FLORIDA STATE UNIVERSITY

Abstract

This study presents a two-phase approach to analyzing and predicting airline flight delays. The first phase consists of a comparative performance analysis between Neo4j and SQL in tracking cumulative flight delays across aircraft tail numbers. The study measures execution time and query efficiency in calculating cumulative delays across multiple aircraft. This database performance comparison provides insights into the scalability and efficiency of graph-based versus relational database approaches for flight delay tracking.

The second phase evaluates three machine learning models for delay prediction: Random Forest, Gradient Boosting, and XGBoost. These models were trained on flight data including basic flight information such as departure times, carrier details, and route distances. The study performed feature engineering to create additional predictors like weekend flight indicators and time-based features. Performance comparison between the three models was conducted using metrics like RMSE and prediction accuracy within various time thresholds. Feature importance analysis across all three models helped identify the most crucial factors in predicting flight delays.

Hypothesis

- Neo4j will outperform SQL for complex queries involving flight relationships
- Flight delays can be predicted based on factors such as day of week and origin airport

Comparative Analysis of Machine Learning Models for Flight Delay Tracking and Prediction

airports, and route distances.

ensure real-world applicability.





FLORIDA STATE UNIVERSITY