

Title: Documentation for East Africa Population Data Analysis

1. **Introduction:** This documentation provides a detailed explanation of the Python script utilized for loading, cleaning, and visualizing a dataset containing population trends in East Africa over several years. The dataset has information ranging from total population, yearly change, median age, fertility rate, urban population among others.

2. **Library Importation:**

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

The necessary libraries for data manipulation and visualization are imported. This includes pandas for data manipulation, numpy for numerical operations, and matplotlib & seaborn for data visualization.

3. **Data Loading:**

```
east_africa = pd.read_csv("150684_original.csv", encoding= 'unicode_escape' )
east_africa
```

The dataset is loaded into a pandas DataFrame using the `pd.read_csv()` function. The `encoding` parameter is set to `'unicode_escape'` to handle any special characters in the dataset.

4. **Initial Data Exploration:**

- Checking the first and last five rows of the dataset to understand the data structure and the types of data contained in the DataFrame.
- Checking the data types and other information like the number of non-null entries in each column using the `info()` method.
- Checking a brief of the statistics using the `describe()` method to understand the distribution of numerical data.
- Listing the column names using the `columns` attribute to understand the field names in the dataset.

5. **Data Cleaning:**

- Dropping an unnecessary column 'Yearly....Change' using the `drop()` method.
- Renaming columns to remove dots and make the column names more readable using the `str.replace()` method.
- Checking the shape of the dataset to understand the number of rows and columns in the DataFrame.
- Checking for unique values in the columns using the `nunique()` method to identify the number of unique values in each column.
- Identifying columns with null values using the `isna().any()` method.

6. **Handling Missing Values:**

- Checking the mode of the column 'EasternAfricaRankwithinAfrica' using the `mode()` method.

- Replacing the null value in 'EasternAfricaRankwithinAfrica' with the mode using the `replace()` method.

7. Data Visualization:

- Creating scatter plots to visualize the relationship between different variables over time, using seaborn's `scatterplot()` method.
- The variables plotted against the 'Year' column are 'Population', 'MedianAge', 'WorldPopulation', and 'UrbanPopulation'.

8. Outlier Detection:

- No outliers were detected in the time series data as per the scatter plots.

9. Data Export:

```
east_africa.to_csv(r'C:\Users\Rick-Royal\Documents\Strathmore University Data  
Science and Analytics\DMSR\150684_cleaned.csv')
```

- The cleaned and processed DataFrame is exported to a CSV file named '150684_cleaned.csv' using the `to_csv()` method.

10. **Conclusion:** The script provided offers a structured approach to loading, cleaning, and visualizing the population dataset of East Africa. Through this script, insights regarding population trends, median age, urban population, and world population in relation to the years can be derived. The script also handles missing data to ensure the dataset is ready for further analysis or modeling.

This documentation provides a step-by-step explanation of the code and its purpose in the data analysis process of East African population data.