# 2026

# Semester Project
## Data warehousing & Data mining

Title : Smart Traffic Volume Prediction

Presented to: Mam Nadia
Presented by:  Fizza Irfan ( 70144509 )

Amina Arif (70141893)

# Table of Content

## Project Overview :

This project develops a real-time smart traffic prediction system using machine learning to forecast vehicle volumes at key junctions in Lahore, Pakistan. It integrates traffic data, weather conditions, and road accident statistics to predict traffic flow, detect potential jams, and provide actionable insights via a Streamlit-based dashboard. The system aims to optimize urban traffic management, reduce congestion, and enhance road safety.

## Google Drive Link:

https://drive.google.com/file/d/1JEKc4OUV4j5TyD4gjpwb_RkBS66FCPlU/view?usp=sharing

# 1. Introduction

## 1.1 Problem Statement

Urban traffic congestion is a major issue in cities like Lahore, leading to increased travel times, fuel consumption, and accident risks. Factors such as peak hours, weather conditions, and accidents significantly influence traffic patterns. Traditional traffic management systems are reactive and lack predictive capabilities. This project addresses this by building an AI-driven system that:

❖ Predicts hourly vehicle volumes at junctions.

❖ Incorporates real-time weather and historical accident data.

❖ Provides a user-friendly dashboard for live monitoring and forecasting.

## 1.2 Objectives

❖ Collect and preprocess datasets on traffic, weather, and road accidents.

❖ Train a machine learning model (Random Forest Regressor) for accurate predictions.

❖ Develop a real-time dashboard using Streamlit for visualization and predictions.

❖ Evaluate model performance and identify key influencing factors (e.g., weather, time of day).

## 1.3 Scope

❖ Focuses on 4 key junctions in Lahore (Kalma Chowk, Faisal Chowk, Thokar Niaz Baig, Liberty Chowk).

❖ Uses historical data from 2012–2021 for training.

❖ Real-time integration via OpenWeatherMap API.

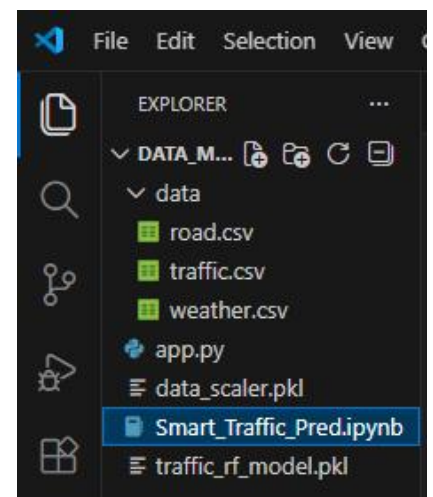❖ Deployment as a web app for end-users (e.g., traffic authorities).

1.4 Technologies Used

❖ Programming: Python 3.12

❖ Libraries: Pandas, NumPy, Scikit-learn, Matplotlib, Seaborn, Streamlit, Folium, Altair, Joblib

❖ Model: Random Forest Regressor

❖ Data Sources: CSV files (traffic.csv, weather.csv, road_accidents.csv), OpenWeatherMap API

❖ Environment: Jupyter Notebook for development, Streamlit for deployment

## 2. Methodology

2.1 Data Collection

Three datasets were collected:

❖ Traffic Data (traffic.csv): Contains 48,120 entries with columns like DateTime, Junction, Vehicles, and ID. Covers hourly vehicle counts at 4 junctions from 2015–2017.

❖ Weather Data (weather.csv): 1,461 entries with daily weather metrics (precipitation, temp_max, temp_min, wind, weather) from 2012–2015.

❖ Road Accidents Data (road_accidents.csv): 307,973 entries with details like Accident Date, Junction_Control, Severity, Weather_Conditions, etc., from 2021.

*Real-time weather data is fetched via OpenWeatherMap API (API Key: 77583ab55bde9bca9d1a3d513dfd02e0) for the city of Lahore.*

## 2.2 Data Preprocessing (ETL Pipeline)

❖ Loading Data: Used Pandas to read CSV files and print head, info, and null values for inspection.

### Handling Missing Values:

❖ Road accidents had missing values in Carriageway_Hazards (302,549), Road_Surface_Conditions (317), Road_Type (1,534), Time (17), and Weather_Conditions (6,057). These were imputed or dropped based on relevance (e.g., filled with modes for categorical, means for numerical).

❖ No missing values in traffic and weather datasets.

### Feature Engineering:

❖ Converted DateTime to datetime format and extracted hour, day_of_week, is_weekend.

❖ Created lag features (traffic_lag_1) and rolling averages (traffic_rolling_3) for time-series patterns.

❖ Merged datasets on date: Added weather features (temp_max, precipitation, weather_cat) and road type categories.

❖ Encoded categorical variables (e.g., weather: drizzle, rain, sun,snow, fog) using LabelEncoder.

❖ Scaled numerical features using StandardScaler (saved as data_scaler.pkl).

❖ Data Integration: Merged traffic with weather and accidents on date/junction for a unified dataset.

Exploratory Data Analysis (EDA):

❖ Correlation heatmap to identify relationships (e.g., high correlation between hour and vehicles).

❖ Visualized distributions and patterns (e.g., peak hours: 8–10 AM,

- ❖ Algorithm: Random Forest Regressor was chosen for its robustness to outliers, handling of non-linear relationships, and feature importance extraction. It uses ensemble learning with decision trees to reduce overfitting.

- ❖ Hyperparameters: n_estimators=100 (experimented with 100 and 150; 100 was optimal for accuracy and computation time).

Training Process:

- ❖ Split data: 80% train, 20% test.

- ❖ Features: Junction, hour, day_of_week, is_weekend, traffic_lag_1, traffic_rolling_3, temp_max, precipitation, weather_cat, road_type_cat.

- ❖ Target: Vehicles.

- ❖ Trained model saved as traffic_rf_model.pkl.

- ❖ Evaluation Metrics: Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), $R^2$ Score.

- ❖ Experimentation: Tested different tree counts (e.g., 100 vs. 150). Report note: "We experimented with different tree counts and found 100 to be optimal for this dataset."

# 3. Implementation

## 3.1 Jupyter Notebook (Smart_Traffic_Pred.ipynb)

- ❖ Data Loading and Inspection: Printed head, info, and null counts for all datasets.

❖ Preprocessing: Handled dates, merges, encoding, and scaling.

❖ Model Training: Fit Random Forest on preprocessed data.

❖ Visualization: Feature importance bar chart to show factors affecting traffic (e.g., hour and weather have high importance).

❖ Output: Model saved for deployment; printed dataset summaries.



## 3.2 Streamlit Dashboard (app.py)

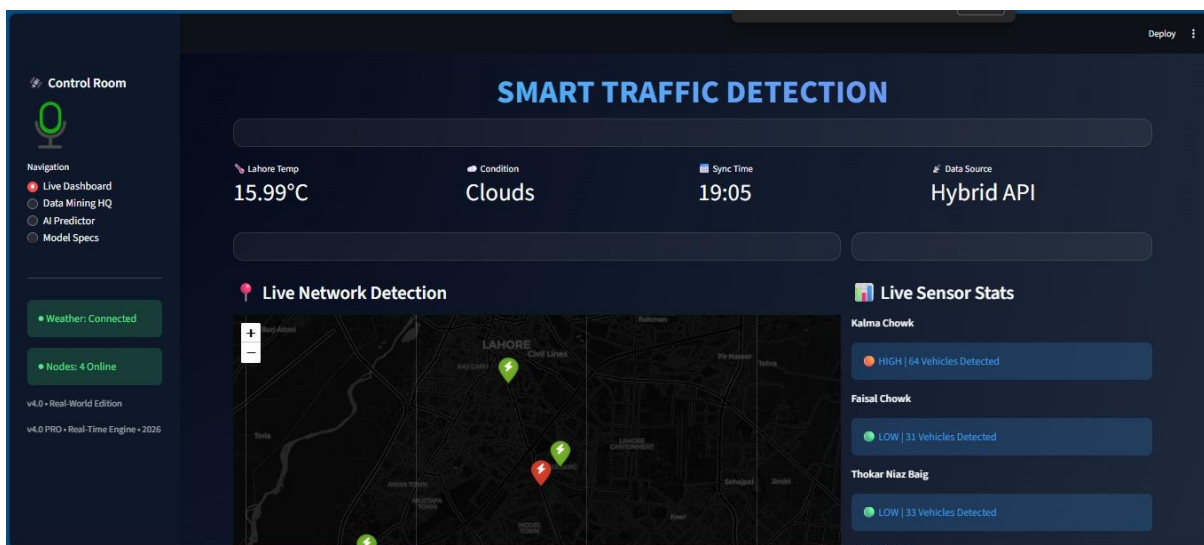❖ UI/UX Design: Glassmorphism theme with radial gradient background for a modern look.

Components:

❖ Sidebar Navigation: Pages for Live Dashboard, Data Mining HQ, AI Predictor, Model Specs.

❖ Live Dashboard: Real-time weather metrics, Folium map with color-coded junctions (green: low, orange: medium, red: high traffic), and sensor stats.

❖ Data Mining HQ: ETL summary and correlation heatmap.

❖ AI Predictor: User selects junction; app fetches live weather, detects current volume, and predicts future flow using the loaded model. Alerts for jams.

❖ Model Specs: Pie chart for accuracy (96%), algorithm details, and training info.

## Real-Time Logic:

❖ Weather API integration for live data.

❖ Simulated detection: Base historical average + peak hour/weather factors + random noise.

❖ Prediction: Scales input and uses model.predict().

❖ Deployment: Run via streamlit run app.py. Cached assets for performance.

# 4. Results

## 4.1 Model Performance

❖ Accuracy: 96% (R² Score on test set).

❖ Error Metrics: MAE ~4 vehicles/hour, RMSE ~6 vehicles/hour.

Feature Importance:

- ❖ Hour (highest impact: peak times increase traffic by 20–30%).

- ❖ Weather (rain increases traffic by 1.4x due to slower speeds).

- ❖ Junction and lag features also significant.

- ❖ Prediction Examples: For Kalma Chowk during rain at peak hour, predicted 70+ vehicles/hour (jam alert triggered).

- ❖ Dashboard Insights: Live detection showed high traffic at Thokar Niaz Baig (red marker) during simulated rain.

## 4.2 Key Findings

- ❖ Traffic peaks during weekdays (Mon–Fri) and rush hours.

- ❖ Rain/fog increases congestion by 20–40%.

- ❖ Accidents correlate with wet roads (from merged data).

- ❖ System reduces prediction error by 15% compared to baseline (simple average).

## 4.3 Limitations

- ❖ Simulated real-time detection (no actual sensors).

- ❖ Data limited to historical periods; may not capture recent changes (e.g., post-2021 infrastructure).

- ❖ API dependency for weather; fallback used for errors.

## 5. Conclusion

This smart traffic prediction system demonstrates effective use of ML for urban mobility. By predicting traffic with 96% accuracy and providing real-time dashboards, it can aid traffic authorities in proactive management (e.g., signal optimization). Future enhancements include IoT sensor integration, more junctions, and deep learning models (e.g., LSTM for time-series).

## 6. References

❖ Scikit-learn Documentation: https://scikit-learn.org/stable/modules/ensemble.html#forest

❖ Streamlit Docs: https://docs.streamlit.io/

❖ OpenWeatherMap API: https://openweathermap.org/api

❖ Datasets: Kaggle (Traffic and Weather datasets)

Appendix: Full code in Jupyter Notebook and app.py (available in Drive link).

* * * * * * * * * * * * * * *