# Predicting Traffic Delays Using U.S. Congestion and Weather Data

Team 26
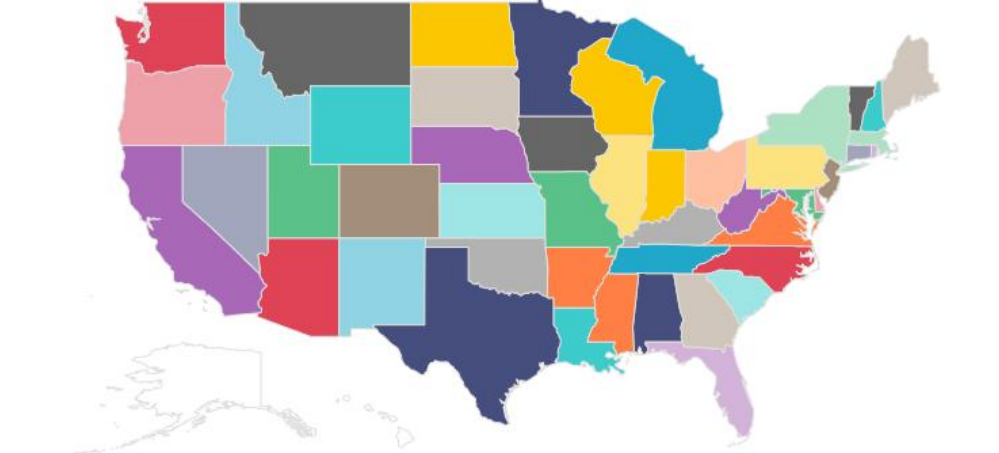
| Name | Role | Innomail |
|---|---|---|
| Alexey Tkachenko | Data engineer | a.tkachenko@innopolis.university |
| Daniil Abrosimov | ML specialist | d.abrosimov@innopolis.university |
| Egor Machnev | Data scientist | e.machnev@innopolis.university |
| Apollinaria Chernikova | Tester and Technical writer | a.chernikova@innopolis.university |

+ CI bot

# Introduction

- Analyze traffic congestion on US roads
- Identify key factors influencing traffic delays (weather, time, location).
- Predict traffic jam duration with reasonable accuracy.
- Develop a big data pipeline from raw CSVs to analytical dashboard.
- Provide decision support for urban planners and transportation authorities.

# Data Description

## US Traffic Congestions (2016-2022)

Comprehensive Dataset of 33 Million U.S. Traffic Congestion Events

**Total rows**

# 33.3M

Total rows

**Rows used for training**

# 1.96M

Rows used for training

# Architecture of data pipeline

| Stage | Input | Processing | Output |
|-------|-------|-----------|--------|
| Stage I | Kaggle ZIP archive (CSV, ~12.8 GB) | - Load and clean data<br>- Import to PostgreSQL<br>- Export to HDFS using Sqoop → convert to Parquet (Snappy-compressed) | Raw lake: /project/warehouse/traffic (Parquet) |
| Stage II | Parquet files from Stage I | - Create external Hive table traffic<br>- Create optimized partitioned & bucketed table traffic_partitioned<br>- Auto-ingest per-state<br>- Generate 14 EDA views (CSV + Hive) | Hive warehouse + /output/dashboard/qX.csv |
| Stage III | 50k-row-per-state sample from traffic_partitioned Hive table. ~2mil rows in total | - Spark MLlib pipeline with custom transformers:<br>- Cyclical time encoder<br>- GeoToECEF spatial transformer<br>- Word2Vec embedding<br>- OneHotEncoding for low cardinality categorical features<br>- Feature hashing for high cardinality features<br>- Grid search for LR and RF models | Trained models + evaluation metrics (RMSE, R², MAE) |
| Stage IV | Analytical results + model outputs | - Load into Apache Superset<br>- Configure dashboards, filters, drilldowns | Interactive BI dashboard |

GitHub

# Development Automation ✨

To streamline development and ensure continuous integration throughout the project lifecycle, we implemented a custom CI script tailored for the cluster environment.

→ Fetches and resets to the latest commit;
→ Executes the main pipeline script (main.sh);
→ Sends logs and outputs to a Telegram chat for real-time team visibility.

# Data preparation

# Data analysis

Initially formulated **14** hypotheses about factors that might influence traffic delays.

**7** were tested and validated through visual analysis using charts and dashboards.

# ML Modeling

→ **Random Forest Regressor**:
- ◆ maxDepth 10;
- ◆ numTrees 50.

→ **Linear Regression**:
- ◆ regParam 0.01;
- ◆ elasticNetParam 0.0 (pure L2 regularization)

| Model | RMSE | $R^2$ | MAE |
|---|---|---|---|
| Random Forest | 2.9146 | 0.5787 | 1.5438 |
| Linear Regression | 3.3102 | 0.4566 | 1.8117 |

# Data presentation



Superset
Dashboard

## Conclusion

Our contributions include:

→ End-to-end automation for data ingestion, partitioning, and transformation;

→ Targeted HQL queries for exploratory data analysis (EDA);

→ A trained machine learning model predicting traffic delays using geographic, temporal, and weather-related features;

→ An interactive Superset dashboard for intuitive data visualization and decision support.

# Reflections on own work

A solid and
awesome team!