMS AND LIBRARIES

**Python**

Python is the primary programming language used for data analysis and processing.

**Pandas**

Pandas is used for data analysis, manipulation, and transformation.

**SQLAlchemy**

SQLAlchemy is used for database interaction and query generation in Python.

**SQL**

SQL is used for interacting with the database and storing the transformed data.

**Pandera**

Pandera is used for data processing and validation to ensure data quality and consistency.

**Matplotlib**

Matplotlib is used for creating basic visualizations of the analyzed data.

**Seaborn**

Seaborn is used for advanced data visualization and enhancing Matplotlib functions.

# DATASET OVERVIEW

THE ANALYSIS OF MORTALITY DATA INVOLVES THE USE OF SEVERAL DATASETS. THESE DATASETS PROVIDE VALUABLE INFORMATION ON BIRTH RATES, DEATH RATES, AND POPULATION STATISTICS, WHICH ARE CRUCIAL FOR UNDERSTANDING AND ANALYZING MORTALITY TRENDS. THE FOLLOWING DATASETS WERE USED IN THIS ANALYSIS:

## Datasets Used

| Dataset | Description |
|---|---|
| Births_df | Provides information on the number of births per year. |
| Income_df | Information showing Size of Household and Median income |
| Bltper_df | Provides life expectancy and mortality rates by year and age. |
| Deaths _df | Includes the number of deaths by year, age, and gender. |
| Population_df | Contains population statistics by year, age, and gender. |
| Per_capita_df | Provides population and per capita income by year. |

| Year | Female | Male | Total |
|------|--------|------|-------|
| 1933 | 1122180 | 1184820 | 2307000 |
| 1934 | 1166072 | 1229928 | 2396000 |
| 1935 | 1158000 | 1219000 | 2377000 |
| 1936 | 1148000 | 1207000 | 2355000 |
| 1937 | 1175000 | 1238000 | 2413000 |
| 1938 | 1217000 | 1280000 | 2497000 |
| 1939 | 1201000 | 1265000 | 2466000 |
| 1940 | 1246000 | 1313000 | 2559000 |
| 1941 | 1316000 | 1387000 | 2703000 |
| 1942 | 1452000 | 1537000 | 2989000 |
| 1943 | 1510000 | 1593000 | 3103000 |
| 1944 | 1430000 | 1509000 | 2939000 |
| 1945 | 1391000 | 1467000 | 2858000 |
| 1946 | 1657000 | 1754000 | 3411000 |
| 1947 | 1857000 | 1960000 | 3817000 |
| 1948 | 1771000 | 1866000 | 3637000 |
| 1949 | 1777000 | 1872000 | 3649000 |
| 1950 | 1768000 | 1863000 | 3631000 |
| 1951 | 1863000 | 1960000 | 3823000 |
| 1952 | 1908000 | 2005000 | 3913000 |
| 1953 | 1931000 | 2034000 | 3965000 |
| 1954 | 1988000 | 2090000 | 4078000 |
| 1955 | 2001000 | 2103000 | 4104000 |
| 1956 | 2056000 | 2162000 | 4218000 |

# Storing Data in SQL (SQLAlchemy)

TO STORE THE ANALYZED MORTALITY DATA, WE WILL BE USING SQLALCHEMY, A PYTHON LIBRARY THAT PROVIDES A SQL TOOLKIT AND OBJECT-RELATIONAL MAPPING (ORM) TOOLS.

SQLALCHEMY ALLOWS US TO CONNECT TO A SQL DATABASE AND STORE OUR PANDAS DATAFRAMES AS SQL TABLES.

| Connect to SQL Database | → | Create SQL Tables | → | Store Pandas DataFrames |

## Connect to SQL Database

The first step is to establish a connection to the SQL database using SQLAlchemy. This involves specifying the necessary connection details, such as the database type, host, port, username, and password.

## Create SQL Tables

Once the connection is established, we can create SQL tables to store our data. We define the table schema, including the column names, data types, and any constraints.

## Store Pandas DataFrames

With the SQL tables created, we can now store our Pandas DataFrames in the SQL database. SQLAlchemy provides methods to easily insert data into the tables, either row by row or in bulk.

# TRANSFORMATION

## Data Loading, Cleaning and Formatting

- Create primary keys to uniquely identify each record.

- Importing CSV's as DataFrames

- Limit the data to the past 50 years to focus on recent trends.

- Format numerical data to ensure consistency and accuracy.

- Round values to the appropriate decimal places for better readability.

## Variable Renaming and Dropping

- Rename variables to provide clearer and more descriptive names.

- Drop unneeded variables that are not relevant to the analysis.

## Incorporating Additional Health Indicators

- Include additional health indicators to provide a more nuanced analysis.

- Explore correlations between birth rates, death rates, population statistics, and income

- Gain insights into the health impact and median household income.

```python
# define income_df_schema
income_df_schema = pa.DataFrameSchema({
    "Year": pa.Column(int, checks=pa.Check(lambda s: s > 1972)),
    "Number (thousands)": pa.Column(int),
    "Median Income Current dollars": pa.Column(int),
    "Median Income 2022 dollars": pa.Column(int),
    "Mean Income Current dollars": pa.Column(int),
    "Mean Income 2022 dollars": pa.Column(int),
    "Average size of household": pa.Column(float),
})

income_dfv = income_df_schema(income_df)
```

# IMAGES

# DATA DESIGN –DIAGRAM



www.quickdatabasediagrams.com

**annual_df**

| Year | INTEGER |
|---|---|
| Female | INTEGER |
| Male | INTEGER |
| Total | INTEGER |
| Numbers_(thousands) | INTEGER |
| Median_Income_Current_dollars | INTEGER |
| Median_Income_2022_dollars | INTEGER |
| Mean_Income_Current_dollars | INTEGER |
| Mean_Income_2022_dollars | INTEGER |
| Average_size_of_household | FLOAT |
| Population_in_thousands | INTEGER |
| Per_capita_income_Current_dollars | INTEGER |
| Per_capita_income_2022_dollars | INTEGER |

**annual_x_age_df**

| Annual_Death_Rate | FLOAT |
|---|---|
| Probability_of_Death | FLOAT |
| Number_of_Deaths | INTEGER |
| Life_Expectancy | INTEGER |
| ID | STRING |
| Female_Deaths | FLOAT |
| Male_Deaths | FLOAT |
| Total_Deaths | FLOAT |
| Year | INTEGER |
| Age | INTEGER |
| Female_population | FLOAT |
| Male_population | FLOAT |
| Total_population | FLOAT |

# DATA PROCESSING AND VALIDATION (PANDERA)

IN THIS PROJECT, WE WILL BE USING PANDERA, A PYTHON LIBRARY, FOR DATA PROCESSING AND VALIDATION. PANDERA ALLOWS US TO DEFINE AND VALIDATE DATA SCHEMAS, ENSURING DATA QUALITY AND CONSISTENCY THROUGHOUT OUR ANALYSIS.

Data Processing and Validation Steps

| Step | Description |
|---|---|
| 1. Data Cleaning | Remove any duplicate or irrelevant data, handle missing values, and standardize data formats. |
| 2. Data Transformation | Apply necessary transformations to the data, such as aggregating or disaggregating data, creating new variables, or merging datasets. |
| 3. Data Validation | Define and validate data schemas using Pandera to ensure data quality and consistency. Validate data types, ranges, and relationships between variables. |
| 4. Data Quality Assurance | Perform quality checks on the processed data to ensure accuracy, completeness, and integrity. Identify and resolve any data quality issues. |

# CORRELATION HEATMAP AND DISTRIBUTION ANALYSIS

## Seaborn Heatmap

- To unveil relationships between income-related variables, we've utilized Seaborn to create a heatmap. This graphical representation allows for a quick and intuitive assessment of correlations, aiding in identifying patterns and potential areas of interest.



Correlation Matrix

# DATA VISUALIZATION (MATPLOTLIB AND SEABORN)

- Positive correlation between "Annual Death Rate" and "Total Population" would suggest that as the total population increases, the annual death rate tends to increase.

- Negative correlation between "Life Expectancy" and "Annual Death Rate" would suggest that as life expectancy increases, the annual death rate tends to decrease.



Correlation Heatmap - Annual Death Rate, Life Expectancy, Year, Age, Total Population

# CORRELATION HEATMAP AND DISTRIBUTION ANALYSIS

## Histograms with Seaborn:

- The distribution analysis involves using Seaborn to construct histograms for both births and life expectancy. These visualizations provide insights into the frequency and distribution of these events, offering a deeper understanding of the underlying trends.

# RESULTS

## Key Findings

•Birth rates have been steadily declining over the past decade, indicating a decrease in fertility rates.
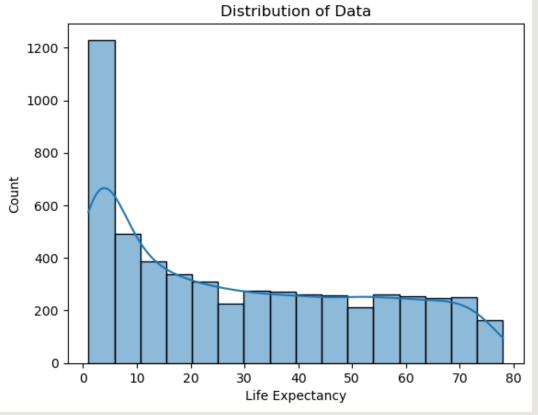
•Death rates have been relatively stable, suggesting improvements in healthcare and medical advancements.

•The population has been aging, with a significant increase in the elderly population.

## Insights

•The decrease in birth rates may have implications for future population growth and workforce dynamics.

•The stable death rates indicate the effectiveness of healthcare systems in maintaining overall health and well-being.

•The aging population presents challenges and opportunities for healthcare, retirement planning, and social support systems.

# Ethics and Privacy

**Data Confidentiality**
- The project places a high priority on data confidentiality, ensuring that sensitive information is protected and only accessible to authorized individuals.

**Responsible Data Handling**
- Data handling practices adhere to ethical standards, including obtaining informed consent, anonymizing data, and securely storing and transmitting data.

**Fairness in Analysis**
- The analysis process is conducted with fairness in mind, avoiding bias and ensuring that all relevant factors are considered.

**Accurate and Responsible Sharing of Findings**
- The project is committed to sharing findings accurately and responsibly, avoiding misinterpretation and providing clear context for the data.

**Privacy, Transparency, and Ethical Standards**
- Privacy, transparency, and ethical standards are respected throughout the project, ensuring that the rights and well-being of individuals are protected.

# SOURCES

**Mortality.org and US Census.gov**

Mortality.org is a comprehensive database that provides access to mortality data from various sources. It includes data on birth rates, death rates, and population statistics from different countries and regions.

**Accessing Data**

To access the data from mortality.org, we utilized web scraping techniques to extract the relevant information. Python libraries such as BeautifulSoup and Requests were used to automate the process and retrieve the data in a structured format.

**Data Analysis**

The extracted data was then transformed and analyzed using Python and Pandas.