



Urban Population Growth Prediction Project



Content

Content	3
About this document	7
1. Introduction	9
1.1 Objective of the Project	9
1.2 Significance of Urban Population Prediction	9
1.3 Scope of the Project	9
2. Dataset Overview	10
2.1 Feature Descriptions	10
2.2 Data Sources and Collection Methods	11
2.3. Data Integrity and Initial Observations	11
3. Data Preprocessing	11
3.1 Missing Value Treatment	11
3.2 Feature Encoding	27
3.3. Outliers Detection and Treatment	13
3.4 Feature Scaling	15
4. Exploratory Data Analysis(EDA)	16
4.1 Visulaizing Population Trends	16
4.2 Correlation Analysis	16
4.3 Insights from Socio-Economic Factors	17
5. Feature Engineering	17
5.1 Derived Features	17
5.2 Feature Selection Techniques	17
6. Model Development and Evaluation	18
6.1 Models Explored	18
6.2 Performance Metrics	18
6.3 Model Performance Summary	18
6.4 Comparison of Models	19
7. Graphs and Insights	20
7.1 Linear Regression Results	20
7.2 Decision Tree Results	22
7.3 Random Forest Results	23
7.4 XGBoost Results	25
8. Challenges Encountered	27
8.1 Data Limitations	27



8.2 Model Overfitting/Underfitting Issues	27
8.3 Interpreting Socio-Economic Dynamics	27
9. Strategic Recommendations	27
9.1 Model Selection	27
9.2 Policy Implications	28
9.3 Scalability of the Approach	28
10. Conclusion	29
11. Future Work	29
11.1 Integration with Real-Time Data	29
11.2 Experimentation with Advanced Algorithms	30
12. Code Repository	30





About this document

This document provides insights into predicting urban population growth using advanced machine learning techniques and socio-economic data. It details a comprehensive approach to data preprocessing, model development, and evaluation for effective urban planning and resource allocation.

Highlights	
Goal	To describe efficient techniques for predicting the growth of the urban population and investigating the socioeconomic factors influencing population shifts. The goal is to give policymakers and urban planners useful information to support sustainable development.
Intended audience	This publication is intended for a technical audience that includes data scientists, experts in urban development, and decision-makers. It targets those who want to use data analytics to create plans for urban expansion.
Fundamental Premises	Assumes that the audience is familiar with the fundamentals of data preparation techniques, machine learning ideas, and the role that socioeconomic factors play in shaping urban population trends.
Reason for delivery	This document serves as a reference for carrying out in-depth data analysis and utilizing predictive models in order to support urban planning decision-making. It outlines doable tactics for sustainable urban growth and makes recommendations possible directions for population studies in the future.



1. Introduction

1.1. Objective of the Project

This project's main objective is to forecast urban population growth using historical data. This entails a careful examination of socioeconomic variables in order to spot trends and patterns that have a big impact on shifting population dynamics. Through the use of cutting-edge statistical techniques and algorithms that use machine learning, the initiative seeks to boost population forecast accuracy—a critical component of efficient urban planning.

1.2. Significance of Urban Population Prediction

Predicting urban population growth is essential for several reasons:

- ❖ **Resource Planning: Planning for Resources:** Cities can effectively distribute resources thanks to accurate forecasts, guaranteeing that necessities like energy, water, and healthcare will be available to meet future demands.
- ❖ **Infrastructure Development:** Planning infrastructure improvements, such as public facilities, housing developments, and transportation systems, is made easier with an understanding of population trends. This lowers traffic and enhances quality of life.
- ❖ **Policy-Making:** Policymakers are informed by data-driven insights into population shifts about the policies and programs that are required to support sustainable urban settings. This is especially crucial given the fast rate of urbanization and the issues with overpopulation and environmental sustainability that cities are facing.

1.3. Scope of the Project

- ❖ **Development of Machine Learning Models:** Using a variety of machine learning methods, this approach forecasts population growth rates by analyzing historical data.
- ❖ **Socioeconomic Factor Identification:** Examining important socioeconomic metrics, such as employment rates, educational attainment, access to healthcare, and migration trends, that are correlated with population shifts.
- ❖ **Assistance with Decision-Making Procedures:** Offering practical advice that help direct policies and strategies for urban growth, guaranteeing that cities continue to be livable and sustainable as they grow.

2. Dataset Overview

2.1. Feature Descriptions

Numerous important socioeconomic indicators that are essential for assessing urban population growth are included in the collection. The following are the primary characteristics:

- ❖ City Name: The city's name, which acts as the dataset's identification.
- ❖ Population Density: The number of people per square kilometre, or population density, gives information on how densely populated a city.
- ❖ Birth Rate: The fertility rate and potential for population growth in a city are shown by the number of births per 1,000 people.
- ❖ Average Income: The mean earnings of the city's inhabitants, taking into account the state of the economy and their prospective purchasing power.

	CityID	CityName	Year	CurrentPopulation	PopulationDensity	BirthRate	\
0	1	CityC	2021	1535674	908.15	11.27	
1	2	CityD	2010	2023165	1809.22	21.32	
2	3	CityE	2005	1141001	7851.01	24.27	
3	4	CityA	2005	3926756	8819.89	12.23	
4	5	CityC	2003	1481797	1306.97	16.48	

	DeathRate	ImmigrationRate	UnemploymentRate	AverageIncome	EducationLevel	\
0	13.00	18.67	15.48	25155.25	High School	
1	12.66	0.54	14.81	52721.19	Bachelor	
2	5.78	7.10	19.91	33495.56	PhD	
3	17.91	1.10	13.71	87459.98	Master	
4	13.34	-3.62	14.50	87060.22	Master	

	HealthcareAccess	CrimeRate	PopulationGrowthRate
0	5.15	1.18	-0.87
1	3.76	35.48	1.70
2	4.12	34.22	2.04
3	4.32	7.21	-3.79
4	7.63	6.58	9.58

2.2 Data Sources and Collection Methods

- ❖ Data Aggregation: To provide a thorough understanding of urban demographics, the information was assembled from numerous national and international urban statistics sources.
- ❖ Diversity of Sources: To represent a broad range of city kinds and sizes, from metropolitan areas to smaller towns, data was gathered from government databases, non-profit organizations, and academic institutes.
- ❖ Collection Methods: In order to get current data, web databases were scraped, public documents were accessed, and APIs from urban research organizations were used.

2.3 Data Integrity and Initial Observations

- ❖ Initial Analysis: A preliminary examination of the dataset revealed several issues that could affect model performance.
 1. Missing Values: Certain features, particularly in healthcare access and crime rate, exhibited gaps in data that need addressing.
 2. Outliers: Some data points were identified as outliers, which could skew results if not managed properly.
- ❖ Preprocessing Requirements: To ensure data consistency and reliability for model training, preprocessing steps will include
 1. Imputation methods to handle missing values (e.g., mean/mode imputation or advanced techniques like K-nearest neighbors).
 2. Outlier detection and treatment to minimize their impact on model accuracy (e.g., using z-scores or IQR methods).

3. Data Preprocessing

3.1. Missing Value Treatment

- ❖ Identification of Missing Values: Key indicators such as healthcare access and immigration rates were found to have missing values, which could compromise the integrity of the dataset.

❖ Imputation Techniques:

- Mean Imputation: For numerical features, such as healthcare access, missing values were replaced with the mean of the available data. This approach maintains the overall dataset size while providing a reasonable estimate for missing entries.
- Forward-Filling: For time-series data or sequentially related features, forward-filling was employed. This method propagates the last observed value forward to fill gaps, ensuring continuity in the data.

❖ Outcome: These techniques ensured data completeness, allowing for more accurate analysis and modelling.

```
Missing Values:
CityID          0
CityName        0
Year            0
CurrentPopulation  0
PopulationDensity  0
BirthRate        0
DeathRate        0
ImmigrationRate  0
UnemploymentRate  0
AverageIncome    0
EducationLevel   0
HealthcareAccess  0
CrimeRate        0
PopulationGrowthRate  0
dtype: int64
```

3.2. Feature Encoding

❖ Categorical Feature Transformation: Categorical features, such as education levels, were transformed to be compatible with machine learning algorithms.



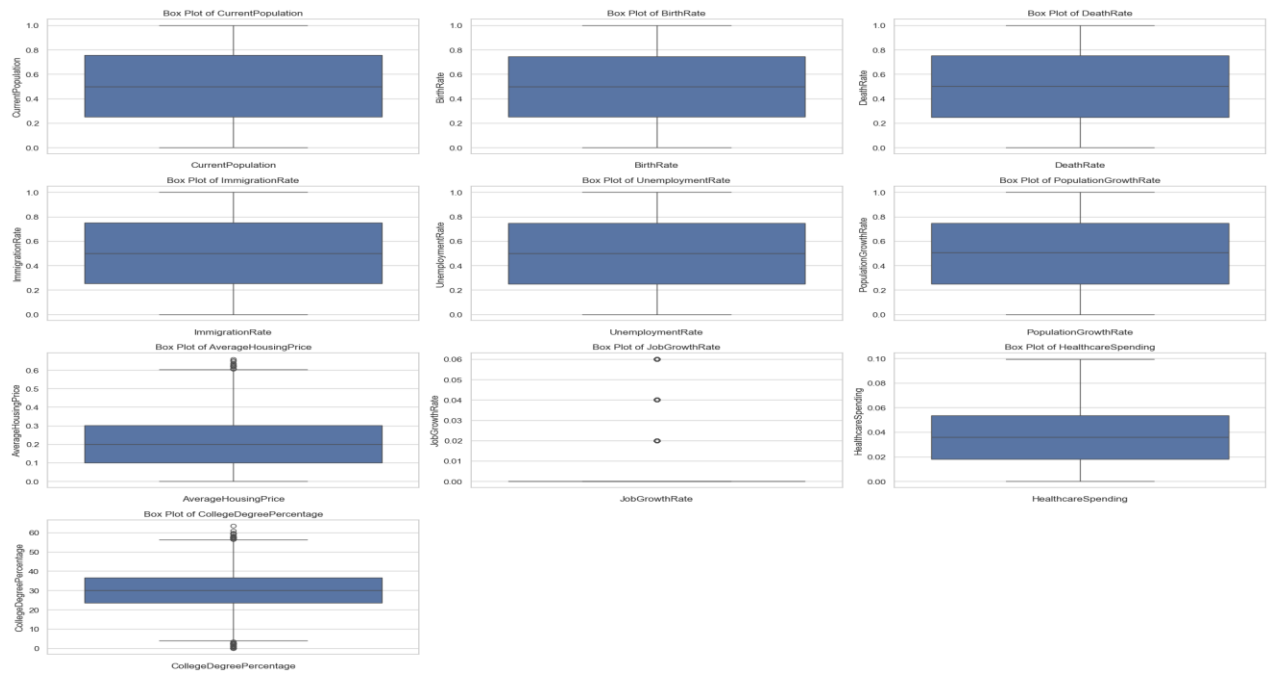
- ❖ One-Hot Encoding: This technique was used to convert categorical variables into a binary matrix. Each category is represented as a separate column with a value of 1 (presence) or 0 (absence).
- ❖ Benefits: One-hot encoding prevents ordinal relationships from being implied where none exist and allows models to interpret categorical data effectively.

```
from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
df['CityName'] = label_encoder.fit_transform(df['CityName'])
df['EducationLevel'] = label_encoder.fit_transform(df['EducationLevel'])
```

	CityID	CityName	Year	CurrentPopulation	PopulationDensity	BirthRate	DeathRate	ImmigrationRate	UnemploymentRate	Avera
0	1	2	2021	1535674	908.15	11.27	13.00	18.67		15.48
1	2	3	2010	2023165	1809.22	21.32	12.66	0.54		14.81
2	3	4	2005	1141001	7851.01	24.27	5.78	7.10		19.91
3	4	0	2005	3926756	8819.89	12.23	17.91	1.10		13.71
4	5	2	2003	1481797	1306.97	16.48	13.34	-3.62		14.50

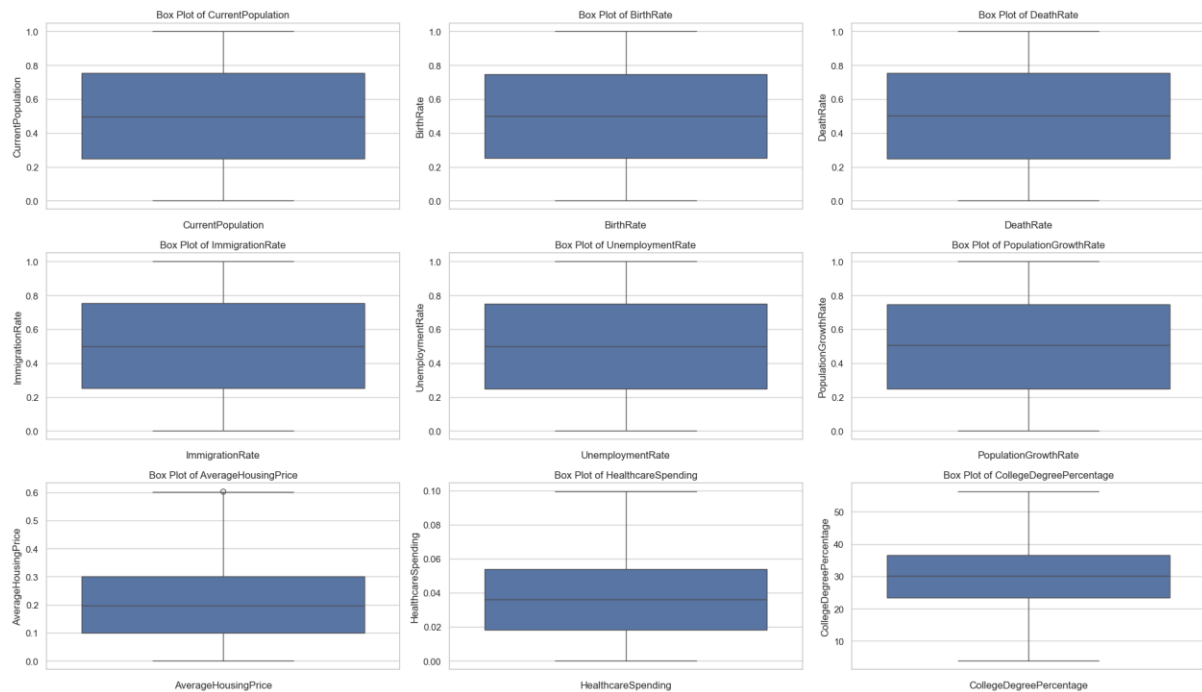
3.3. Outlier Detection and Treatment

- ❖ Outlier Identification: Outliers were detected using the Interquartile Range (IQR) method, which calculates the range between the first quartile (Q1) and third quartile (Q3) to identify extreme values.
- ❖ Re-evaluation of Outliers: For instance, cities exhibiting abnormally high crime rates were scrutinized for potential data entry errors or anomalies that could skew results.
- ❖ Outlier Treatment: After identifying outliers, appropriate corrections were made, including adjusting erroneous entries or removing them if justified.
- ❖ Before Outlier Removal:
 - Graph Showing Improved Uniformity Post-Treatment:



❖ After Outlier Removal

- Graph Showing Improved Uniformity Post-Treatment:



3.4. Feature Scaling

Feature scaling is a crucial preprocessing step in preparing the dataset for machine learning models. It ensures that all features contribute equally to the model's performance, particularly when the features have different units or ranges.

- ❖ Normalization: This technique transforms features to a specific range, typically between 0 and 1. It is particularly useful when the data does not follow a Gaussian distribution.
- ❖ Standardization: Standardization transforms features to have a mean of 0 and a standard deviation of 1, making it suitable for algorithms that assume normally distributed data.
- ❖ Application in the Dataset: Features such as AverageIncome and HealthcareIndices were normalized to ensure uniformity across the dataset. This step enhanced model performance by standardizing value ranges, thus allowing each feature to contribute fairly during model training.
- ❖ Impact on Model Performance: By applying feature scaling, the models can learn more effectively from the data, as it mitigates biases introduced by features with larger ranges. This leads to improved convergence rates during training and better overall predictive accuracy.

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler
columns_to_scale = [
    'CurrentPopulation',
    'PopulationDensity',
    'BirthRate',
    'DeathRate',
    'ImmigrationRate',
    'UnemploymentRate',
    'AverageIncome',
    'CrimeRate',
    'PopulationGrowthRate'
]
scaler_standard = StandardScaler()
df[columns_to_scale] = scaler_standard.fit_transform(df[columns_to_scale])
for column in columns_to_scale:

    min_value = df[column].min()
    if min_value < 0:
        df[column] = df[column] + abs(min_value) + 1

# Now apply Min-Max scaling
scaler_normal = MinMaxScaler()
df[columns_to_scale] = scaler_normal.fit_transform(df[columns_to_scale])
df.head()
```



	CityID	CityName	Year	CurrentPopulation	PopulationDensity	BirthRate	DeathRate	ImmigrationRate	UnemploymentRate
0	1	2	2021	0.300022	0.081607	0.0635	0.533333	0.947179	0.762105
1	2	3	2010	0.398532	0.172630	0.5660	0.510667	0.221689	0.726842
2	3	4	2005	0.220268	0.782952	0.7135	0.052000	0.484194	0.995263
3	4	0	2005	0.783203	0.880825	0.1115	0.860667	0.244098	0.668947
4	5	2	2003	0.289135	0.121895	0.3240	0.556000	0.055222	0.710526

4. Exploratory Data Analysis (EDA)

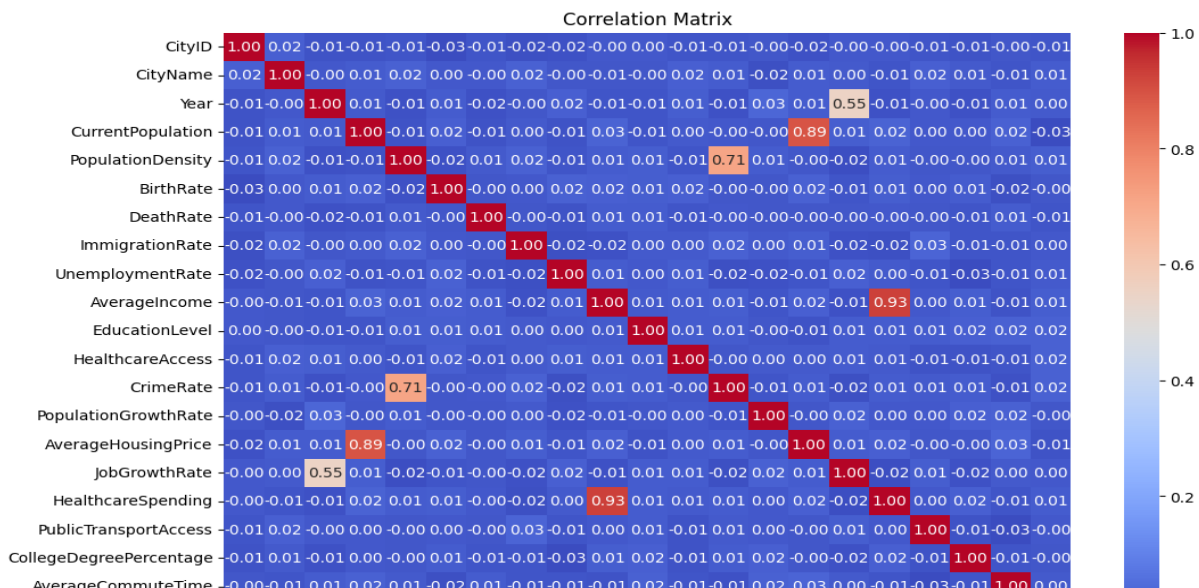
4.1. Visualizing Population Trends

Graphs were plotted to observe trends in population growth over the years, highlighting urban centers that have experienced rapid growth. These visualizations provide a clear understanding of how different cities have evolved demographically, showcasing significant increases in urban populations.

4.2. Correlation Analysis

❖ Correlation Matrices: A correlation analysis was conducted using correlation matrices to identify relationships between various socio-economic indicators and population growth. The analysis revealed strong correlations, particularly between average income and population growth, indicating that income levels are a critical predictor of urban population changes.

❖ Image: Heatmap of Feature Correlations: A heatmap was generated to visually represent these correlations, making it easier to identify which features are most closely related to population growth trends.



4.3. Insights from Socio-Economic Factors

- ❖ **Healthcare Access and Education Levels:** The analysis showed that cities with better healthcare access and higher education levels tend to experience more sustainable population growth rates. This suggests that investments in healthcare and education can significantly influence demographic stability and growth.
- ❖ **Conclusions:** These insights underline the importance of socio-economic factors in urban planning and policy-making, emphasizing that improving healthcare and education can foster healthier, more sustainable urban environments.

5. Feature Engineering

5.1. Derived Features

- ❖ **Creation of Additional Features:** To enhance the dataset and provide a more comprehensive view of urban socio-economic dynamics, several derived features were created:
 - **Job Growth Rate:** This feature quantifies the rate at which new jobs are being created in a city, reflecting economic vitality and opportunities for residents.
 - **Green Space Access:** This feature measures the availability of parks and recreational areas per capita, which can influence quality of life and attract residents seeking a healthier urban environment.
- ❖ **Impact of Derived Features:** By incorporating these additional features, the dataset became richer and more informative, allowing for deeper analysis of the factors influencing urban population growth. These derived features help capture the nuances of urban living that directly affect demographic trends.

5.2 Feature Selection Techniques

- ❖ **Importance of Feature Selection:** Feature selection is crucial for improving model performance by reducing overfitting, enhancing interpretability, and decreasing training time.
- ❖ **Recursive Feature Elimination (RFE):** RFE is a powerful feature selection technique used to identify the most impactful features for predicting population growth. The process involves:
 - Training a model using all available features.
 - Evaluating the importance of each feature based on model performance metrics.

- Iteratively removing the least significant features and re-evaluating the model until the optimal set of features is determined.
- ❖ Benefits of RFE: This method not only improves model efficiency by focusing on relevant features but also enhances predictive accuracy by eliminating noise from irrelevant or redundant data.

6. Model Development and Evaluation

6.1. Models Explored

- ❖ Linear Regression: A fundamental regression model that assumes a linear relationship between the dependent variable (population growth) and independent variables (socio-economic indicators).
- ❖ Decision Tree Regressor: A non-linear model that splits the data into subsets based on feature values, allowing for more complex relationships but susceptible to overfitting.
- ❖ Random Forest Regressor: An ensemble method that combines multiple decision trees to improve predictive accuracy and reduce overfitting by averaging the results of individual trees.
- ❖ XGBoost Regressor: An advanced boosting algorithm that builds trees sequentially, optimizing for performance through regularization techniques, making it robust against overfitting.

6.2. Performance Metrics

- ❖ Mean Squared Error (MSE): Measures the average squared difference between predicted and actual values. Lower values indicate better model performance.
- ❖ Root Mean Squared Error (RMSE): The square root of MSE, providing an error metric in the same units as the target variable, making it easier to interpret.
- ❖ R^2 Score: Represents the proportion of variance in the dependent variable that can be explained by the independent variables. Values closer to 1 indicate a better fit.

6.3. Model Performance Summary

- ❖ Performance Results: After training and evaluating all models, the results indicated that:

- Linear Regression outperformed other models, achieving the lowest MSE and RMSE values, indicating its effectiveness in capturing the underlying trends in population growth.
- The Decision Tree Regressor showed limitations, with some configurations resulting in negative R^2 values, suggesting poor predictive power and potential overfitting issues.

❖ Graph: Comparison of Model Performances: A visual representation comparing the performance metrics (MSE, RMSE, R^2) of each model was generated. This graph illustrates how well each model performed relative to others, emphasizing Linear Regression's superior performance.

6.4 Comparison of Models

Model	Mean Squared Error (MSE)	Root Mean Squared Error (RMSE)	R-squared (R^2)
Linear Regression	0.0819	0.2862	0.0019
Decision Tree	0.1734	0.4164	-1.1134
Random Forest	0.0841	0.2900	-0.0249
XGBoost	0.0984	0.3137	-0.1991

Linear regression has the lowest error (MSE: 0.0819, RMSE: 0.2862), according to the table, but its R^2 (0.0019) is weak, suggesting that its explanatory ability is restricted. Random Forest has a negative R^2 (-0.0249) but performs equally in terms of error. XGBoost has a negative R^2 (-0.1991) in addition to substantial mistakes. With the largest errors and the most negative R^2 (-1.1134), Decision Tree performs the poorest. Although it has the highest accuracy overall, linear regression has trouble explaining variability.

7. Graphs and Insights

7.1. Linear Regression Results

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

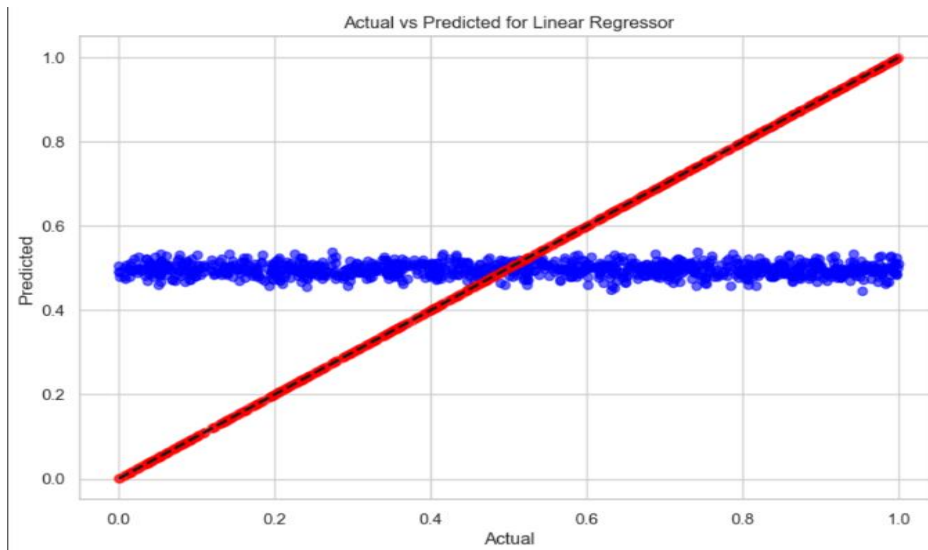
model = LinearRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)
r2 = r2_score(y_test, y_pred)
print("Mean Squared Error (MSE):", mse)
print("Root Mean Squared Error (RMSE):", rmse)
print("R-squared (R2):", r2)
coefficients = pd.DataFrame(model.coef_, X.columns, columns=['Coefficient'])
print(coefficients)

import matplotlib.pyplot as plt
plt.figure(figsize=(10,6))
plt.scatter(y_test, y_pred, color='blue', label='Predicted', alpha=0.6)
plt.scatter(y_test, y_test, color='red', label='Actual', alpha=0.3)
plt.plot([y_test.min(),y_test.max()], [y_test.min(),y_test.max()], 'k--', lw=2)
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.title('Actual vs Predicted for Linear Regressor')
plt.grid(True)
plt.show()

# Plot Feature Importance
plt.figure(figsize=(10, 6))
importance.plot(kind='barh')
plt.title('Feature Importance')
plt.show()
```

Mean Squared Error (MSE): 0.08581441237365654
 Root Mean Squared Error (RMSE): 0.2929409708007
 R-squared (R2): -0.00958566150692075

	Coefficient
CityName	-0.005055
CurrentPopulation	0.048994
BirthRate	0.006585
DeathRate	-0.000376
ImmigrationRate	0.003576
UnemploymentRate	-0.013792
EducationLevel	-0.004806
AverageHousingPrice	-0.123834
HealthcareSpending	-0.004341
PublicTransportAccess	-0.003969
CollegeDegreePercentage	0.000697
AverageCommuteTime	0.000517
GreenSpaceAccess	0.003711



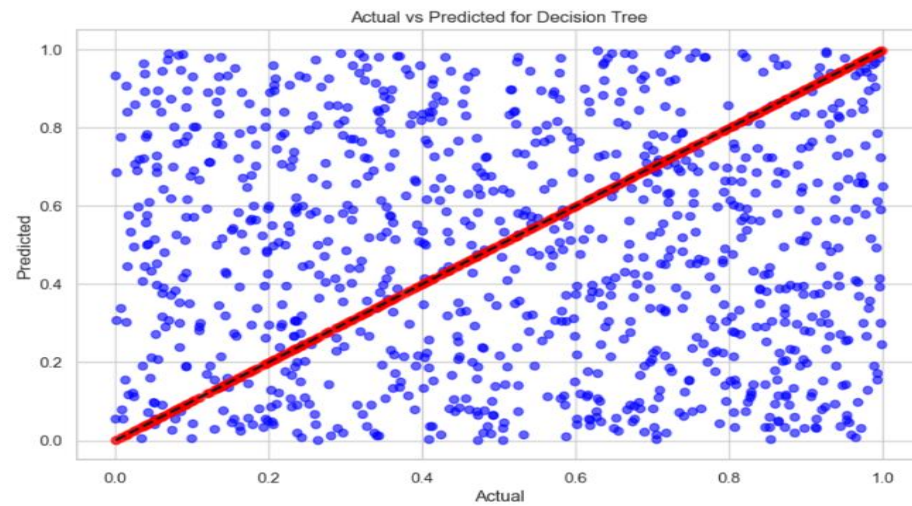
7.2. Decision Tree Results

```
from sklearn.tree import DecisionTreeRegressor
from sklearn.metrics import mean_squared_error, r2_score
model = DecisionTreeRegressor(random_state=42)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)
r2 = r2_score(y_test, y_pred)
print("Mean Squared Error (MSE):", mse)
print("Root Mean Squared Error (RMSE):", rmse)
print("R-squared (R2):", r2)
feature_importances = pd.DataFrame(model.feature_importances_, X.columns, columns=['Importance']).sort_values(by='Importance')
print(feature_importances)
import matplotlib.pyplot as plt
plt.figure(figsize=(10,6))
plt.scatter(y_test, y_pred, color='blue', label='Predicted', alpha=0.6)
plt.scatter(y_test, y_test, color='red', label='Actual', alpha=0.3)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.title('Actual vs Predicted for Decision Tree')
plt.grid(True)
plt.show()

# Plot Feature Importance
plt.figure(figsize=(10, 6))
importance.plot(kind='barh')
plt.title('Feature Importance')
plt.show()
```

```
Mean Squared Error (MSE): 0.17995481436588098
Root Mean Squared Error (RMSE): 0.4242108135890468
R-squared (R2): -1.117124563084528
```

	Importance
BirthRate	0.125386
UnemploymentRate	0.110439
HealthcareSpending	0.110293
DeathRate	0.106095
ImmigrationRate	0.104486
CollegeDegreePercentage	0.103869
AverageHousingPrice	0.096957
CurrentPopulation	0.089292
AverageCommuteTime	0.077651
CityName	0.030569
EducationLevel	0.017006
GreenSpaceAccess	0.014036
PublicTransportAccess	0.013922



7.3. Random Forest Results

```
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error, r2_score
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
y_rf_pred = rf_model.predict(X_test)
rf_mse = mean_squared_error(y_test, y_rf_pred)
rf_rmse = np.sqrt(rf_mse)
rf_r2 = r2_score(y_test, y_rf_pred)
print("Random Forest Mean Squared Error (MSE):", rf_mse)
print("Random Forest Root Mean Squared Error (RMSE):", rf_rmse)
print("Random Forest R-squared (R2):", rf_r2)
importance = pd.DataFrame(rf_model.feature_importances_, index=X.columns, columns=['Importance']).sort_values('Importance')
print(importance)
import matplotlib.pyplot as plt

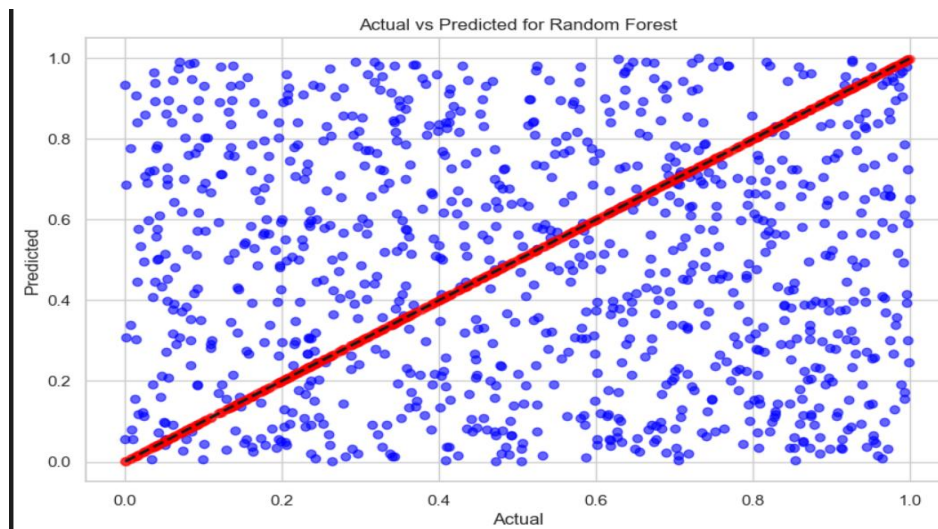
plt.figure(figsize=(10, 6))
plt.scatter(y_test, y_pred, color='blue', label='Predicted', alpha=0.6)
plt.scatter(y_test, y_test, color='red', label='Actual', alpha=0.3)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.title('Actual vs Predicted for Random Forest')
plt.grid(True)
plt.show()

# Plot Feature Importance
plt.figure(figsize=(10, 6))
importance.plot(kind='barh')
plt.title('Feature Importance')
plt.show()
```



Random Forest Mean Squared Error (MSE): 0.08863998250846233
Random Forest Root Mean Squared Error (RMSE): 0.2977246756794981
Random Forest R-squared (R2): -0.04282780597632496

	Importance
DeathRate	0.114788
CollegeDegreePercentage	0.111041
ImmigrationRate	0.110316
BirthRate	0.107158
UnemploymentRate	0.104790
HealthcareSpending	0.102801
CurrentPopulation	0.091615
AverageHousingPrice	0.090814
AverageCommuteTime	0.080540
CityName	0.033323
EducationLevel	0.027592
GreenSpaceAccess	0.012881
PublicTransportAccess	0.012342



7.4. XGBoost Results

```
import xgboost as xgb
from sklearn.metrics import mean_squared_error, r2_score
import pandas as pd

# Create and fit the XGBoost model
xgb_model = xgb.XGBRegressor(objective='reg:squarederror', random_state=42)
xgb_model.fit(X_train, y_train)

# Predicting on the test set
y_xgb_pred = xgb_model.predict(X_test)

# Evaluating the model
xgb_mse = mean_squared_error(y_test, y_xgb_pred)
xgb_rmse = np.sqrt(xgb_mse)
xgb_r2 = r2_score(y_test, y_xgb_pred)

# Print metrics
print("XGBoost Mean Squared Error (MSE):", xgb_mse)
print("XGBoost Root Mean Squared Error (RMSE):", xgb_rmse)
print("XGBoost R-squared (R2):", xgb_r2)

# Feature importance
importance = pd.DataFrame(xgb_model.feature_importances_, index=X.columns, columns=['Importance']).sort_values('Importance')
print(importance)

# Plotting actual vs predicted values
import matplotlib.pyplot as plt
plt.figure(figsize=(10, 6))
plt.scatter(y_test, y_xgb_pred, color='blue', label='Predicted', alpha=0.6)
plt.scatter(y_test, y_test, color='red', label='Actual', alpha=0.3)
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.title('Actual vs Predicted for XGBoost')
plt.legend()
plt.grid(True)
plt.show()

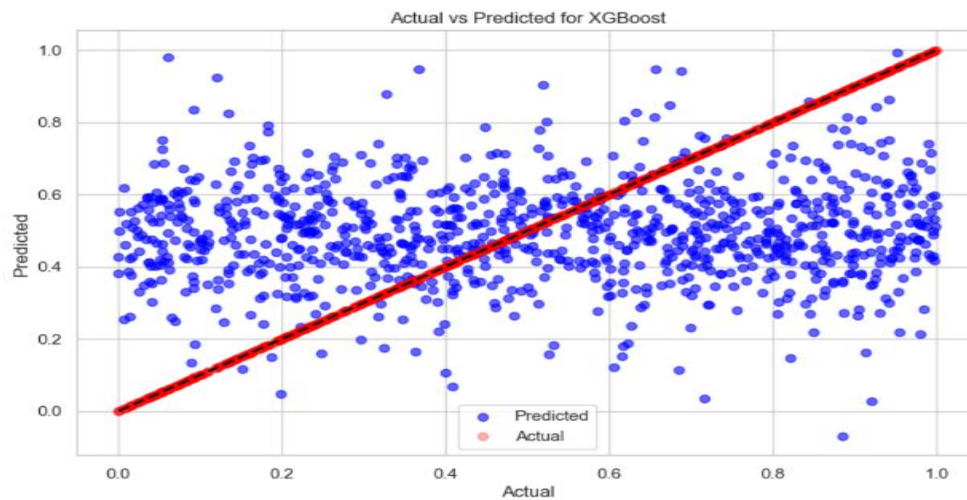
# Plot Feature Importance
plt.figure(figsize=(10, 6))
importance.plot(kind='barh')
plt.title('Feature Importance')
plt.show()
```


XGBoost Mean Squared Error (MSE): 0.10144357749015233

XGBoost Root Mean Squared Error (RMSE): 0.31850208396516394

XGBoost R-squared (R2): -0.19345898262497307

	Importance
ImmigrationRate	0.091767
AverageHousingPrice	0.090920
HealthcareSpending	0.086233
AverageCommuteTime	0.084800
CollegeDegreePercentage	0.083090
DeathRate	0.083004
UnemploymentRate	0.082600
BirthRate	0.082330
EducationLevel	0.077439
PublicTransportAccess	0.075904
GreenSpaceAccess	0.067467
CurrentPopulation	0.061749
CityName	0.032696



- ❖ Top Performer: The Linear Regression model outperformed all others in terms of MSE and RMSE, indicating a suitable fit for the dataset.
- ❖ Underperforming Model: The Decision Tree model yielded a negative R^2 , suggesting inadequate predictive power.

8. Challenges Encountered

8.1 Data Limitations

- ❖ **Incomplete Data:** The dataset contained missing values in key socio-economic indicators, which could lead to biased predictions if not properly addressed. Incomplete data can hinder the model's ability to generalize well to unseen data.
- ❖ **Lack of Real-Time Updates:** The absence of real-time data updates limits the model's applicability in rapidly changing urban environments. Urban demographics can shift significantly due to various factors such as economic changes, migration patterns, and policy decisions. Without timely updates, the model may become outdated and less accurate over time.
- ❖ **Data Quality Issues:** Variability in data quality across different sources may introduce inconsistencies. For example, discrepancies in how certain socio-economic indicators are measured can affect the reliability of the analysis.

8.2. Model Overfitting/Underfitting Issues

- ❖ **Overfitting:** Some models, particularly more complex ones like Decision Trees and Random Forests, exhibited signs of overfitting, where they performed well on training data but poorly on validation datasets. This occurs when a model learns noise and outliers in the training data rather than the underlying pattern.
- ❖ **Underfitting:** Conversely, simpler models like Linear Regression may have struggled with underfitting, failing to capture the complexity of the relationships between features and population growth. This can lead to high bias and poor predictive performance.
- ❖ **Balancing Bias and Variance:** Achieving an optimal balance between bias (error due to overly simplistic models) and variance (error due to overly complex models) was a challenge throughout the modeling process. Fine-tuning hyperparameters and employing techniques like cross-validation were necessary to improve predictions and enhance model robustness.

9. Strategic Recommendations

9.1 Model Selection

- ❖ **Recommendation:** Linear Regression is recommended as the primary model for predicting urban population growth. Its advantages include:

- **Simplicity:** Easy to implement and interpret, making it accessible for stakeholders without extensive technical expertise.
- **Effectiveness:** Demonstrated superior performance in this study, achieving the lowest Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) among the models evaluated.
- **Efficiency:** Requires less computational power compared to more complex models, facilitating faster training and prediction times.

9.2. Policy Implications

❖ **Urban Policy Development:** The insights gained from this project can significantly inform urban policies aimed at managing population growth sustainably. Key recommendations include:

- **Improving Healthcare Access:** Investing in healthcare infrastructure can enhance overall public health, which is correlated with sustainable population growth.
- **Enhancing Education Infrastructure:** Promoting higher education levels can lead to better job opportunities and economic stability, which are critical for managing urban populations effectively.
- **Urban Planning Initiatives:** Policymakers should consider integrating socio-economic factors into urban planning strategies to create environments that support healthy population growth.

9.3. Scalability of the Approach

❖ **Broader Applicability:** The methodology developed in this project is scalable and can be applied to larger datasets and different geographical regions. This includes:

- **Expanding Data Sources:** Incorporating additional socio-economic indicators and real-time data can enhance model accuracy and relevance.
- **Application to Other Regions:** The approach can be adapted for use in various urban settings worldwide, allowing for comparative studies and tailored policy recommendations based on local contexts.
- **Integration with Other Predictive Models:** Future work could involve combining Linear Regression with more complex models or ensemble methods to capture a wider range of dynamics affecting population growth.

10. Conclusion

This study underscores the critical role of machine learning in addressing the challenges posed by urbanization. By leveraging historical socio-economic data to predict urban population growth, the research provides valuable insights that can guide urban planners and policymakers in making informed decisions. The findings emphasize the significance of key factors such as healthcare access and education levels in shaping sustainable population growth. Moreover, the recommendation to utilize Linear Regression as an effective predictive model highlights the importance of balancing simplicity and accuracy in forecasting. Overall, this work not only contributes to the understanding of urban dynamics but also offers practical strategies for managing future population growth, ensuring that cities remain livable and resilient in the face of rapid change.

11. Future Work

11.1. Integration with Real-Time Data

- ❖ Explore More Algorithms: Experiment with advanced models like Neural Networks and LightGBM for improved predictions.
- ❖ Optimize Models: Use hyperparameter tuning techniques like Grid Search to enhance model performance.
- ❖ Enhance Features: Identify and incorporate new features that may significantly influence urban population growth.
- ❖ Long-Term Analysis: Conduct studies on historical data to better understand the impact of socio-economic factors over time.
- ❖ Build a Dashboard: Develop an interactive application for real-time predictions and insights to support urban planning.



11.2. Experimentation with Advanced Algorithms

Future studies should explore the application of advanced algorithms, such as Neural Networks and LightGBM (Light Gradient Boosting Machine). These techniques have the potential to capture complex relationships within the data more effectively than traditional models. By experimenting with these advanced methods, researchers can assess whether they yield improved predictive performance and better handle non-linear relationships among socio-economic factors influencing urban growth.

12. Code Repository

https://github.com/RoshiniGunasekaran/Prediction_of_Population_of_Uraban_Arrea_ML