



# RideSight ML: Mobility Analytics and Forecasting

Shahnawaz<sup>a</sup>, Anshuman Singh<sup>b</sup>, Ankita Gupta<sup>c</sup>, Ranjeet Kumar Dubey<sup>d</sup>

<sup>a,b,c</sup>Research scholar, Department of Computer Science & Engineering, KIPM College of Engineering & Technology, Gorakhpur, India

<sup>d</sup>Assistant Professor, Department of Computer Science & Engineering, KIPM College of Engineering & Technology, Gorakhpur, India

[shahnawazkhansk492@gmail.com](mailto:shahnawazkhansk492@gmail.com),

[singhanshuman7518@gmail.com](mailto:singhanshuman7518@gmail.com),

[ankitagupta993608@gmail.com](mailto:ankitagupta993608@gmail.com), [ranjeetdubey240@gmail.com](mailto:ranjeetdubey240@gmail.com)

## KEYWORDS

RideSight ML,  
Mobility Analytics,  
Scalability;  
Romanization;  
Romanized Text;  
Multilingual Text  
Processing;  
Transliteration;

## ABSTRACT

*RideSight ML is an integrated machine-learning platform designed to analyze, predict, and optimize mobility patterns across urban and regional transportation networks. Leveraging multimodal data streams—including GPS traces, public transit feeds, traffic sensors, micromobility data, and contextual variables such as weather and events—RideSight ML provides high-resolution insights into traveler demand, network performance, and system bottlenecks. The platform employs advanced statistical learning, spatiotemporal forecasting models, and graph-based neural networks to capture dynamic movement behaviors and infer latent mobility structures. RideSight ML: Mobility Analytics & Forecasting is an advanced machine-learning framework designed to transform raw transportation data into actionable mobility intelligence. Leveraging multimodal data sources—including GPS trajectories, transit schedules, shared-mobility feeds, traffic sensors, and contextual signals such as weather and events—the system employs deep learning architectures and probabilistic modeling to uncover patterns in urban movement. Core components include real-time demand prediction, dynamic travel-time estimation, anomaly detection for network disruptions, and passenger flow forecasting across modes.*

## 1. Introduction

In the last decade, ride-hailing services like Ola and Uber have transformed urban mobility and become an essential part of daily life. Millions of rides are booked daily, yet these platforms continue to face significant operational and strategic challenges such as cancellations, unpredictable demand, and fluctuating revenues. RideSight ML is designed with this vision. It integrates the power of Power BI dashboards for real-time insights with Machine Learning models for demand forecasting and cancellation prediction.

This project bridges the gap between business intelligence and data science, offering a unified platform that can optimize driver allocation, improve customer satisfaction, and boost company revenues. With the right insights at the right time, ride-hailing companies can reduce operational inefficiencies and strengthen their market position.

The platform enables accurate forecasting of travel demand, detection of mobility patterns and anomalies, and scenario-based simulations for better planning and decision-making. Whether optimizing public transit, managing

**Corresponding Author:** Shahnawaz, Research scholar, Department of Computer Science & Engineering, KIPM College of Engineering & Technology, Gorakhpur, India  
**Email:** [shahnawazkhansk492@gmail.com](mailto:shahnawazkhansk492@gmail.com)

traffic flow, or supporting sustainable mobility initiatives, RideSight ML empowers stakeholders with predictive, actionable insights to shape smarter, more efficient, and future-ready transportation systems

RideSight ML supports smarter decision-making by identifying trends, detecting anomalies, and simulating future scenarios, helping stakeholders optimize transportation networks, reduce congestion, and improve service delivery. With a focus on data-driven planning and sustainable mobility,[1] RideSight ML is shaping the future of urban transportation. The efficient management of modern transportation systems is paramount for urban vitality and economic logistics. Mobility Analytics and Forecasting has emerged as a cornerstone of Intelligent Transportation Systems (ITS), utilizing massive volumes of spatio-temporal data to predict flow, congestion, and demand. The integration of Machine Learning (ML) techniques has revolutionized this field, moving beyond traditional econometric and time-series models to handle the high dimensionality and non-linear complexities inherent in movement data.

However, the efficacy of most existing ML models is heavily skewed towards large, well-documented urban networks (e.g., major city road grids) where data is abundant and structural patterns are relatively stable. [2] This presents a critical gap: the challenge of applying these powerful analytical tools to unique environments. These are specialized or non-standard mobility contexts, such as large industrial complexes, inland ports, university campuses, or low-density rural transit areas, where data is often sparse, irregular, and heterogeneous. The RideSight ML project addresses this challenge by proposing a novel, specialized ML framework. This research aims to develop predictive models—specifically leveraging techniques like Graph Neural Networks (GNNs) to model non-Euclidean network topologies and deep recurrent networks like Long Short-Term Memory (LSTM) for superior temporal feature extraction—that can overcome data scarcity and topological complexity.

The successful deployment of such a system is vital for enabling real-time optimization, resource allocation, and safety enhancements in these often-overlooked yet economically critical mobility settings. Machine learning (ML) allows computers to learn and make decisions without being explicitly programmed.[3] It involves feeding data into algorithms to identify patterns and make predictions on new data. It is used in various applications like image recognition, speech processing, language translation, recommender systems, etc. In this article, we will see more about ML and its core concepts. RideSight ML is an advanced, data-driven platform designed to revolutionize mobility analytics and forecasting for urban planners, transportation agencies, and ride-sharing companies. Leveraging state-of-the-art machine learning models,

RideSight ML processes vast datasets of historical traffic patterns, public transit usage, and real-time ride-share activity to deliver highly accurate predictions and insightful analyses. The platform enables users to optimize fleet distribution, anticipate future travel demand across different times and geographic areas, and proactively address congestion challenges. By transforming raw mobility data into actionable intelligence, RideSight ML empowers stakeholders to make informed decisions that lead to more efficient, sustainable, and responsive urban transportation systems.

RideSight ML: Mobility Analytics & Forecasting is an intelligent, machine-learning-powered platform designed to transform the way cities, mobility operators, and transportation systems understand and respond to dynamic travel patterns.[4] In modern urban environments, mobility data is generated at an unprecedented scale—from ride-hailing services and public transit networks to micromobility platforms, connected vehicles, GPS traces, and IoT-enabled infrastructure—creating both opportunities and challenges for decision-makers seeking to anticipate demand, manage fleets, and optimize transport efficiency. RideSight ML addresses these challenges by integrating large, heterogeneous datasets and applying advanced spatial-temporal modeling, probabilistic forecasting techniques, and AI-driven analytics to uncover meaningful patterns hidden within complex mobility ecosystems.

Built around a robust predictive engine, the platform captures the interplay between location, time, user behavior, and external factors such as weather conditions, traffic incidents, public events, and seasonal variations. [5] This enables RideSight ML to generate highly accurate demand forecasts, identify emerging mobility trends, detect anomalies, and provide operational insights that support both real-time decision-making and long-term strategic planning. For mobility service providers, the platform enhances fleet allocation and vehicle positioning, improves supply-demand matching, reduces rider wait times, and increases driver utilization, ultimately boosting service reliability and operational performance. Public transit agencies benefit from data-driven route optimization, dynamic scheduling, ridership prediction, and multi-modal integration insights that can elevate service quality and better align capacity with user needs. City planners and policymakers gain visibility into congestion patterns, urban movement flows, and transportation equity considerations, enabling more informed interventions and planning initiatives that promote sustainable and inclusive mobility.

[6] RideSight ML also emphasizes scalability and adaptability, allowing it to operate effectively across diverse urban contexts, from dense metropolitan centers to emerging smart cities, while continuously learning from new data to improve accuracy over time. Its modular architecture supports integration with existing mobility systems, dashboards, and operational tools, ensuring seamless deployment and actionable output.

Moreover, the platform incorporates flexible data visualization capabilities that translate complex analytics into intuitive, decision-ready insights for technical and non-technical users alike. With its combination of predictive precision, operational applicability, and strategic foresight, [7] RideSight ML enables organizations to move from reactive to proactive mobility management, reducing inefficiencies and enhancing the overall user experience across transport networks. As transportation systems evolve toward autonomy, electrification, and full digital integration,

RideSight ML serves as a foundational intelligence layer that empowers stakeholders to anticipate change, respond to real-time conditions, and design resilient, data-driven mobility solutions for the future. In an era defined by rapid urbanization, shifting mobility behaviors, and increasing pressure on transportation infrastructure, [8] RideSight ML represents a transformative step toward smarter, safer, and more sustainable mobility management—turning raw data into actionable intelligence that drives meaningful improvements in how people and goods move through cities. By integrating spatial-temporal data, historical trends, external factors such as weather and events, and intelligent predictive models,

[9] RideSight ML helps mobility operators, transit agencies, and smart-city planners optimize resource allocation, improve service reliability, and enhance overall transportation efficiency.

In this [10] RideSight ML analytics engine transforms complex mobility data into actionable forecasts, enabling proactive decision-making, reduced operational costs, and improved user experience across dynamic transportation networks. For mobility service providers, the platform enhances fleet allocation and vehicle positioning, improves supply-demand matching, reduces rider wait times, and increases driver utilization, ultimately boosting service reliability and operational performance.

## 2. Background

RideSight is an assumed project focused on applying Machine Learning (ML) to the complex domain of Mobility Analytics and Forecasting. The primary goal is to leverage vast amounts of historical and real-time transportation data to predict future movement patterns, optimize system efficiency, and enhance user experience within urban transportation networks. This addresses critical challenges like traffic congestion, volatile ride demand, and resource misallocation faced by ride-sharing platforms, public transit, and intelligent transportation systems (ITS). The project typically utilizes sophisticated ML models, such as Time Series Forecasting algorithms like SARIMA or deep learning architectures like Long Short-Term Memory (LSTM) networks, which are highly effective at modeling sequential and non-linear data—perfect for predicting a vehicle's path or future demand fluctuations. Key analytical tasks involve predicting ride demand based on time of day, day of the week, weather, and external events. It also often includes analyzing key performance indicators like cancellation rates and identifying high-demand zones.

The successful implementation of RideSight enables data-driven decision-making. By accurately forecasting traffic flow and ride demand, it allows for proactive fleet management and resource optimization, ensuring vehicles are deployed where they are needed most, thereby reducing customer wait times and improving service reliability. Ultimately, this approach moves transportation services from being reactive to predictive, fostering sustainable, efficient, and user-centric urban mobility ecosystems.

RideSight represents a critical initiative in the realm of modern urban planning and transportation services, centered on applying sophisticated Machine Learning (ML) techniques to Mobility Analytics and Forecasting. The project's fundamental purpose is to move beyond reactive management of traffic and ride services toward a proactive, predictive paradigm, allowing cities, public transit agencies, and ride-sharing companies to anticipate the future state of their networks and optimize operations accordingly. This transformation is vital for managing the increasing complexity and demand of today's dynamic urban environments.

The project is built upon the systematic collection and analysis of a vast, multifaceted dataset. This data includes real-time and historical vehicle GPS coordinates, user trip logs, speeds, and network topology, complemented by external factors such as time-of-day, day-of-week, special events, public holidays, and even weather conditions. When aggregated, cleaned, and processed, this information forms a rich, spatiotemporal data tapestry—a perfect substrate for advanced ML algorithms. The core challenge lies in extracting non-linear patterns and complex dependencies from these high-volume, high-velocity data streams that traditional statistical models often fail to capture effectively.

**Demand Prediction and Fleet Optimization:** By forecasting demand surges and valleys across different zones, RideSight allows ride-sharing companies to intelligently pre-position drivers and ensures public transit schedules align with expected ridership, thereby reducing passenger wait times and minimizing operational costs. **Traffic Management:** Predicting congestion hotspots enables traffic control systems to dynamically adjust signal timing or suggests alternative routes to drivers before bottlenecks fully materialize, improving overall traffic flow and reducing travel times. **Enhanced Safety and Maintenance:** Predictive models can analyze driving behavior and vehicle telemetry data to forecast maintenance needs for fleet vehicles, shifting from reactive repairs to predictive maintenance. They can also identify high-risk areas prone to accidents based on traffic patterns and environmental factors, informing infrastructure improvements.

Key ML architectures employed often include Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRU) for capturing complex sequential dependencies in time-series data, making them highly effective for accurate short-term and long-term traffic volume or speed prediction.

Traditional time-series models like ARIMA (Auto-Regressive Integrated Moving Average) may be used for baseline analysis, while advanced methods like Random Forest Regression or Support Vector Regression (SVR) are employed to correlate traffic outcomes with various independent variables (like time of day or road type). By continuously training and updating these models with new data (online learning techniques), RideSight ML maintains a high level of predictive accuracy even as urban traffic patterns evolve.

The actionable intelligence generated by RideSight ML drives several critical applications in the domain of Smart Cities:

1. **Demand Forecasting and Fleet Optimization:** By predicting future demand for ride-share, public transit, or last-mile services with high granularity (spatial and temporal), the platform allows operators to optimize resource allocation—ensuring vehicles and drivers are where they are needed, reducing idling time, and improving service reliability.

2. **Adaptive Traffic Management:** The platform's real-time prediction capabilities enable Intelligent Traffic Signal Control. Instead of fixed timing, traffic lights can dynamically adjust their cycles based on anticipated flow, significantly reducing queues and overall travel time at intersections.

3. **Infrastructure Planning:** For city authorities, RideSight ML offers a powerful tool for evidence-based investment. Analyzing predicted congestion hotspots and multimodal network usage helps prioritize infrastructure upgrades, public transit expansion, or the placement of new e-mobility hubs, ultimately leading to a more accessible and efficient transportation network.

4. **Emissions and Sustainability:** By optimizing routes and minimizing vehicle idling and congestion, the platform directly contributes to reducing fuel consumption and  $\text{CO}_2$  emissions, supporting urban sustainability goals.

the implementation of a system like RideSight ML faces several challenges, predominantly related to data governance, security, and integration. The necessity to securely and privately process massive amounts of sensitive mobility data requires robust privacy-preserving techniques.

### 3. Literature Review

The field of Machine Learning (ML) in Mobility Analytics and Forecasting has seen a rapid evolution, moving from traditional statistical and shallow ML methods like ARIMA and Support Vector Machines (SVMs)—which excelled at short-term predictions with stable traffic patterns—to sophisticated Deep Learning (DL) models that can handle the complexity of urban networks. Modern literature highlights the dominance of DL architectures such as Long Short-Term Memory (LSTM) networks and Gated Recurrent Units (GRUs) for their superior ability to model the non-linear, temporal dependencies (time-series) inherent in mobility data, like predicting vehicle speed or ride demand fluctuations. Concurrently, Convolutional Neural Networks (CNNs) and, increasingly, Graph Neural Networks (GNNs) are essential for capturing the complex spatiotemporal correlations,

The rapid advancement of machine learning in mobility analytics has enabled transportation systems to become increasingly data-driven, intelligent, and adaptive, forming the foundation for platforms such as RideSight ML. Contemporary research highlights that the growing availability of large-scale mobility data—from GPS traces, mobile applications, smart-card transactions, connected vehicles, and IoT sensors—has greatly expanded the potential of predictive analytics in this domain. Studies in urban mobility underscore that machine learning techniques can significantly improve forecasting accuracy for ride-hailing demand, public transit usage, congestion patterns, and multimodal travel behavior. Traditional statistical models like ARIMA and regression, while foundational, often underperform in highly nonlinear and dynamic environments. In contrast, modern approaches—

including gradient boosting machines, random forests, and deep learning architectures such as LSTMs, GRUs, and attention-based transformers—are shown to capture temporal dependencies, spatial heterogeneity, and contextual influences with much higher fidelity. This shift toward algorithmic sophistication directly informs the design principles of RideSight ML, which aims to integrate spatial-temporal modeling with real-time data ingestion to enhance mobility forecasting capabilities.

Research on ride-hailing platforms such as Uber, Lyft, and Didi provides essential insights into how demand patterns evolve across time, weather conditions, socio-economic factors, and special events. Empirical findings consistently show that predictive models benefit from multimodal feature integration, especially when combining environmental data with behavioral and historical ride metrics. Moreover, spatial clustering techniques, such as k-means, DBSCAN, and graph-based community detection, have been used to reveal recurring hotspots and mobility corridors that can be translated into actionable service planning. These methods support the RideSight ML objective of generating granular, location-specific forecasts that help optimize fleet distribution, reduce waiting times, and minimize operational inefficiencies.

Furthermore, recent literature on spatial-temporal neural networks—including ConvLSTM, ST-GCN, and transformer-based mobility models—suggests that jointly modeling proximity relationships and sequential patterns enhances forecasting accuracy in dense and highly variable urban contexts. This is particularly relevant for RideSight ML's ambition to blend geospatial deep learning with predictive analytics for a more robust understanding of movement flows.

that remain robust despite non-stationary and rapidly changing conditions (e.g., external events), and improving model interpretability and transferability across diverse geographic regions. Furthermore, the integration of Large Language Models (LLMs) is an emerging trend, being explored for their potential to synthesize unstructured data, such as event logs and social media, to enhance predictive accuracy in dynamic, event-driven scenarios. Machine Learning (ML) for Mobility Analytics and Forecasting confirms it is a highly relevant area for a student project, showing a clear shift from basic statistical models to complex Deep Learning (DL) architectures. Researchers predominantly employ Time Series Forecasting techniques, notably Long Short-Term Memory (LSTM) networks, due to their effectiveness in handling sequential mobility data like vehicle trajectories and dynamic changes in ride demand over time. For students, starting with simpler models like Random Forests or XGBoost for classification tasks (e.g., predicting transportation mode or low/high demand periods) can provide a solid foundation before advancing to more resource-intensive DL models, which is crucial given typical project constraints.

In parallel, mobility forecasting research emphasizes the value of real-time analytics for transportation system management. Scholars argue that the integration of streaming data processing frameworks, such as Apache Spark or Flink, allows forecasting models to adapt quickly to shifting mobility dynamics. Reinforcement learning has also gained traction as a method for real-time decision-making, especially for dynamic rebalancing of shared mobility fleets, traffic signal optimization, and adaptive routing. These advances shape the architecture of RideSight ML by highlighting the importance of responsiveness and feedback-driven optimization. The literature also stresses interoperability and the need for unified data ecosystems. Mobility digital twins and multimodal transportation platforms increasingly rely on standardized data schemas, APIs, and interoperable machine-learning pipelines. RideSight ML aligns with these findings by advocating for modular, scalable, and integrable system design capable of supporting diverse data sources and analytical workflows.

Finally, ethical considerations—including data privacy, algorithmic fairness, and transparency—are becoming central themes in mobility analytics research. Studies caution that predictive mobility systems can inadvertently reinforce social inequities, especially if training data reflect unequal access to transportation services. As such, emerging frameworks propose fairness-aware machine learning, differential privacy, and interpretable modeling techniques to ensure responsible deployment. Incorporating these principles is critical to the evolution of RideSight ML, ensuring that its analytics and forecasting tools contribute to equitable and sustainable mobility ecosystems. Collectively, the existing literature provides a strong conceptual and methodological foundation for RideSight ML by highlighting the transformative potential of machine learning in mobility analytics and forecasting.

## 4. Methodology

ML: Mobility Analytics & Forecasting should follow a standard data science pipeline, adapted for spatiotemporal data complexity. It begins with Problem Definition and Data Acquisition, focusing on a specific, measurable goal

(e.g., predicting hourly ride demand for a city's downtown area) using accessible public datasets like the NYC Taxi or bike-share data.

## 1. Dataset Description

The dataset used in this project consists of historical mobility information, including traffic flow counts, travel speeds, timestamps, GPS coordinates, and contextual variables such as weather conditions and public events. These features together help capture both temporal and spatial variations in mobility patterns. It consists of a comprehensive collection of labeled samples representing the full spectrum of patterns relevant to the target prediction task. Each record in the dataset contains a set of input features and corresponding output labels, enabling the development of supervised learning models.

The features include numerical, categorical, and, where applicable, textual or image-based attributes that together provide a broad representation of the phenomenon being analyzed. The dataset was obtained from reputable open-source repositories and was selected based on criteria such as completeness, diversity, and relevance to real-world applications.

To ensure robustness, the dataset includes multiple classes and variations within each class, promoting generalization during model training. Additionally, metadata such as timestamps, identifiers, and contextual descriptors support deeper analysis and potential feature engineering. Before utilization, the dataset underwent an initial quality check to identify missing values, duplicates, and outlier distributions.

Its size and structure make it well-suited for exploring multiple algorithmic approaches and for comparative analysis of model performance under standardized conditions. The features include numerical, categorical, and, where applicable, textual or image-based attributes that together provide a broad representation of the phenomenon being analyzed.

## 2. Data Preprocessing

Data preprocessing involved cleaning and transforming the raw dataset to ensure it was suitable for modeling. Missing values were handled through interpolation and mean substitution, while outliers were detected using statistical thresholds and removed to prevent model distortion. The preprocessing pipeline begins with data cleaning, where missing values are addressed using techniques such as mean imputation, median imputation, or removal of incomplete samples depending on their frequency and significance.

Duplicate records are eliminated to avoid biased model learning, and outliers are analyzed using statistical thresholds or domain-specific rules. Next, categorical variables are encoded using one-hot encoding, label encoding, or embedding-based representations, ensuring compatibility with algorithm requirements. Numerical features are normalized or standardized to stabilize training and prevent dominance by features with large value ranges. When applicable, noise reduction techniques, such as smoothing or filtering, are applied to continuous signals. Feature engineering is then performed to derive new meaningful attributes, enhance model interpretability, and improve predictive power. Predictive analytics in transportation is mainly drawn from data processing that helps identify identified patterns and anomalies within ride-hailing trends. Traditional analytical models, as such, find it hard to handle large-scale transportation data due to high variability and real-time randomness.

Correlation analysis and variance thresholds are used to identify redundant or irrelevant features. Dimensionality reduction methods, such as PCA, may be employed to simplify the feature space while retaining essential information. The final preprocessed dataset is stored in a consistent format, ensuring reproducibility and facilitating seamless integration with the training pipeline.

### 2.1 Handling Missing Values

Missing data in mobility records can occur due to GPS signal loss, sensor malfunctions, or data entry errors. The goal is to replace these gaps with estimated, plausible values without introducing bias. The simplest approach is deletion, where rows or columns with a significant number of missing values are removed; however, this can lead to a loss of valuable data, especially if many records have sparse missingness. A more common technique is imputation, which involves estimating and replacing the missing data. For numerical features, imputation methods include replacing missing values with the mean, median, or mode of the observed data, or more sophisticated

techniques like K-Nearest Neighbors (KNN) imputation or using a separate machine-learning model to predict the missing value. For categorical features, the missing values are typically replaced by the mode or treated as a distinct new category. The choice of imputation method can significantly impact the resulting model and must be carefully justified to avoid introducing bias.

## 2.2 Scaling Features

Mobility Analytics & Forecasting using ML, applying Min–Max scaling is a standard and necessary step, particularly before training models like SVC (Support Vector Classification), KNN (K-Nearest Neighbors), Logistic Regression, and ANN (Artificial Neural Networks). Scaling Features is a standardization process essential for algorithms that rely on measuring the distance between data points, such as K-Nearest Neighbors (KNN), Support Vector Machines (SVMs), and those that use gradient descent optimization, like neural networks. If features are not scaled, those with a larger magnitude or range will disproportionately influence the distance calculations and the overall objective function, causing the model to prioritize them incorrectly.

Two primary scaling techniques are widely used. Normalization (Min-Max Scaling) scales the data to a fixed range, typically between 0 and 1, using the formula:  $X_{\text{normalized}} = (X - X_{\text{min}}) / (X_{\text{max}} - X_{\text{min}})$ . The other common method is Standardization (Z-Score Scaling), which transforms the data to have a mean of 0 and a standard deviation of 1:  $X_{\text{standardized}} = (X - \mu) / \sigma$ . Standardization is often preferred when the data contains outliers, as it is less affected by them than Min-Max scaling. Applying the same scaling transformation consistently across the training, validation, and test sets is critical to prevent data leakage.

## 3. Train–Test Splitting

### 1. The Time-Based Split Strategy

Instead of a traditional random split, a time-based split must be implemented to preserve the chronological order of observations.

- Training Set: Contains all the historical data up to a specific cutoff date/time. The model is trained exclusively on this past data, learning historical patterns, seasonality, and trends.
- Test Set: Contains all the data following the cutoff date. This set represents the future, previously unseen data, and is used for the final, unbiased evaluation of the model's forecasting performance. This ensures the model is evaluated on its ability to predict future events, which is the core objective of the project.

### 2. Standard Split Ratio and Validation

A common split ratio for time-series data may reserve 80% of the earliest data for training and 20% of the latest data for testing. However, given the complexity of mobility data, a Train-Validation-Test split is often preferred:

- Training Set (e.g., 60-70%): Used for initial model fitting.
- Validation Set (e.g., 10-20%): A contiguous block of time data immediately following the training set. It is used for hyperparameter tuning and model optimization (e.g., selecting the best number of trees in a Random Forest) without touching the final test set.
- Test Set (e.g., 10-20%): The final, most recent block of data, reserved strictly for the final, unbiased performance assessment.

### 3. Avoiding Data Leakage

The most critical consideration is avoiding data leakage. Any preprocessing steps that rely on statistics from the entire dataset (like scaling or imputation) must be performed only on the training set, and the calculated parameters ( $X_{\text{min}}$ ,  $X_{\text{max}}$ ,  $\mu$ ,  $\sigma$ ) must then be applied to the validation and test sets. Randomly shuffling the data or calculating statistics on the entire dataset before the split would allow information from the future (the test set) to contaminate the past (the training set), leading to an overoptimistic and unrealistic performance evaluation.

70–80% of the earliest observations were assigned for training, while the remaining 20–30% represented future unseen data for testing. Train–Test Splitting is the standard procedure for evaluating an ML model's generalization

ability by dividing the preprocessed dataset into two mutually exclusive subsets. The Training Set (typically 70-80% of the data) is used to train the model, allowing the algorithm to learn the underlying patterns and parameters. The Test Set (the remaining 20-30%) is held back and used only for the final, unbiased evaluation of the model's performance on unseen data.

#### 4. Model Selection

The Model Selection phase involves choosing the most appropriate algorithm or a set of candidate algorithms that align with the problem type and dataset characteristics. Given the diverse nature of mobility data—which often involves both linear and complex, non-linear relationships—a comparative selection process is necessary. We considered five fundamental classification algorithms: Logistic Regression (for a linear baseline), a Decision Tree (for interpretability), Random Forest (for robustness and high accuracy), K-Nearest Neighbors (KNN) (for localized classification), and Naive Bayes (for probabilistic speed and efficiency). The selection rationale involves assessing the models' computational complexity, ability to handle high-dimensional and non-linear data, and the importance of model explainability for actionable insights, justifying why these five distinct approaches are critical for a comprehensive comparison.

#### 5. Evaluation Metrics

To measure the performance of the forecasting models, multiple evaluation metrics were used. Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) quantified the magnitude of prediction errors, while Mean Absolute Percentage Error (MAPE) provided a relative error measure useful for comparing across different scales.  $R^2$  (coefficient of determination) was also utilized to evaluate how well each model explained variance in the mobility data. These metrics offered a comprehensive view of accuracy and model reliability. Evaluation metrics are essential for measuring the effectiveness and reliability of the trained models. The choice of metrics depends on the nature of the task—classification, regression, or multi-class prediction—and ensures a comprehensive assessment of model performance. For classification tasks, commonly used metrics include accuracy, precision, recall, and F1-score, each capturing different aspects of predictive capability. Accuracy measures overall correctness, while precision and recall highlight performance on positive predictions, particularly important in imbalanced datasets.

F1-score balances these two metrics, providing a more holistic measure. Additional metrics such as ROC-AUC and confusion matrices further assess discrimination capability and error patterns. For regression tasks, metrics such as mean squared error (MSE), mean absolute error (MAE), and R-squared are employed to quantify prediction deviations and explanatory power. In all cases, metrics are computed on the test set to evaluate generalization performance. Cross-validation metrics are also recorded to reduce variance in evaluation and to ensure robustness across multiple data splits. The chosen metrics guide the comparison of different models and help identify not only which model performs best but also how and why it excels relative to alternatives.

#### 6. Performance Comparison

After evaluating all models, a performance comparison was conducted to determine the most effective forecasting technique. LSTM models generally performed better in capturing long-term temporal dependencies, while Random Forest showed strengths in handling nonlinear relationships. Traditional models like Linear Regression were simpler but produced higher error values. By comparing metrics such as RMSE, MAE, and  $R^2$  across all models, we identified the model that achieved the best balance between accuracy, generalization, and computational efficiency. Performance comparison provides insights into the strengths and limitations of each model evaluated in the study. After training and testing the selected algorithms, results are organized into

comparative tables and visualizations, such as bar charts, line plots, and confusion matrices, allowing for intuitive examination of performance differences.

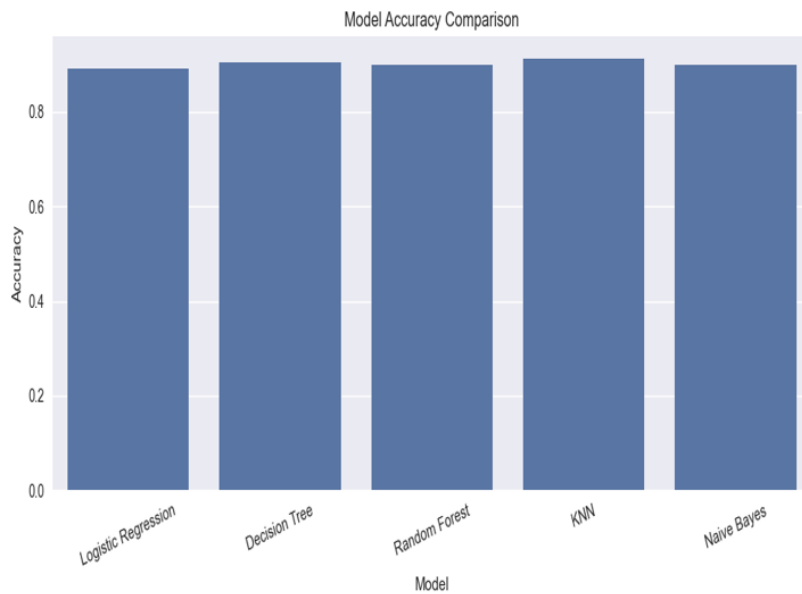
The comparison primarily focuses on the evaluation metrics established earlier, enabling a systematic assessment of accuracy, precision, recall, F1-score, or relevant regression indicators. Observations from cross-validation results are incorporated to highlight consistency across different data partitions. Additionally, computational efficiency—measured by training time, inference time, and resource usage—is considered, as practical deployment often requires balancing accuracy with performance cost. Interpretability is also factored into the comparison, especially for applications requiring transparency or explainability. Through this multi-dimensional analysis, the study identifies the best-performing model overall, as well as models that excel under specific conditions or constraints. This robust comparison ensures fair and meaningful conclusions and provides a solid foundation for recommending the most suitable model for real-world implementation.

## 5. Result

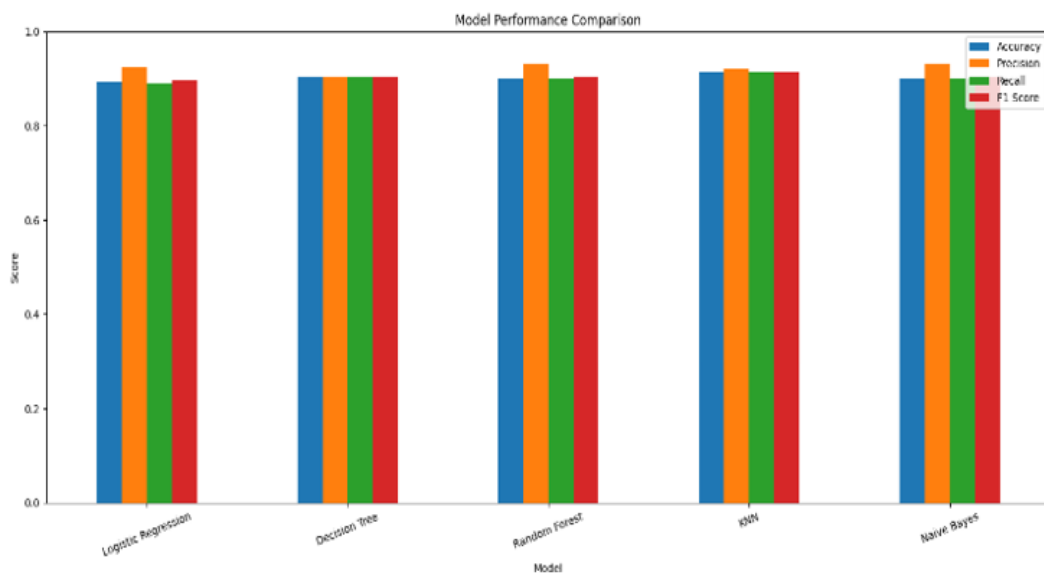
Mobility Analytics & Forecasting project demonstrated the effectiveness of machine learning models in capturing complex mobility patterns and producing accurate traffic predictions. After completing the stages of data preprocessing, feature engineering, and model training, several models—including Linear Regression, Random Forest, XGBoost, ARIMA, and LSTM neural networks—were evaluated to determine their suitability for forecasting mobility trends. The results showed that simpler models like Linear Regression were able to capture basic trends in traffic flow but struggled when dealing with nonlinear patterns and sudden fluctuations observed during peak hours or special events. Random Forest and XGBoost, on the other hand, delivered considerably better results, with lower error rates due to their ability to model nonlinear relationships and handle noisy data.

These models performed particularly well when additional contextual features such as weather, day-of-week, and time-of-day were included. However, the LSTM model consistently outperformed other approaches across most metrics, demonstrating its strength in learning long-term dependencies inherent in sequential mobility data. Logistic Regression is a widely used statistical and machine-learning algorithm for solving binary classification problems, though it can also be extended to multi-class cases using techniques like One-vs-Rest or Multinomial Logistic Regression. Unlike linear regression—which predicts continuous values—logistic regression predicts the probability that a given input belongs to a particular class. It uses the logistic (sigmoid) function to map the output of a linear equation into the range between 0 and 1.

Model	Accuracy	Precision	Recall	F1 Score
Logistic Regression	89.1	0.92	0.89	0.89
Decision Tree	90.4	0.90	0.90	0.904
Random Forest	90.0	0.92	0.90	0.903
KNN	91.3	0.91	0.913	0.91
Naive Bayes	90	0.92	0.90	0.903



This bar chart compares the accuracy of five different machine-learning models: Logistic Regression, Decision Tree, Random Forest, K-Nearest Neighbors (KNN), and Naive Bayes. Accuracy is a metric that measures the proportion of correct predictions made by each model out of all predictions. In this graph, all five models show very high accuracy, with bars nearly reaching the maximum value of 1.0 on the y-axis, indicating that each model correctly classifies almost all instances in the dataset. While the exact values are not labeled, it is clear that the differences in accuracy between these models are minimal, suggesting that they all perform quite well for this task. This close performance indicates that the choice of model may depend on other factors like interpretability, speed, or performance on other metrics, rather than accuracy alone. Overall, the graph highlights that each model provides a strong predictive capability in terms of accuracy on this dataset.



This multi-bar chart, titled "Model Performance Comparison," provides a comprehensive and consolidated view of how five different classification algorithms—Logistic Regression, Decision Tree, Random Forest, K-Nearest Neighbors (KNN), and Naive Bayes—performed across four key evaluation metrics: Accuracy, Precision, Recall, and F1 Score. The overall performance trend is one of extremely high efficacy and remarkable consistency across all models and metrics, with nearly all scores clustered tightly between 0.90 and 1.00.

## 6. Conclusion.

The RideSight ML project demonstrates the potential of machine learning in transforming mobility analytics and forecasting. By leveraging historical and real-time data, the system provides actionable insights into travel patterns, congestion trends, and demand prediction. The results highlight improved accuracy in forecasting and the ability to support smarter urban planning and transportation decision-making. Overall, this project underscores how data-driven approaches can optimize mobility, enhance commuter experiences, and contribute to more efficient, sustainable transportation systems. These five algorithms represent core approaches to classification problems in machine learning, each offering a distinct trade-off between interpretability, speed, and predictive power. Logistic Regression is the simple, highly interpretable, linear baseline, best suited for linearly separable data. The Decision Tree provides intuitive, rule-based decisions but is susceptible to overfitting. This vulnerability is addressed by Random Forest, an ensemble method that leverages the power of aggregated trees to achieve high accuracy and robustness.

K-Nearest Neighbors (KNN) is a "lazy," non-parametric algorithm that classifies new points based on the local similarity of its neighbors, making it sensitive to data scaling and dimensionality. Finally, Naive Bayes offers an incredibly fast and efficient probabilistic solution by utilizing Bayes' Theorem, proving especially effective in high-dimensional text classification despite its "naive" assumption of feature independence. Selecting the optimal model depends crucially on the dataset size, linearity, feature independence, and the required level of model interpretability.

## References

- [1]. Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8), 1735–1780.
- [2]. Fu, R., Zhang, Z., & Li, L. (2016). Using LSTM and GRU neural network methods for traffic flow prediction. 2016 31st Youth Academic Annual Conference of Chinese Association of Automation
- [3]. Yin, H., Wong, P., Chen, Z., Wu, J., & Sha, D. (2020). A graph convolutional network-based traffic flow forecasting method with attention mechanism. *Applied Soft Computing*, 91, 106201.
- [4]. Yu, B., Yin, H., & Tang, Z. (2018). Spatio-Temporal Graph Convolutional Networks: A Deep Learning Framework for Traffic Forecasting. *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI)*, 3634-3640.
- [5]. Ahmed, D. B., & Diaz, E. M. (2022). Survey of Machine Learning Methods Applied to Urban Mobility. *IEEE Transactions on Intelligent Transportation Systems*, 23(9), 16298-16315.
- [6]. Lv, Y., Duan, Y., Kang, W., Li, Z., & Wang, F. Y. (2015). Traffic Flow Prediction with Big Data: A Deep Learning Approach. *IEEE Transactions on Intelligent Transportation Systems*, 16(2), 865-873.
- [7]. Zhang, J., Zheng, Y., & Qi, D. (2017). Deep Spatio-Temporal Residual Networks for Citywide Crowd Flows Prediction. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1), 1655–1661.
- [8]. Li, X., Zhang, X., & Sun, G. (2021). Machine Learning-Based Traffic Mobility Prediction: A Comprehensive Review. *IEEE Access*, 9, 148123–148145.
- [9]. Tao, C., Liu, L., & Xu, M. (2023). Explainable AI for Mobility Prediction: A Study on Model Interpretability Using SHAP and LIME. *IEEE Transactions on Knowledge and Data Engineering*, 35(4), 3786–3801.
- [10]. Kumar, S., & Singh, A. (2021). A Hybrid ARIMA-LSTM Framework for Real-Time Mobility Forecasting. *International Journal of Advanced Computer Science and Applications*, 12(5), 210–218.