

Pakistan Sub-Division Population Analysis & Prediction

1. Problem Definition

The purpose of this project is to analyze the population distribution across Pakistan's sub-divisions and predict whether a given sub-division is urban-dominated or rural-dominated. This can support government planning, resource allocation, and development projects by identifying urbanizing areas.

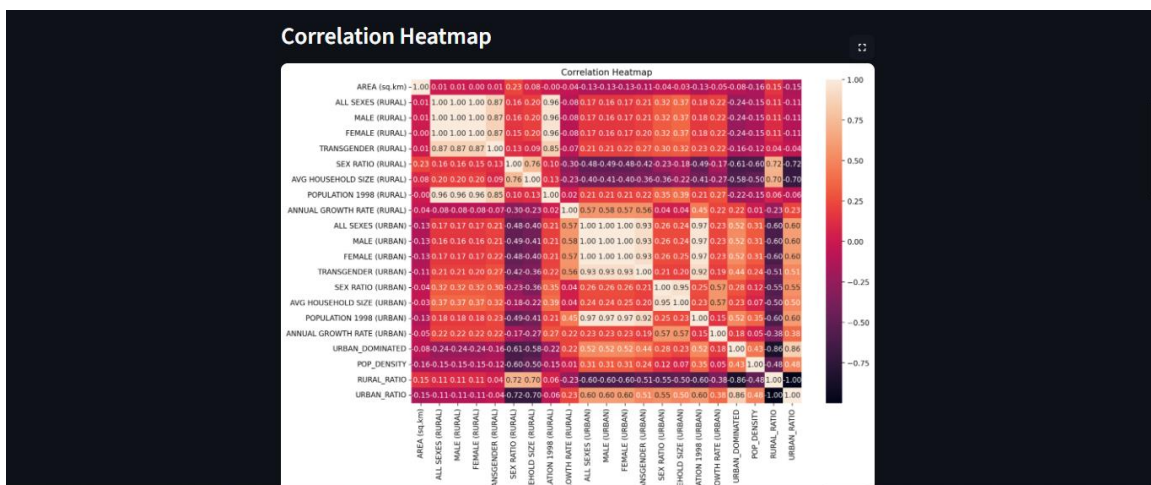
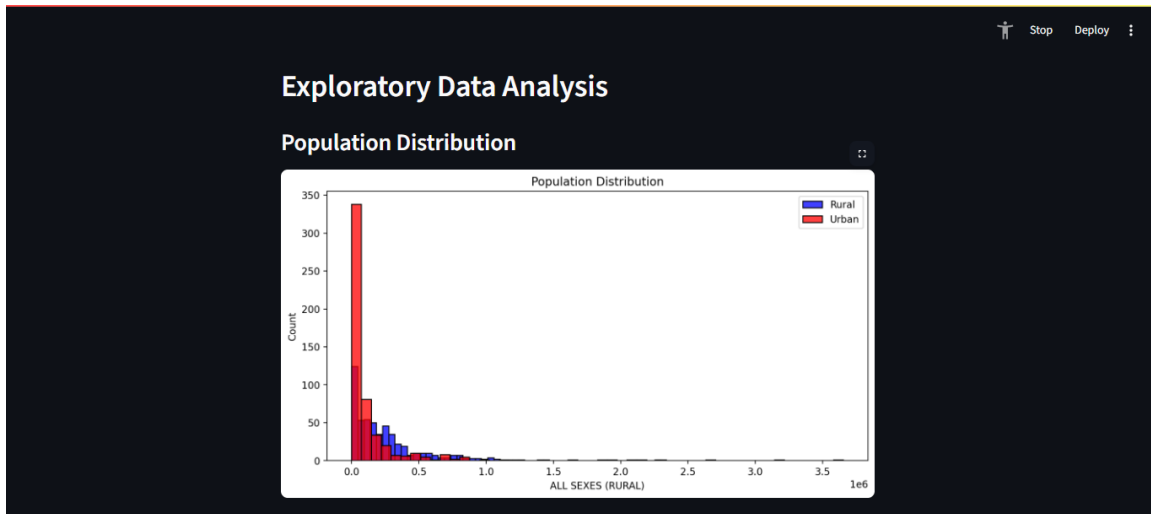
2. Data Preprocessing

The dataset contains demographic and geographic attributes of sub-divisions in Pakistan, including population figures for rural and urban areas, area size, sex ratio, and household sizes. The following steps were taken during preprocessing:

- Removed leading/trailing spaces from column names.
- Replaced missing values with 0.
- Engineered new features:
 - * URBAN_DOMINATED: Binary flag if urban > rural population.
 - * POP_DENSITY: Total population divided by area.
 - * RURAL_RATIO and URBAN_RATIO: Proportional values of population distribution.

3. Exploratory Data Analysis (EDA)

Visual analysis was performed using histograms and heatmaps. Histograms showed the distributions of rural and urban populations, while correlation heatmaps highlighted relationships between variables such as area, population, sex ratios, and derived ratios.



4. Model Building

Two models were trained to predict whether a sub-division is urban or rural dominated:

- Logistic Regression: A linear model used on standardized data.
- Random Forest Classifier: An ensemble method trained on original feature values.

Features used:

- Area, Rural and Urban Populations
- Sex Ratios (Rural & Urban)
- Avg Household Sizes (Rural & Urban)
- POP_DENSITY, RURAL_RATIO, URBAN_RATIO

5. Model Evaluation

Classification reports were generated for both models, including precision, recall, and F1-score. Additionally, GridSearchCV was used to tune hyperparameters for the Random Forest model, optimizing for accuracy. Best parameters and cross-validation score were displayed in the Streamlit app.

Deploy ⋮

Logistic Regression Classification Report

	precision	recall	f1-score	support			
0	0.98	1.00	0.99	89			
1	1.00	0.88	0.94	17			
accuracy				0.98	106		
macro avg				0.99	0.94	0.96	106
weighted avg				0.98	0.98	0.98	106

Random Forest Classification Report

	precision	recall	f1-score	support			
0	0.98	1.00	0.99	89			
1	1.00	0.88	0.94	17			
accuracy				0.98	106		
macro avg				0.99	0.94	0.96	106
weighted avg				0.98	0.98	0.98	106

Random Forest Hyperparameter Tuning (GridSearchCV)	
Best Random Forest Params:	
<pre>{ "max_depth": 5 "n_estimators": 50 }</pre>	
Best CV Score: 1.0	

Predict Urban or Rural Dominance for a Sub-Division ☺☺

Area (sq.km)

1000.00

-

+

Rural Population

50000.00

-

+

Urban Population

50000.00

-

+

Sex Ratio (Rural)

100.00

-

+

Sex Ratio (Urban)

100.00

-

+

The screenshot shows a Streamlit web application interface with a dark theme. At the top right, there is a 'Deploy' button and a menu icon. The main content area contains four input fields, each with a label and a value: 'Sex Ratio (Rural)' with '100.00', 'Sex Ratio (Urban)' with '100.00', 'Avg Household Size (Rural)' with '6.00', and 'Avg Household Size (Urban)' with '6.00'. Each input field has a minus sign on the left and a plus sign on the right. Below these fields is a red 'Predict' button. At the bottom, a green box displays the prediction: 'Prediction: Urban Dominated'.

6. Deployment with Streamlit

A user-friendly Streamlit web app was developed, allowing users to input sub-division statistics and receive a prediction on whether the region is urban or rural dominated. The app includes data visualization, model comparison, and real-time prediction.