

A

**REAL TIME RESEARCH PROJECT REPORT ON
SUBWAY PASSENGER FLOW FORECASTING USING MACHINE
LEARNING**

*Submitted in partial fulfilment for the award of the degree of
BACHELOR OF TECHNOLOGY*

**In
COMPUTER SCIENCE AND ENGINEERING ARTIFICIAL INTELLIGENCE AND
MACHINE LEARNING**

By

GAJULA NIKHIL : 22Q91A6685

HARSH R BAGTHARIA : 22Q91A6691

MADHARAPU NITHIN KUMAR : 22Q91A6697

NANGNURI NIKHITHA : 22Q91A66A2

Under the guidance of

DR. K. MADAN MOHAN

Assistant Professor, Dept. of CSE(AIML)



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING ARTIFICIAL
INTELLIGENCE AND MACHINE LEARNING
MALLA REDDY COLLEGE OF ENGINEERING**

(Approved by AICTE-Permanently Affiliated to JNTU-Hyderabad) Accredited by NBA & NAAC,

Recognized section 2(f) & 12(B) of UGC New Delhi ISO

9001:2015 certified Institution

Maisammaguda, Dhulapally (Post via Kompally), Secunderabad 500100

2023-2024

MALLA REDDY COLLEGE OF ENGINEERING

(Approved by AICTE-Permanently Affiliated to JNTU-Hyderabad)

Accredited by NBA & NAAC, Recognized section 2(f) & 12(B) of UGC New Delhi ISO
9001:2015 certified Institution

Maisammaguda, Dhulapally (Post via Kompally), Secunderabad 500100

DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING



CERTIFICATE

This is to certify that the Real Time Research Project report on “SUBWAY PASSENGER FLOW FORECASTING USING MACHINE LEARNING” is successfully done by the following students of Department of Computer Science & Engineering (Artificial Intelligence and Machine Learning) of our college in partial fulfilment of the requirement for the award of B. Tech degree in the year 2023-2024. The results embodied in this report have not been submitted to any other University for the award of any diploma or degree.

GAJULA NIKHIL	: 22Q91A6685
HARSH R BAGTHARIA	: 22Q91A6691
MADHARAPU NITHIN KUMAR	: 22Q91A6697
NANGNURI NIKHITHA	: 22Q91A66A2

Submitted for the viva-voce examination held on: _____

INTERNAL GUIDE

Mr. M. Sakthivel
Assistant Professor

HOD

Mr. Vignesh Janarthanan
Assistant Professor

Internal Examiner

External Examiner

DECLARATION

We, Gajula Nikhil, Harsh R Bagtharia, Madharapu Nithin Kumar, R Nangnuri Nikhitha with Regd.no. 22Q91A6685, 22Q91A6691, 22Q91A6697, 22Q91A66A2 are hereby declaring that the real time research project report entitled “SUBWAY PASSENGER FLOW FORECASTING USING MACHINE LEARNING” has done by us under the guidance of Dr. K Madan Mohan, Assistant Professor, Department of Computer Science Engineering (Artificial Intelligence and Machine Learning) is submitted in the partial fulfilment of the requirements for the award of degree of **BACHELOR OF TECHNOLOGY in COMPUTER SCIENCE AND ENGINEERING ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING.**

The Results embedded in this project report have not been submitted to any other University or institute for the award of any degree or diploma.

Signature of the Candidate

GAJULA NIKHIL	: 22Q91A6685
HARSH R BAGTHARIA	: 22Q91A6691
MADHARAPU NITHIN KUMAR	: 22Q91A6697
NANGNURI NIKHITHA	: 22Q91A66A2

DATE:

PLACE: Maisammaguda

ACKNOWLEDGEMENT

First and foremost, we would like to express our immense gratitude towards our institution Malla Reddy College of Engineering, which helped us to attain profound technical skills in the field of Computer Science & Engineering, there by fulfilling our most cherished goal.

We are pleased to thank **Sri Ch. Malla Reddy**, our Founder, Chairman **MRGI**, **Sri Ch. Mahender Reddy**, Secretary, **MRGI** for providing this opportunity and support throughout the course.

It gives us immense pleasure to acknowledge the perennial inspiration of **Mr. Vignesh Janarthanan** HOD, **Mr. C. Dinesh** Project Coordinator of Computer Science Engineering (Artificial Intelligence and Machine Learning) Department for their inspiration adroit guidance and constructive criticism for successful completion of our degree.

We would like to thank **Dr. K. Madan Mohan** Assistant Professor our internal guide, for his valuable suggestions and guidance during the exhibition and completion of this project.

Finally, we avail this opportunity to express our deep gratitude to all staff who have contribute their valuable assistance and support making our project success.

GAJULA NIKHIL	: 22Q91A6685
HARSH R BAGTHARIA	: 22Q91A6691
MADHARAPU NITHIN KUMAR	: 22Q91A6697
NANGNURI NIKHITHA	: 22Q91A66A2

ABSTRACT

With the escalating demand for efficient urban transportation systems, accurate forecasting of subway passenger traffic has become paramount for both commuters and city planners. This project introduces a Django-based web application designed to predict subway traffic using various input parameters such as stop name, line, hour, and day. The core of the application is powered by a pre-trained machine learning model capable of forecasting the number of passengers entering and leaving each subway stop.

The significance of this project lies in its potential to enhance passenger experience by providing real-time traffic predictions, which can help commuters plan their journeys more effectively and avoid congested stations. Additionally, the application supports subway operators in optimizing their operations, such as scheduling trains more efficiently and managing resources better to accommodate passenger flows.

Key features of the application include an intuitive user interface that allows users to input their desired parameters through a simple form. The backend processes these inputs and utilizes the machine learning model to generate accurate predictions, which are then displayed in a clear and accessible format. This not only makes the tool user-friendly but also ensures that the information is easily interpretable for users of all backgrounds.

Moreover, the application's predictive capabilities are built upon a robust dataset encompassing historical passenger flow data, line-specific information, and temporal factors. By integrating these diverse data sources, the machine learning model achieves high accuracy in its forecasts. This integration of data is facilitated by advanced techniques such as embedding methods and hierarchical attention mechanisms, which allow the model to capture complex patterns and relationships within the data.

Through this project, we aim to contribute to the broader vision of smart city initiatives by enabling more informed decision-making in urban transit systems. The insights generated by the application can be invaluable for municipal authorities in designing better infrastructure and policies to meet the growing needs of urban populations. Furthermore, the visualization of predicted traffic patterns helps in understanding the impact of different factors on subway usage, thereby supporting strategic planning and development.

TABLE OF CONTENTS

ACKNOWLEDGEMENT	I
ABSTRACT.....	II
TABLE OF CONTENTS	III
CHAPTER - 1 INTRODUCTION.....	1
INTRODUCTION	2
METHODOLOGY	3
PURPOSE OF PROJECT	5
SCOPE	6
CHAPTER - 2 LITERATURE SURVEY	8
LITERATURE SURVEY	9
CHAPTER - 3 SYSTEM ANALYSIS	10
EXISTING SYSTEM	11
PROPOSED SYSTEM.....	13
CHAPTER - 4 SYSTEM REQUIREMENTS.....	15
SOFTWARE REQUIREMENTS SPECIFICATIONS	16
REQUIREMENT SPECIFICATION DOCUMENT	16
FUNCTIONAL REQUIREMENTS	16
NON-FUNCTIONAL REQUIREMENTS	16
SOFTWARE REQUIREMENTS	17
CHAPTER - 5 SOFTWARE ARCHITECTURE	18
SOFTWARE ARCHITECTURE OVERVIEW	19
CHAPTER - 6 SYSTEM DESIGN	21
<i>Class Diagram</i>	22
<i>Sequence Diagram</i>	23
<i>Activity Diagram</i>	24
<i>Use Case Diagram</i>	25
<i>Data Flow Diagram</i>	26
<i>State Diagram</i>	27
<i>Component Diagram</i>	27
<i>Flowchart Diagram</i>	28
<i>Deployment Diagram</i>	29

CHAPTER - 7 SYSTEM IMPLEMENTATION	30
TECHNOLOGIES	31
SOURCE CODE.....	32
WORKING OF SYSTEM	42
SOFTWARE INSTALLATION	43
ALGORITHMS	44
MODULES	45
<i>Detailed Description of Each Module</i>	<i>46</i>
CHAPTER - 8 TESTING.....	47
INTRODUCTION TO TESTING	48
<i>Introduction to Testing</i>	<i>48</i>
WHO DOES TESTING?	48
FUNCTIONAL TESTING	48
SOFTWARE TESTING LIFECYCLE.....	49
TESTING THE AI MODEL FOR SUBWAY PASSENGER FLOW FORECASTING	49
CHAPTER - 9 SCREENSHOTS & RESULT	51
SCREENSHOTS & RESULTS	52
<i>User Interface Description</i>	<i>52</i>
<i>Expected Results</i>	<i>54</i>
<i>Example Scenario</i>	<i>57</i>
CHAPTER - 10 ADVANTAGES AND LIMITATIONS	59
ADVANTAGES:	60
LIMITATIONS:	60
CHAPTER - 11 FUTURE ENHANCEMENTS	62
FUTURE ENHANCEMENTS	63
CHAPTER - 12 CONCLUSION	63
CONCLUSION	64
REFERENCES	66
REFERENCES.....	67
<i>References.....</i>	<i>67</i>
<i>Libraries and Frameworks</i>	<i>67</i>
<i>Dataset.....</i>	<i>67</i>

LIST OF FIGURES

Figure 1: Existing System	11
Figure 2: Proposed System.....	13
Figure 3: Software Architecture	19
Figure 4: Class Diagram.....	22
Figure 5: Sequence Diagram	23
Figure 6: Activity Diagram	24
Figure 7: Use Case Diagram.....	25
Figure 8: Data Flow Diagram.....	26
Figure 9: State Diagram.....	27
Figure 10: Component Diagram.....	27
Figure 11: Flowchart Diagram	28
Figure 12: Deployment Diagram.....	29
Figure 13: Homepage	52
Figure 14: Form Inputs.....	53
Figure 15: Submit Button	53
Figure 16: Result Page	54
Figure 17: Result Page (Forecasted Traffic)	54
Figure 18: Result Page (Day Prediction Bar Graph)	55
Figure 19 Result Page (Week Prediction Bar Graph).....	55
Figure 20: Result Page (Day Prediction Pie Chart).....	56
Figure 21: Result Page (Week Prediction Pie Chart)	56
Figure 22: Example Scenario (Home Page)	57
Figure 23: Example Scenario (Complete Result Page)	57
Figure 24: Example Scenario (Complete Result Page)	58

CHAPTER - 1

INTRODUCTION

INTRODUCTION

Urbanization and the rapid expansion of rail transit systems have led to a significant increase in subway ridership in cities worldwide. As more people rely on subways for their daily commute, the need for efficient and reliable transportation systems becomes critical. Accurate forecasting of subway passenger flow is essential for managing the operational aspects of transit systems, optimizing resource allocation, and enhancing the overall commuter experience.

Traditional methods of passenger flow prediction often fall short in addressing the dynamic and complex nature of urban transit systems. With advancements in data science and machine learning, there is a growing opportunity to develop more sophisticated models that can accurately predict passenger traffic based on various factors.

This project presents a Django-based web application designed to predict subway passenger traffic using a pre-trained machine learning model. By allowing users to input parameters such as stop name, line, hour, and day, the application provides real-time forecasts of the number of passengers entering and leaving each subway stop. These predictions are crucial for improving journey planning for commuters and optimizing operational decisions for transit authorities.

The integration of historical data, line-specific information, and temporal factors into the prediction model enables a comprehensive analysis of passenger flow patterns. Advanced techniques such as embedding methods and hierarchical attention mechanisms enhance the model's accuracy, making it a valuable tool for both daily operations and long-term strategic planning in urban transit systems.

METHODOLOGY

The methodology of the subway passenger flow forecasting project involves several key steps, from data preprocessing to deploying the predictive model within a Django web application. Below is a detailed breakdown of the methodology:

1. Data Preprocessing

1. Loading the Dataset:

- **Library Used:** Pandas
- The dataset containing historical subway passenger flow data is loaded into the system using Pandas. This dataset includes information on passenger counts entering and leaving each station at various times.

2. Data Cleaning and Sorting:

- The loaded dataset is cleaned to handle missing values, outliers, and inconsistencies. It is then sorted based on relevant features such as station name, line, hour of the day, and day of the week.

3. Feature Engineering:

- Additional features are derived from the existing data to enhance the predictive power of the machine learning models. This may include time-based features (e.g., peak hours, weekdays vs. weekends) and external factors (e.g., weather conditions).

2. Model Training

1. Choosing Algorithms:

- **Libraries Used:** Scikit-learn, XGBoost
- Two main machine learning algorithms are utilized for training:
 - **Scikit-learn:** Provides various algorithms for initial modeling and evaluation.
 - **XGBoost:** A gradient boosting algorithm known for its high performance and accuracy in prediction tasks.

2. Training the Model:

- The cleaned and engineered dataset is split into training and testing sets.
- Models are trained on the training set using both Scikit-learn and XGBoost algorithms. Hyperparameter tuning is performed to optimize model performance.

3. Model Evaluation:

- The trained models are evaluated on the testing set to assess their accuracy and generalizability. Metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared are used for evaluation.

4. Model Persistence:

- **Library Used:** Joblib
- The best-performing model is serialized and saved using Joblib for later use in predictions. This ensures that the model can be efficiently loaded and utilized within the web application.

3. Web Application Development

1. Django Framework:

- **Library Used:** Django
- A Django web application is developed to serve as the interface for users to input station details and receive passenger flow predictions.

2. Creating Models:

- A Django model (`TrafficPrediction`) is defined to handle the storage and management of prediction data. This model includes fields for station name, line, hour, day, and predicted passenger counts (both entering and leaving).

3. User Input Form:

- A form is created in Django to allow users to input details such as station name, line, hour, and day. This form collects the necessary input data for making predictions.

4. Prediction Logic:

- Upon form submission, the input data is processed and fed into the trained AI model (loaded using Joblib) to generate predictions.
- The predicted passenger counts are then stored in the `TrafficPrediction` model for further use and display.

4. Data Visualization

1. Chart.js Integration:

- **Library Used:** Chart.js
- The predicted passenger flow data is visualized using Chart.js. Bar charts and pie charts are created to display:
 - The number of passengers arriving and leaving at a specific hour.
 - Hourly predictions for each day of the week.
 - Overall daily trends.

2. Interactive Graphs:

- The visualizations are interactive, allowing users to explore the data and understand passenger flow patterns intuitively.

PURPOSE OF PROJECT

The purpose of this project is to develop a comprehensive and user-friendly web application for predicting subway passenger flow using machine learning techniques. This project aims to address the challenges of managing urban transit systems by providing accurate and real-time forecasts of passenger traffic at various subway stations. By leveraging historical data and advanced predictive models, the application serves multiple stakeholders, including commuters, transit authorities, and urban planners. The key objectives of the project are as follows:

1. Enhance Commuter Experience:

- **Objective:** To provide commuters with reliable information about passenger traffic at subway stations, enabling them to make informed decisions about their travel plans.
- **Benefit:** Helps reduce wait times, avoid overcrowded trains, and improve overall commuting experience.

2. Optimize Transit Operations:

- **Objective:** To assist transit authorities in managing and optimizing subway operations by forecasting passenger volumes at different times and locations.
- **Benefit:** Allows for better resource allocation, schedule adjustments, and infrastructure planning, leading to more efficient and reliable transit services.

3. Support Urban Planning:

- **Objective:** To provide valuable insights for urban planners and city officials to understand and anticipate transit demand patterns.
- **Benefit:** Facilitates informed decision-making for future transit developments, infrastructure investments, and policy-making.

4. Utilize Advanced Machine Learning:

- **Objective:** To apply state-of-the-art machine learning algorithms (Scikit-learn, XGBoost) to predict subway passenger flow with high accuracy.
- **Benefit:** Ensures that the predictions are reliable and can adapt to changing patterns and external factors.

5. Integrate Data Visualization:

- **Objective:** To present prediction results in an intuitive and interactive manner using Chart.js.
- **Benefit:** Enhances user understanding of traffic patterns through visual representations such as bar charts and pie charts, making complex data accessible and actionable.

6. Promote Real-Time Decision Making:

- **Objective:** To provide real-time processing and prediction capabilities, allowing users to receive immediate feedback based on current input data.
- **Benefit:** Enables timely decisions and adjustments, improving the responsiveness of both commuters and transit authorities to dynamic conditions.

SCOPE

The scope of the subway passenger flow forecasting project encompasses several key aspects, focusing on both technical implementation and practical application within urban transit management. Here's an outline of the project's scope:

1. Data Collection and Preprocessing:

- **Data Sources:** Utilizing historical data on subway passenger flow, including details such as station names, subway lines, time of day, and day of the week.
- **Data Cleaning:** Handling missing values, outliers, and inconsistencies to ensure data quality and reliability.
- **Feature Engineering:** Deriving additional features from the dataset to enhance prediction accuracy, such as time-based trends and external factors.

2. Machine Learning Model Development:

- **Algorithm Selection:** Implementing machine learning algorithms, specifically Scikit-learn and XGBoost, for predicting passenger flow based on historical patterns.
- **Model Training:** Splitting the dataset into training and testing sets, optimizing hyperparameters, and evaluating model performance using metrics like Mean Absolute Error (MAE) and R-squared.
- **Model Persistence:** Saving and loading the trained model using Joblib for efficient deployment within the Django web application.

3. Web Application Development with Django:

- **User Interface Design:** Creating a user-friendly interface where users can input details such as station name, subway line, hour, and day to obtain predictions.
- **Form Handling:** Developing forms in Django to capture user inputs and validate data before processing.
- **Integration of Predictive Models:** Incorporating the trained machine learning model to generate real-time predictions based on user inputs.
- **Visualization:** Utilizing Chart.js to visualize prediction results through interactive graphs and charts, displaying passenger flow trends by hour and day.

4. Deployment and Scalability:

- **System Architecture:** Deploying the Django web application on scalable cloud platforms such as AWS or Heroku to handle varying loads and ensure reliability.
- **Performance Optimization:** Monitoring application performance and implementing optimizations to maintain responsiveness during peak usage periods.
- **Security Considerations:** Implementing security measures within the application to protect user data and ensure secure interactions.

5. User Engagement and Feedback:

- **User Testing:** Conducting usability testing to gather feedback on the application's interface, functionality, and performance.
- **Iterative Improvement:** Incorporating user feedback to refine features, enhance usability, and improve overall user experience.
- **Documentation and Support:** Providing comprehensive documentation for users and administrators, including instructions for usage, troubleshooting, and updates.

6. Future Enhancements:

- **Real-Time Data Integration:** Enhancing predictions by integrating real-time data sources such as weather conditions and transit schedules.
- **Advanced Analytics:** Implementing advanced analytics features for deeper insights into passenger behavior and traffic patterns.
- **Geospatial Visualization:** Expanding visualization capabilities with geospatial mapping to visualize passenger flow across subway networks.

CHAPTER - 2

LITERATURE SURVEY

LITERATURE SURVEY

1. Model Generation and Origin-Destination (OD) Matrix Forecasting for Metro Systems

- Authors: Daniel Sanz Sobrino, Javier Andión, Juan C. Dueñas, José M. Del Álamo, Felix Cuadrado
- Conference: 2023 7th International Young Engineers Forum (YEF-ECE), pp.112-117, 2023
- Abstract: This paper addresses the problem of forecasting Origin-Destination (OD) matrices for metro systems. The authors propose a model that generates OD matrices by incorporating various data sources and leveraging machine learning techniques. The model focuses on predicting the flow of passengers between different stations, considering historical data and external factors. Experimental results show that the proposed method improves the accuracy of OD matrix forecasting, which is crucial for efficient metro system planning and operations.

2. Dynamic Spatio-Temporal Multi-Scale Representation for Bus Ridership Prediction

- Authors: Lilan Peng, Xiu Wang, Hongchun Lu, Xiangyu Guo, Tianrui Li, Shenggong Ji
- Conference: 2023 International Joint Conference on Neural Networks (IJCNN), pp.1-9, 2023
- Abstract: This paper introduces a novel approach for predicting bus ridership using dynamic spatio-temporal multi-scale representations. The model integrates spatial and temporal information at multiple scales to capture complex patterns in bus ridership data. By employing deep learning techniques, the model effectively handles the dynamic nature of bus passenger flows, leading to improved prediction accuracy. The proposed method is validated through extensive experiments on real-world bus ridership datasets.

3. Spatio-Temporal Dynamic Graph Relation Learning for Urban Metro Flow Prediction

- Authors: Peng Xie, Minbo Ma, Tianrui Li, Shenggong Ji, Shengdong Du, Zeng Yu, Junbo Zhang
- Journal: IEEE Transactions on Knowledge and Data Engineering, vol.35, no.10, pp.9973-9984, 2023
- Abstract: This paper presents a spatio-temporal dynamic graph relation learning model for predicting urban metro passenger flows. The model constructs dynamic graphs to represent relationships between metro stations over time and employs graph neural networks to learn these relationships. The approach captures both spatial and temporal dependencies, leading to accurate metro flow predictions. The model's performance is demonstrated on real-world metro datasets, showing significant improvements over traditional methods.

4. MVOPFAD: Multiview Online Passenger Flow Anomaly Detection

- Authors: Erlong Tan, Haoyang Yan, Kaiqi Zhao, Xiaolei Ma, Zhenliang Ma, Yuchuan Du
- Journal: IEEE Sensors Journal, vol.24, no.9, pp.14668-14681, 2024
- Abstract: This paper introduces MVOPFAD, a multiview online passenger flow anomaly detection system for urban transit networks. The system integrates multiple data sources, including passenger flow data and external factors, to detect anomalies in real-time. By using a multiview learning approach, MVOPFAD captures diverse perspectives on passenger flow patterns, enabling timely and accurate detection of anomalies. The system's effectiveness is validated through experiments on large-scale urban transit datasets.

5. COV-STFormer for Short-Term Passenger Flow Prediction During COVID-19 in Urban Rail Transit Systems

- Authors: Shuxin Zhang, Jinlei Zhang, Lixing Yang, Chengcheng Wang, Ziyao Gao
- Journal: IEEE Transactions on Intelligent Transportation Systems, vol.25, no.5, pp.3793-3811, 2024
- Abstract: This paper proposes COV-STFormer, a model designed for short-term passenger flow prediction in urban rail transit systems during the COVID-19 pandemic. The model incorporates spatio-temporal transformers to capture the impact of the pandemic on passenger flow patterns. By integrating COVID-19-related data and traditional transit data, COV-STFormer provides accurate and robust predictions. Experimental results on real-world datasets demonstrate the model's effectiveness in handling the unique challenges posed by the pandemic.

CHAPTER - 3

SYSTEM ANALYSIS

EXISTING SYSTEM

In traditional approaches to subway passenger flow forecasting, methods typically rely on simplistic statistical models or manual data collection processes. Statistical models often use basic time-series analysis or regression techniques, which may not effectively capture the complexities of urban transit dynamics, such as peak hours, seasonal variations, and external factors like events or weather conditions. Manual data collection involves labor-intensive processes of periodic observations at subway stations, which are prone to errors and may not provide real-time insights needed for proactive management. Visualization of data in the existing system often lacks interactivity and dynamic updating capabilities, limiting its utility in understanding and responding to changing passenger flow patterns effectively.

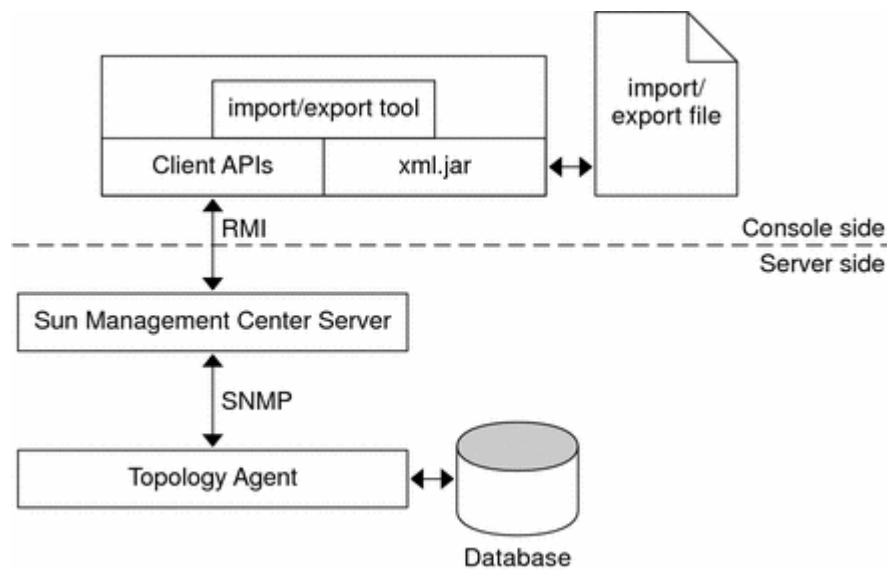


Figure 1: Existing System

Challenges in the Existing System

- Scalability Issues:

As the number of subway stations and the volume of passengers increase, existing systems may struggle to scale and maintain accuracy across a larger dataset.

- Data Quality and Availability:

Ensuring high-quality and up-to-date data is a significant challenge. Existing systems often lack mechanisms to automatically update and validate the data they rely on.

- User Experience:

Many current systems are not user-friendly, making it difficult for transit operators and commuters to interact with the system and understand the predictions.

Disadvantages of Existing System

While existing systems for predicting subway passenger flow offer foundational insights, they come with several significant disadvantages that limit their effectiveness and reliability. These disadvantages are critical to understand when considering improvements and innovations in transit management systems.

1. Limited Accuracy:

Simplistic Models: Most existing systems use basic statistical methods such as linear regression or moving averages, which cannot capture the complex patterns and fluctuations in passenger flow. This results in less accurate predictions.

Static Data Utilization: Reliance on historical data without dynamic updates leads to outdated predictions that do not reflect current conditions.

2. Inadequate Data Integration:

Narrow Data Scope: Existing systems often do not incorporate a broad range of relevant data, such as real-time passenger counts, weather conditions, special events, and social factors. This limits their ability to make comprehensive predictions.

Lack of Contextual Factors: Many systems fail to account for external factors like public holidays, sporting events, or other large gatherings that significantly impact passenger flow.

3. Lack of Real-Time Capabilities:

Delayed Predictions: Without real-time data processing, these systems cannot provide timely predictions. This is crucial for managing unexpected surges in passenger numbers during peak hours or emergency situations.

Static Reporting: Existing systems often produce static reports rather than continuous, real-time updates, making them less useful for immediate decision-making.

4. Poor Scalability:

Inflexibility: As the transit network grows and the number of passengers increases, existing systems struggle to scale. This can lead to performance issues and decreased prediction accuracy.

Maintenance Challenges: Updating and maintaining these systems to handle larger datasets and new stations is often complex and resource-intensive.

5. User Interaction Limitations:

Non-Interactive Interfaces: Many systems do not offer user-friendly interfaces for inputting specific parameters or interacting with the prediction model. This limits their accessibility and usefulness for end-users.

Lack of Customization: Users cannot easily customize the input variables to get tailored predictions, reducing the practical applicability of the system.

6. Data Quality Issues:

Inconsistent Data: Existing systems often struggle with maintaining high-quality, consistent data inputs, leading to unreliable predictions.

Manual Updates: Many systems require manual updates for data, which can introduce errors and delays.

PROPOSED SYSTEM

The proposed subway passenger flow forecasting system aims to leverage advanced machine learning algorithms and modern data processing techniques to overcome the limitations of the existing system. By using Scikit-learn and XGBoost, the system can train robust predictive models on comprehensive datasets of historical passenger flow data. This approach allows for the integration of multiple variables such as station-specific attributes, time of day, day of the week, and external factors, enhancing the accuracy and granularity of predictions. Data preprocessing tasks, handled efficiently with Pandas, ensure data quality and relevance for model training.

The deployment of trained AI models within a Django web application environment, facilitated by Joblib, enables real-time predictions based on user inputs. This provides commuters and transit authorities with immediate access to forecasted passenger counts, visualized through interactive charts and graphs using Chart.js. These visualizations offer insights into hourly and daily passenger flow trends, empowering stakeholders to make informed decisions regarding resource allocation, service planning, and operational adjustments. Real-time data integration further enhances the system's adaptability to changing transit conditions, ensuring that forecasts remain responsive and accurate in dynamic urban environments.

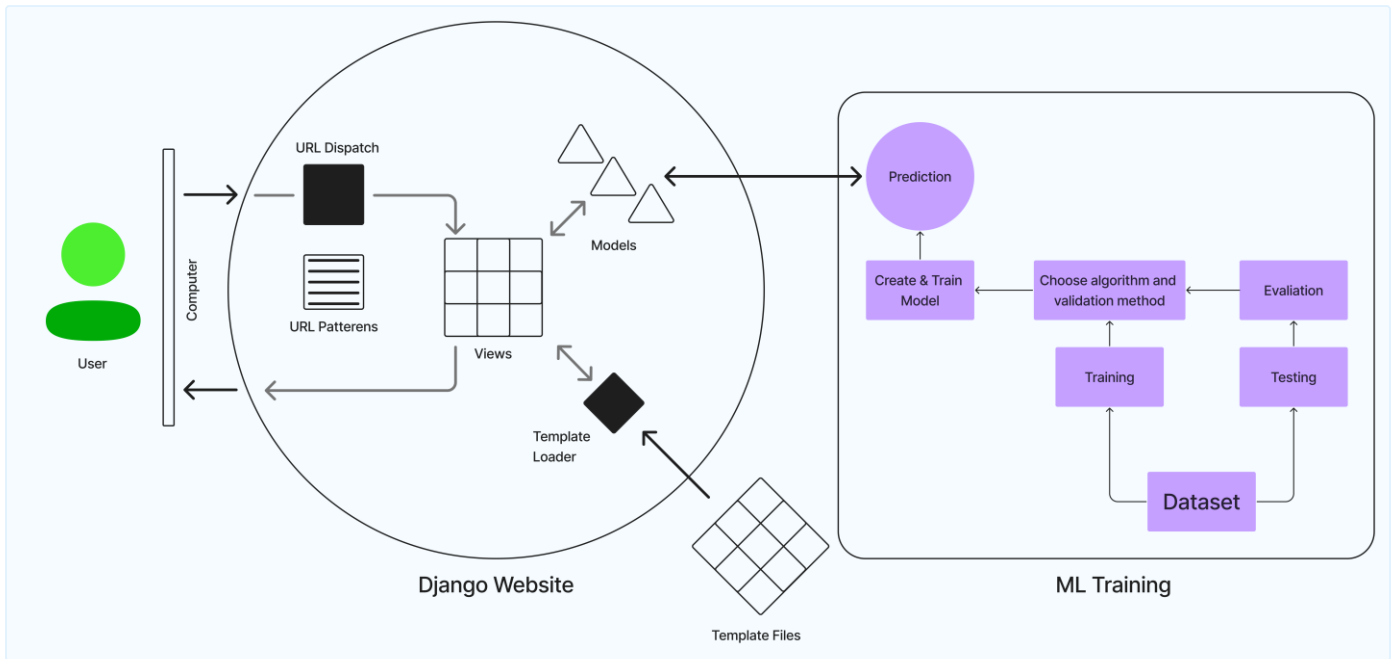


Figure 2: Proposed System

Key Features of the Proposed System

1. Advanced Machine Learning Model:

Dynamic Prediction: The system uses a pre-trained machine learning model capable of adapting to new data and trends, providing more accurate and up-to-date predictions.

Complex Data Integration: By incorporating diverse datasets including historical passenger counts, weather conditions, event schedules, and temporal factors, the model can capture complex interactions and patterns in passenger flow.

2. Real-Time Data Processing:

Continuous Updates: The system processes real-time data, ensuring that predictions reflect current conditions and allow for immediate response to sudden changes in passenger traffic.

Responsive Interface: Users can input parameters at any time and receive real-time predictions, making the tool highly interactive and responsive.

3. User-Friendly Interface:

Intuitive Forms: The application provides a simple, intuitive form for users to input stop name, line, hour, and day details. The design prioritizes ease of use and accessibility.

Clear Visualization: Predictions are displayed in a clear and accessible format, helping users easily interpret the data and make informed decisions.

4. Scalability and Flexibility:

Scalable Architecture: The system is built to handle an increasing number of stations and larger datasets, ensuring it can scale with the growth of the subway network.

Flexible Data Integration: New data sources can be easily integrated into the model, allowing for continuous improvement and adaptation.

5. Comprehensive Data Utilization:

Embedding Techniques: The use of embedding methods allows the system to combine different kinds of data effectively, enhancing the accuracy of predictions.

Hierarchical Attention Mechanism: This mechanism models the hierarchical relationship between subway lines and stations, further improving prediction accuracy.

6. Enhanced Predictive Insights:

Impact Analysis: The system can visualize the impact of different stations and external factors on traffic, providing valuable insights for transit authorities and city planners.

Customizable Predictions: Users can customize input variables to obtain tailored predictions, making the system more practical and relevant to individual needs.

Benefits of the Proposed System

- **Improved Accuracy:** By utilizing advanced machine learning techniques and integrating diverse data sources, the proposed system provides more accurate and reliable predictions compared to existing models.
- **Real-Time Responsiveness:** The ability to process and reflect real-time data ensures that predictions are timely and relevant, allowing for better management of sudden changes in passenger flow.
- **Enhanced User Experience:** The user-friendly interface and clear visualization of data make the system accessible and easy to use for both commuters and transit authorities.
- **Operational Efficiency:** Accurate predictions enable better resource allocation, optimized train schedules, and improved overall efficiency in subway operations.
- **Strategic Planning Support:** The insights generated by the system can inform long-term planning and infrastructure development, supporting the broader vision of smart city initiatives.

CHAPTER - 4

SYSTEM REQUIREMENTS

SOFTWARE REQUIREMENTS SPECIFICATIONS

The Software Requirements Specification (SRS) document for the "Subway Passenger Flow Forecasting Using Machine Learning" project serves to outline the functional and non-functional requirements, as well as the software and hardware prerequisites necessary for its development and deployment. This document bridges the communication gap between stakeholders, ensuring that client and user needs are accurately captured and translated into actionable development guidelines.

Requirement Specification Document

The SRS document defines the nature and scope of the project, providing a comprehensive overview of its purpose, functional capabilities, and quality attributes. It establishes a framework for software development, ensuring clarity and alignment among all parties involved.

Functional Requirements

Functional requirements specify the specific behaviors and functionalities that the system must exhibit to meet the needs of its users and stakeholders. These requirements dictate the system's inputs, processing logic, and expected outputs.

- **FR1: User Input Form**
 - Users can input details of a subway station including stop name, line, hour, and day via a web form.
 - Inputs trigger the AI model to predict the number of passengers arriving and exiting the station for the specified parameters.
- **FR2: Prediction Visualization**
 - Predicted passenger counts are visualized as bar graphs and pie charts.
 - Graphs depict passenger flow for each hour of the day and day of the week using Chart.js, facilitating easy interpretation and analysis.

Non-Functional Requirements

Non-functional requirements define the quality attributes and constraints of the system, focusing on aspects such as performance, usability, reliability, and security.

- **NFR1: Performance**
 - The system should process predictions in real-time with minimal latency, ensuring responsiveness even under high user load.
- **NFR2: Usability**
 - The user interface should be intuitive and user-friendly, allowing users to easily interact with prediction results and visualizations.
- **NFR3: Reliability**
 - The system should be robust and reliable, capable of handling continuous operation without unexpected downtime.
- **NFR4: Security**
 - Strong authentication and authorization mechanisms should be implemented to protect user data and system integrity.

Software Requirements

The software requirements specify the necessary tools, frameworks, libraries, and environments required for developing and deploying the system.

- **Operating System:**
 - Development: Linux (Ubuntu 20.04 or later), macOS, or Windows 10.
 - Production: Linux (preferred for stability and security).
- **Programming Languages:**
 - Python 3.8 or later for implementing Django backend and machine learning components.
- **Frameworks and Libraries:**
 - Django 2.2 for web application development.
 - Pandas for data manipulation and preprocessing.
 - Scikit-learn and XGBoost for training machine learning models.
 - Joblib for saving and loading trained models.
- **Database:**
 - Utilizes datasets NYC Subway Traffic Data 2017-2021, from Kaggle.com.
- **Front-End Technologies:**
 - HTML/CSS, JavaScript for frontend development.
 - Chart.js for interactive and dynamic visualization of prediction results.

Environmental Conditions

- **Development Environment:**
 - Should support Python, Django, and related libraries.
 - Testing should replicate the production environment to ensure accurate validation of system functionalities.
- **Production Environment:**
 - High availability setup with proper configurations for real-time data processing and user interaction.

Safety and Security Requirements

- **Authentication and Authorization:**
 - Implement robust mechanisms to secure user interactions and prevent unauthorized access.
- **Security Audits:**
 - Regular audits and updates to address vulnerabilities and maintain system security.

Software Quality Attributes

- **Maintainability:**
 - Modular design and clear documentation to facilitate future updates and maintenance.
- **Scalability:**
 - Horizontal scalability to handle increased data volumes and user traffic.
- **Usability:**
 - Intuitive user interface with clear visualizations to aid in decision-making.
- **Performance:**
 - Efficient data processing and analysis to ensure timely and accurate prediction of passenger flow.

CHAPTER - 5

SOFTWARE ARCHITECTURE

SOFTWARE ARCHITECTURE OVERVIEW

The architecture diagram provided illustrates the interaction and data flow between different components in a subway passenger flow forecasting system. Below is a detailed explanation, similar to what would be found in a technical documentation:

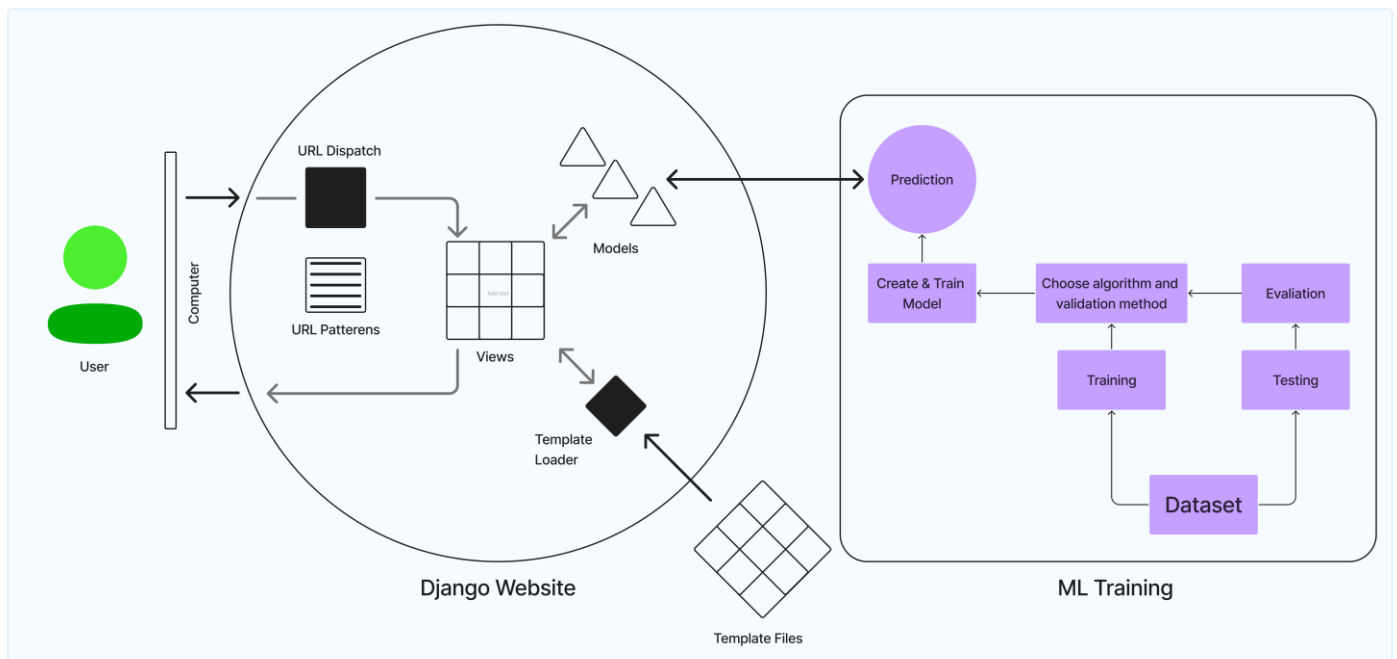


Figure 3: Software Architecture

Components of the System:

1. Models:

- **TrafficPrediction Model:** Defines the structure for storing subway passenger flow predictions. It includes fields for stop name, line number, hour, day, and predicted counts for entering and leaving passengers.

```
class TrafficPrediction(models.Model):
    stop_name = models.CharField(max_length=100)
    line = models.IntegerField()
    hour = models.IntegerField()
    day = models.IntegerField()
    enter_count = models.FloatField()
    leave_count = models.FloatField()

    def __str__(self):
        return f"{self.stop_name} {self.hour}:{self.day}"
```

2. Views:

- **predict_traffic View:** Handles user requests for predicting subway passenger flows. It loads a pre-trained machine learning model, processes user input, performs predictions, and stores results in the database.

```
def predict_traffic(request):
    if request.method == 'POST':
        # Process user input, load model, perform predictions, and
        # store results
        ...
        return render(request, 'traffic/prediction_result.html',
            {'prediction': prediction, 'week_results': week_results,
            'day_results': day_results})
    return render(request, 'traffic/predict_traffic.html')
```

3. Templates:

- **predict_traffic.html**: Provides a form for users to input details such as stop name, line, hour, and day.
- **prediction_result.html**: Displays prediction results using bar graphs and pie charts generated with Chart.js to visualize passenger flows.

4. Machine Learning Model:

- **sklearn and XGBoost**: Used for training and deploying the machine learning model. Predictions are made based on input data (stop name, line, hour, day) provided by users via the Django application.

This software architecture overview highlights the structured approach of using Django's MVT pattern along with machine learning integration and visualization techniques to deliver a robust and user-friendly subway passenger flow forecasting system.

CHAPTER - 6

SYSTEM DESIGN

Class Diagram

Class Diagram Overview: The class diagram depicts the static structure of the system, showing classes, their attributes, methods, and relationships. Here’s an overview based on your project:

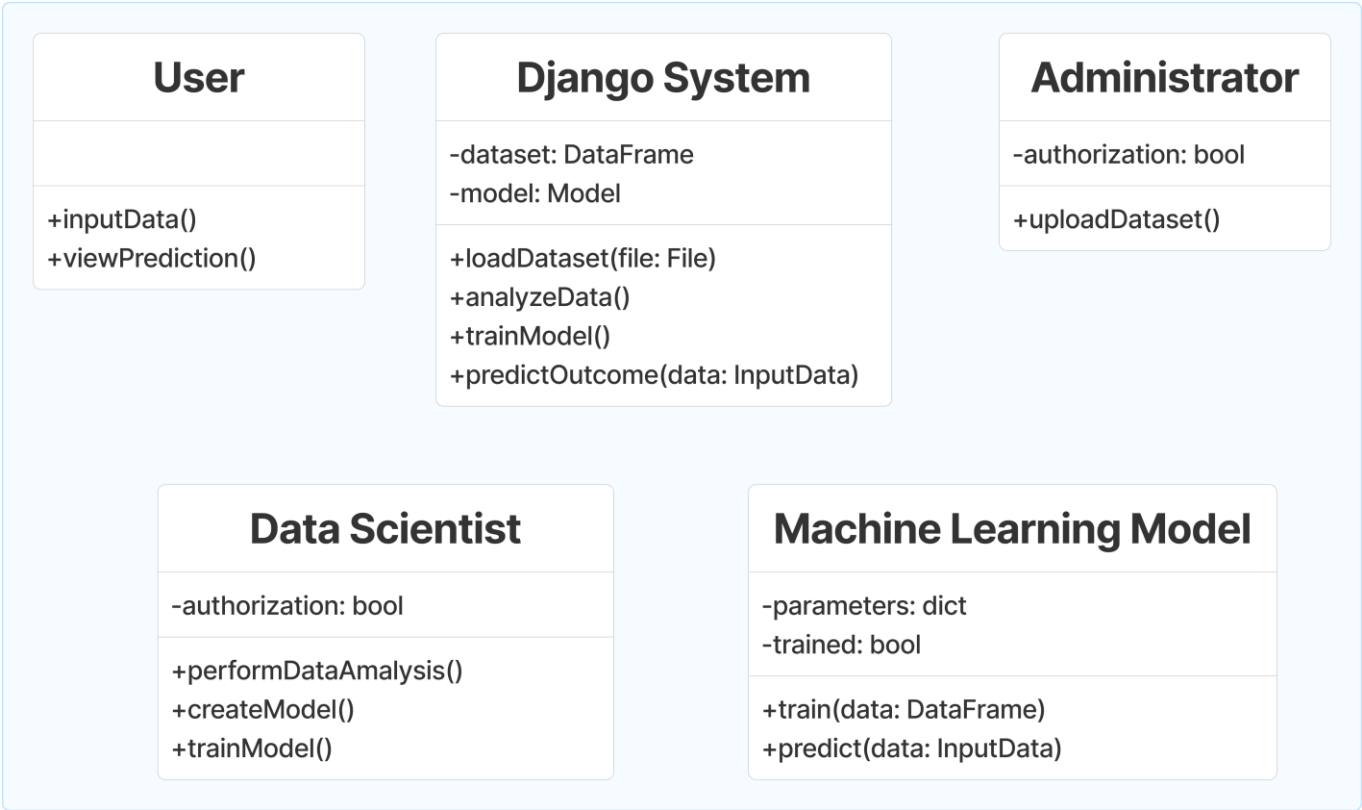


Figure 4: Class Diagram

- **Entities:**
 - **User:** Represents individuals interacting with the system.
 - **Administrator:** Manages system administration tasks such as dataset management.
 - **Data Scientist:** Responsible for developing and managing machine learning models.
- **Attributes and Methods:**
 - **Dataset:** Represents the dataset used for analysis and training.
 - **Model:** Represents the machine learning model used for prediction.
 - **InputData:** A class or structure defining the format of input data for predictions.
 - **ViewPrediction():** Method allowing users to view prediction results.
 - **UploadDataset():** Method enabling administrators to upload datasets into the system.
 - **LoadDataset(file):** Method for loading datasets into the system for analysis and training.
 - **TrainModel():** Method for training the machine learning model using the dataset.
 - **PredictOutcome(data):** Method for predicting outcomes based on input data.

Sequence Diagram

Sequence Diagram Overview: A sequence diagram illustrates the interactions between objects in a sequential order, showing the flow of messages and actions between different components of the system.

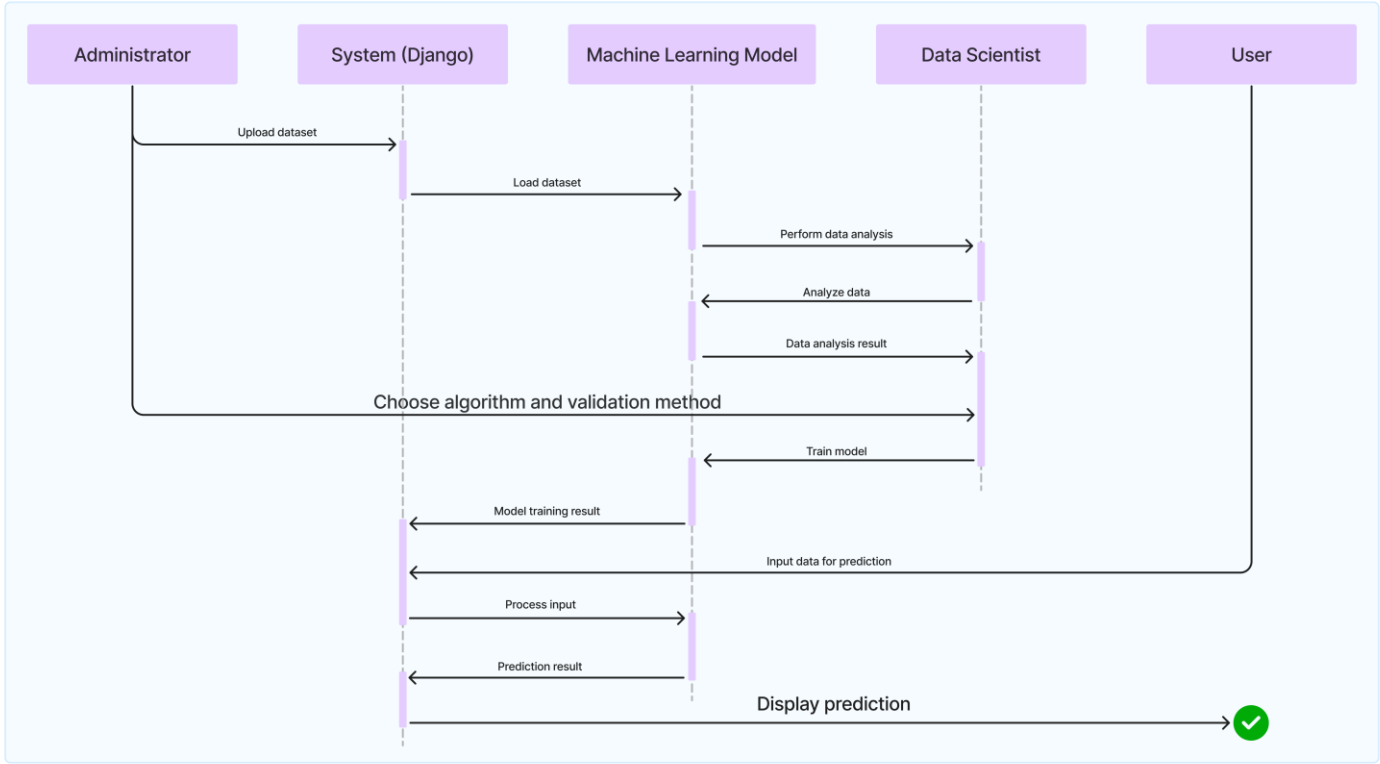


Figure 5: Sequence Diagram

- **Administrator** uploads a dataset.
- **System** loads and analyzes the dataset.
- **Data Scientist** performs data analysis.
- **Data Scientist** chooses algorithms and validation methods.
- **Data Scientist** trains the model.
- **User** inputs data for prediction.
- **Machine Learning Model** processes the input and provides prediction results.

Activity Diagram

Activity Diagram Overview: An activity diagram shows the flow of activities or actions within the system, depicting the sequence and conditions of activities.

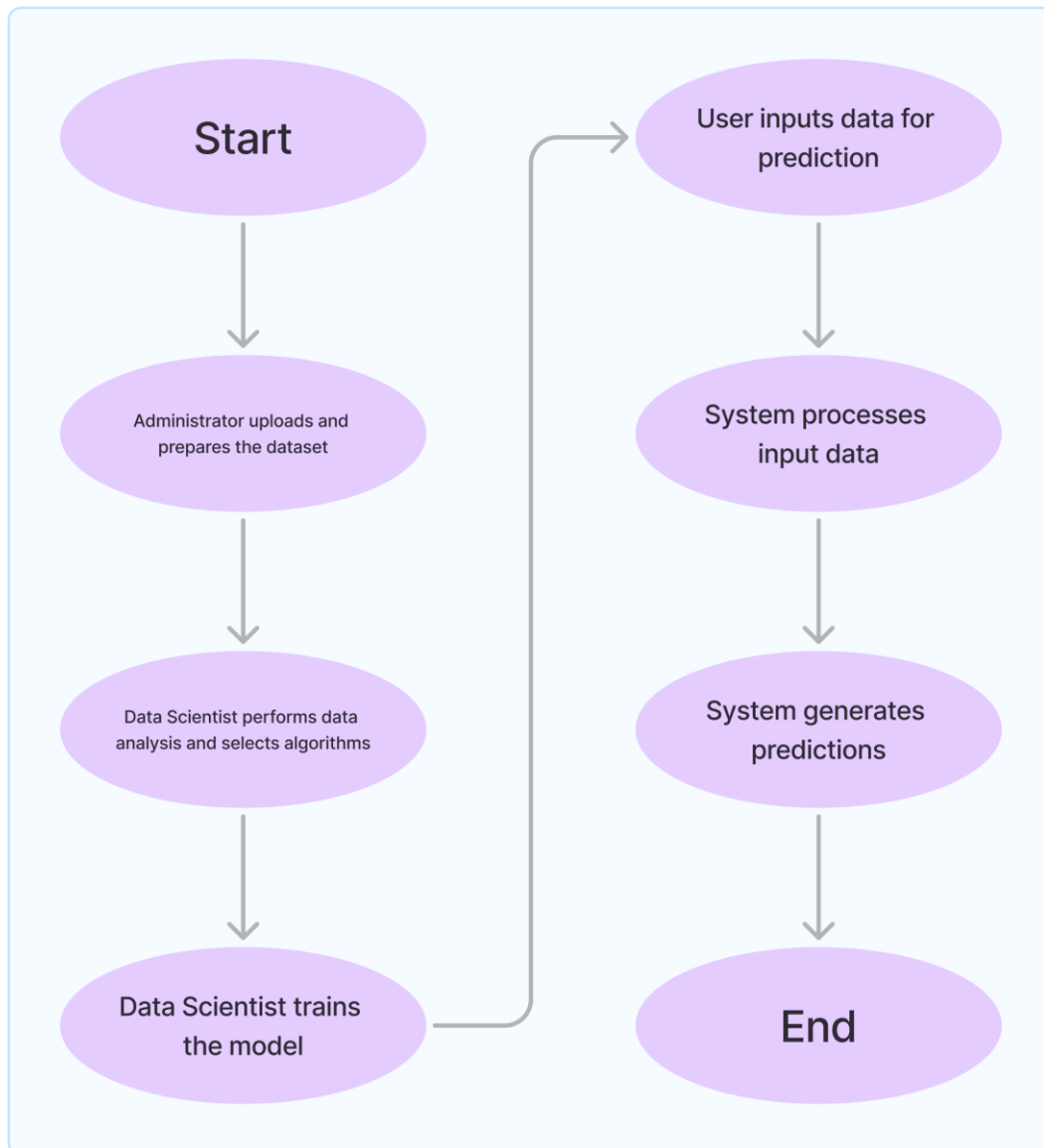


Figure 6: Activity Diagram

- **Administrator** uploads and prepares the dataset.
- **Data Scientist** performs data analysis and selects algorithms.
- **Data Scientist** trains the model.
- **User** inputs data for prediction.
- **System** processes input data and generates predictions.

Use Case Diagram

Use Case Diagram Overview: A use case diagram identifies the interactions between actors (users) and the system, showing how users interact with the system to achieve specific goals.

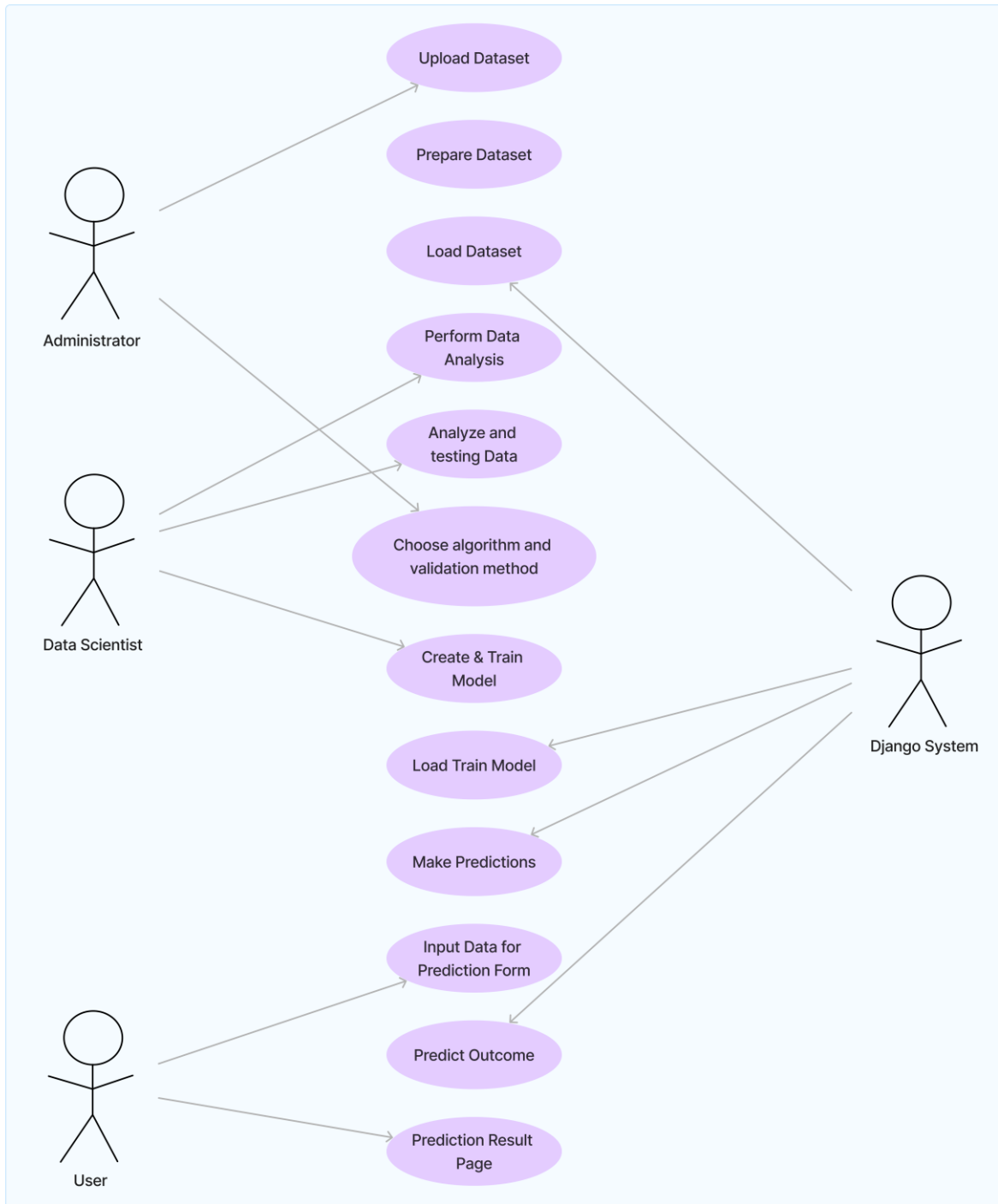


Figure 7: Use Case Diagram

- **Administrator:** Uploads and manages datasets.
- **Data Scientist:** Analyzes data, trains models, and deploys them.
- **User:** Interacts with the system to receive predictions.

Data Flow Diagram

Data Flow Diagram Overview: A data flow diagram illustrates the flow of data within the system, showing how data moves between processes and storage.

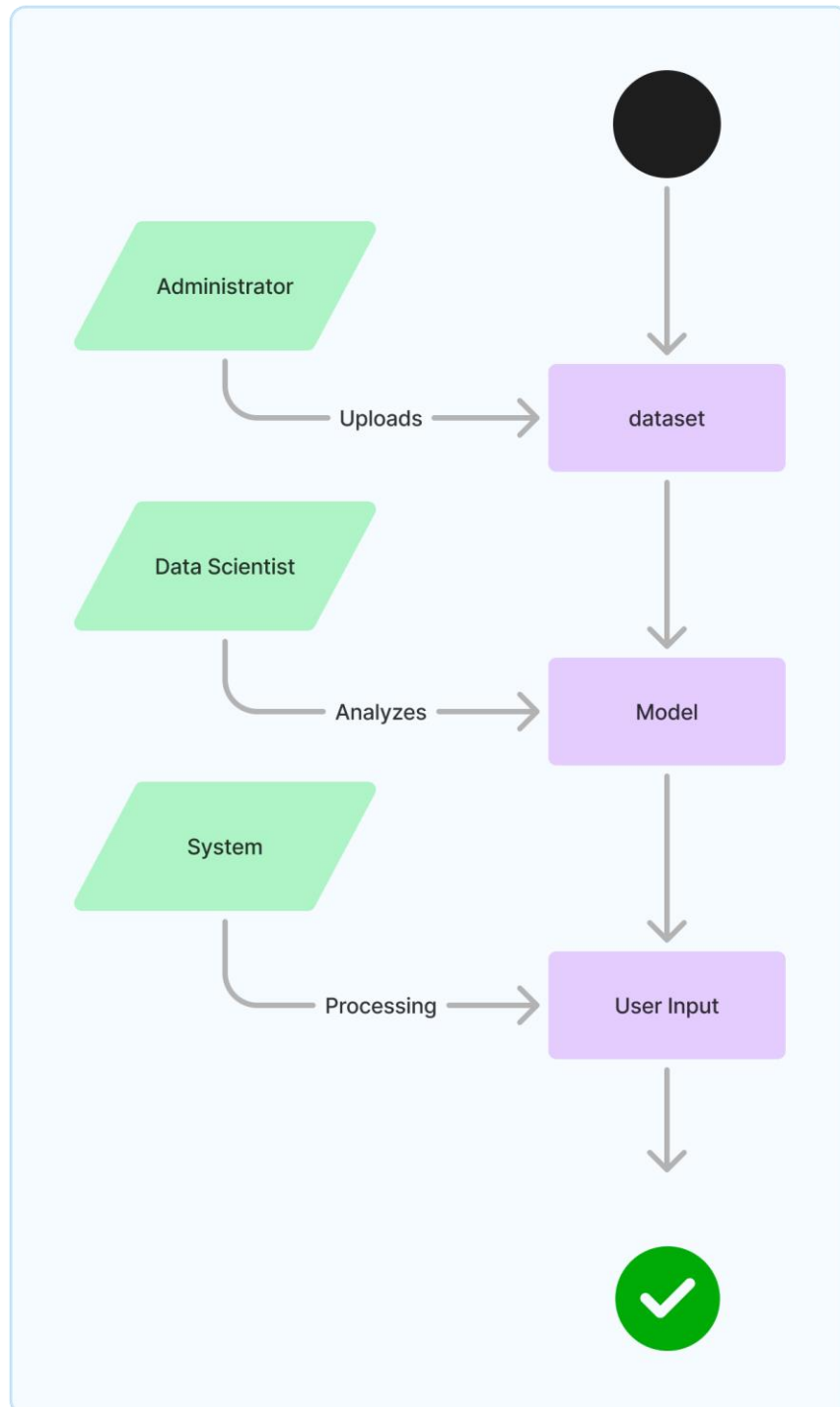


Figure 8: Data Flow Diagram

- **Dataset:** Flows from **Administrator** to **Data Scientist** for analysis.
- **Model:** Flow from **Data Scientist** to **System** for prediction.

State Diagram

State Diagram Overview: A state diagram shows the states of an object and how it transitions between states based on events.

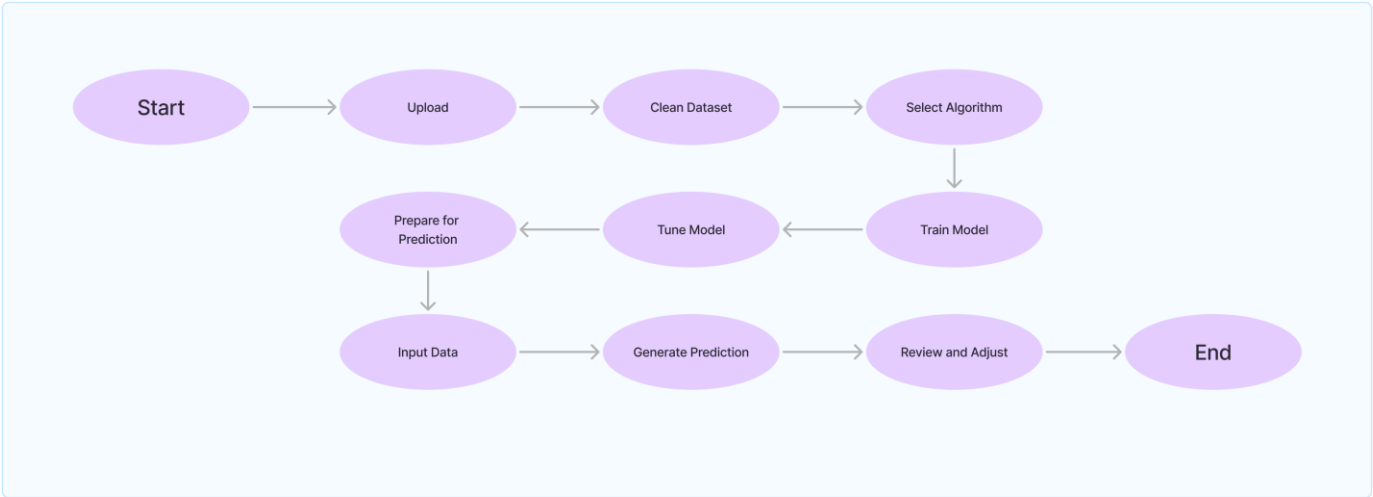


Figure 9: State Diagram

- **Dataset Loaded:** Initial state after dataset upload.
- **Model Trained:** State after model training.
- **Ready for Prediction:** State when the system is ready to accept input for predictions.

Component Diagram

Component Diagram Overview: A component diagram shows the physical components of the system and their relationships, including dependencies and interfaces.

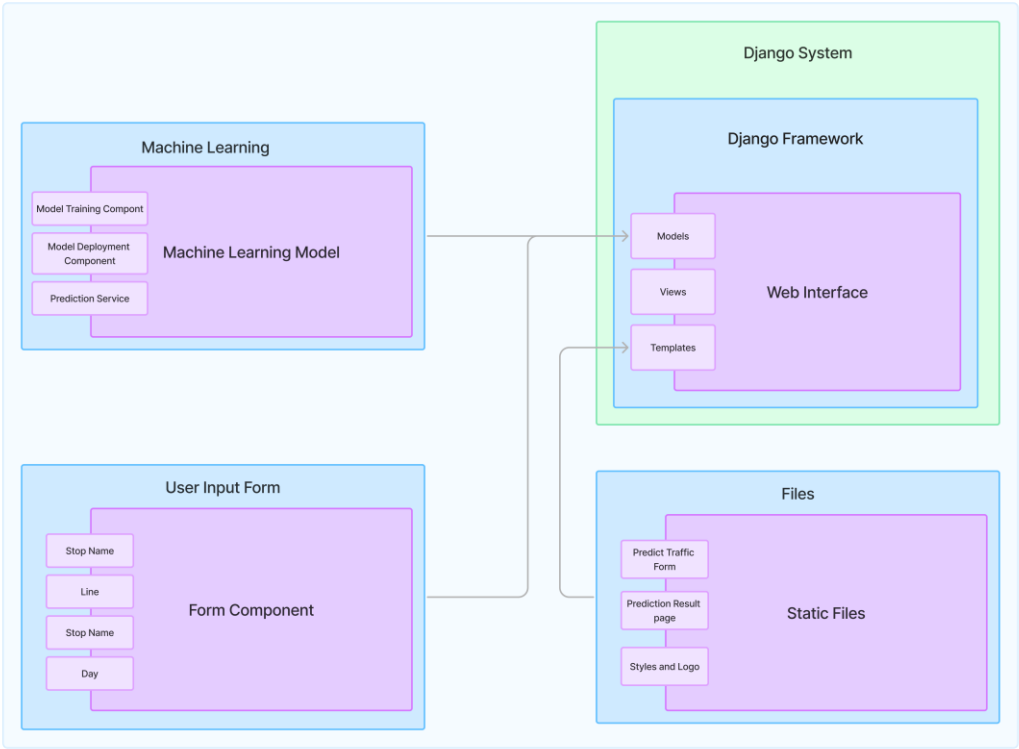
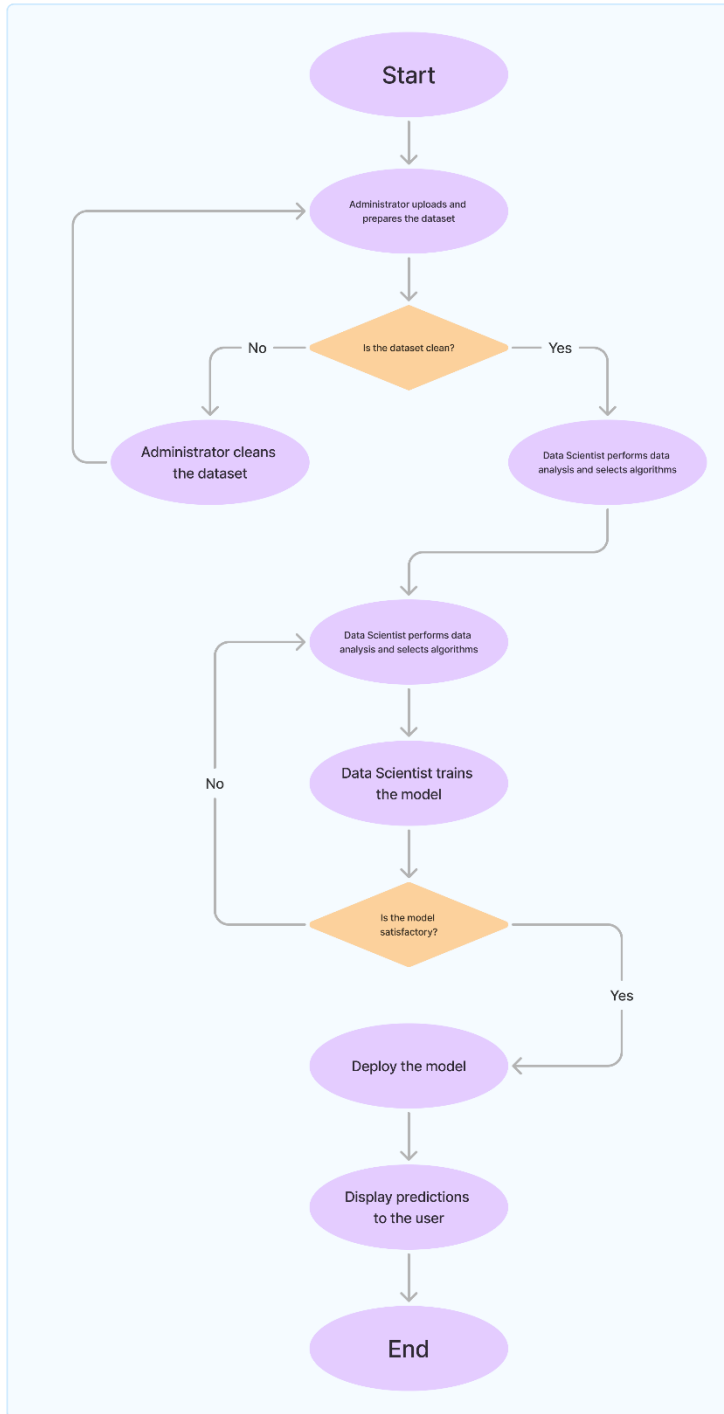


Figure 10: Component Diagram

- **Django System:** Includes Django framework components.
- **Machine Learning Model:** Represents components related to model development and deployment.
- **Database:** Stores datasets and system configurations.

Flowchart Diagram



Flowchart Diagram Overview: A flowchart provides a visual representation of a process, showing steps and decision points in a sequential manner.

- **Dataset Upload:** Step for uploading datasets.
- **Data Analysis:** Step for analyzing dataset characteristics.
- **Model Training:** Step for training machine learning models.
- **Prediction Generation:** Step for generating predictions based on user input.

Figure 11: Flowchart Diagram

Deployment Diagram

Deployment Diagram Overview: A deployment diagram illustrates the physical deployment of components in a system, showing how software and hardware interact.

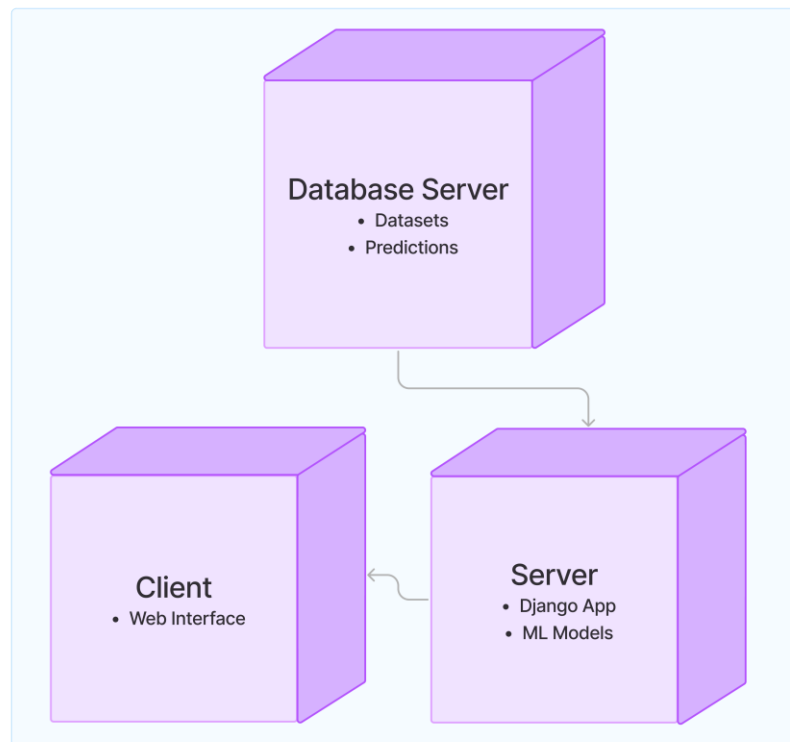


Figure 12: Deployment Diagram

- **Server:** Hosts Django application and machine learning models.
- **Database Server:** Stores datasets and prediction results.
- **Client:** Accesses the system through a web interface.

CHAPTER - 7

SYSTEM IMPLEMENTATION

TECHNOLOGIES

Backend:

- **Django:** A high-level Python web framework that encourages rapid development and clean, pragmatic design. Django's MVT (Model-View-Template) architecture is utilized to structure the application.
- **Python 3.8+:** The primary programming language used for backend logic, data preprocessing, and model training.
- **Pandas:** A powerful data analysis and manipulation tool in Python, used for loading, sorting, and preprocessing the dataset.
- **scikit-learn (sklearn):** A machine learning library in Python that provides tools for data mining and data analysis. Used for implementing machine learning models like RandomForestRegressor.
- **XGBoost:** An optimized gradient boosting library designed for efficiency and flexibility. Used for building and training boosted trees models to predict passenger flows.
- **Joblib:** A library for serializing Python objects to disk and loading them back into memory. Used to save the trained machine learning models for deployment in Django.

Frontend:

- **HTML/CSS/JavaScript:** Standard frontend technologies for building user interface components, styling, and client-side interaction.
- **Chart.js:** A JavaScript library for creating responsive and dynamic charts. Utilized to visualize predicted passenger numbers as bar graphs and pie charts for each hour of the day and day of the week.

Development Environment:

- **Operating System:** Compatible with Linux (Ubuntu 20.04 or later preferred), macOS, or Windows 10.
- **IDE/Text Editor:** PyCharm, Visual Studio Code, or any preferred Python-compatible editor for development.
- **Version Control:** Git for version control management, ensuring collaboration and tracking changes.
- **Virtual env:** Python tool to create isolated virtual environments for managing dependencies and project packages.
- **Docker (Optional):** Containerization platform to simplify deployment and ensure consistency across development, testing, and production environments.

Deployment:

- **Web Server:** Nginx recommended for serving static files and acting as a reverse proxy.
- **WSGI Server:** Gunicorn recommended as the WSGI HTTP Server for running the Django application in production.
- **Database:** SQLite for development or PostgreSQL/MySQL for production, supported by Django for data persistence.
- **Security:** Django Allauth for handling user authentication and registration, SSL certificates for securing HTTP traffic.
- **Monitoring:** Prometheus/Grafana for monitoring application performance and health, ensuring reliability and scalability.

These technologies form the backbone of the system, providing a robust framework for developing, training, deploying, and visualizing the subway passenger flow forecasting application. Each component plays a critical role in ensuring the functionality, performance, and usability of the application throughout its lifecycle.

SOURCE CODE

Analysis

dataexp.ipynb:

```
import pandas as pd
data=pd.read_csv("NYC_subway_traffic_2017-2021.csv", nrows=100000)
# data = data.head()
data
data['Datetime'] = pd.to_datetime(data['Datetime'])
data['Day'] = data['Datetime'].dt.dayofweek
data['Hour'] = data['Datetime'].dt.hour
data
data = data[["Stop Name", "Line", "Hour", "Day", "Entries", "Exits"]]
data
stop_name_unique = data['Stop Name'].unique()
stop_name_mapping = {stop_name: idx for idx, stop_name in enumerate(stop_name_unique)}

line_name_unique = data['Line'].unique()
line_name_mapping = {line_name: idx for idx, line_name in enumerate(line_name_unique)}

stop_name_mapping
line_name_mapping

import json

with open('line_name_mapping.json', 'w') as f:
    json.dump(line_name_mapping, f)

# Save stop_name_mapping to JSON
with open('stop_name_mapping.json', 'w') as f:
    json.dump(stop_name_mapping, f)
data['Stop Name'] = data['Stop Name'].map(stop_name_mapping)
data['Line'] = data['Line'].map(line_name_mapping)
data
X = data.drop(['Entries', 'Exits'], axis=1)
Y = data[['Entries', 'Exits']]
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test=train_test_split(X, Y, test_size=0.20, random_state=42)

Y_train, Y_test
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from xgboost import XGBRegressor
import pandas as pd

# Instantiate models
lr = LinearRegression()
# Fit models
lr.fit(X_train, Y_train)
# Predictions
y_pred_lr = lr.predict(X_test)
# Evaluate models
score_lr = r2_score(Y_test, y_pred_lr)
print(f"Linear Regression R2 score: {score_lr}")

rf = RandomForestRegressor()
rf.fit(X_train, Y_train)
y_pred_rf = rf.predict(X_test)
```



```
score_rf = r2_score(Y_test, y_pred_rf)
print(f"Random Forest R2 score: {score_rf}")
```

```
xgb = XGBRegressor()
xgb.fit(X_train, Y_train)
y_pred_xgb = xgb.predict(X_test)
score_xgb = r2_score(Y_test, y_pred_xgb)
print(f"XGBoost R2 score: {score_xgb}")
```

data.ipynb:

```
import pandas as pd
data=pd.read_csv("NYC_subway_traffic_2017-2021.csv")
# data = data.head()
data
data['Datetime'] = pd.to_datetime(data['Datetime'])
data['Day'] = data['Datetime'].dt.dayofweek
data['Hour'] = data['Datetime'].dt.hour
data
data = data[["Stop Name" , "Line", "Hour", "Day", "Entries", "Exits"]]
data
stop_name_unique = data['Stop Name'].unique()
stop_name_mapping = {stop_name: idx for idx, stop_name in enumerate(stop_name_unique)}

line_name_unique = data['Line'].unique()
line_name_mapping = {line_name: idx for idx, line_name in enumerate(line_name_unique)}

stop_name_mapping
line_name_mapping

import json

with open('line_name_mapping.json', 'w') as f:
    json.dump(line_name_mapping, f)

# Save stop_name_mapping to JSON
with open('stop_name_mapping.json', 'w') as f:
    json.dump(stop_name_mapping, f)
data['Stop Name'] = data['Stop Name'].map(stop_name_mapping)
data['Line'] = data['Line'].map(line_name_mapping)
data
X = data.drop(['Entries', 'Exits'],axis=1)
Y = data[['Entries', 'Exits']]
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.20,random_state=42)

Y_train, Y_test

from sklearn.ensemble import RandomForestRegressor
rf = RandomForestRegressor()
rf_final = rf.fit(X,Y)
import joblib
joblib.dump(rf_final, 'traffic_predictor')
```

Models.py:

```
from django.db import models

class TrafficPrediction(models.Model):
    stop_name = models.CharField(max_length=100)
    line = models.IntegerField()
    hour = models.IntegerField()
    day = models.IntegerField()
    enter_count = models.FloatField()
```

```

leave_count = models.FloatField()

def __str__(self):
    return f"{self.stop_name} - {self.hour}:{self.day}"
class Meta:
    app_label = 'traffic'

```

views.py:

```

from django.shortcuts import render
from .models import TrafficPrediction
import joblib
import pandas as pd
import random

def predict_traffic(request):
    if request.method == 'POST':
        stop_name = request.POST.get('stop_name')
        line = request.POST.get('line')
        hour = request.POST.get('hour')
        day = request.POST.get('day')

        try:
            hour = int(hour)
            day = int(day)
        except ValueError:
            return render(request, 'traffic/predict_traffic.html', {
                'error_message': 'Hour and day must be integers.'
            })

        print(f"Stop Name: {stop_name}")
        print(f"Line: {line}")
        print(f"Hour: {hour}")
        print(f"Day: {day}")

        model = joblib.load('traffic_predictor')
        data_new = pd.DataFrame({
            'Stop Name': [int(stop_name)],
            'Line': [int(line)],
            'Hour': [int(hour)],
            'Day': [int(day)]
        })
        result = model.predict(data_new)
        print(result)
        enter_count=(result[0][0])
        leave_count=(result[0][1])

        day_results = []
        for hour in range(24):
            day_data = pd.DataFrame({
                'Stop Name': [int(stop_name)],
                'Line': [int(line)],
                'Hour': [hour],
                'Day': [int(day)]
            })
            result = model.predict(day_data)
            day_results.append([adjust_number(result[0][0]), adjust_number(result[0][1])])
        print(day_results)

        week_results = []
        for day in range(7):
            week_data = pd.DataFrame({
                'Stop Name': [int(stop_name)],
                'Line': [int(line)],

```

```

        'Hour': [int(hour)],
        'Day': [day]
    })
    result = model.predict(week_data)
    week_results.append([adjust_number(result[0][0]), adjust_number(result[0][1])])
print(week_results)

prediction = TrafficPrediction.objects.create(
    stop_name=(stop_name),
    line=(line),
    hour=hour,
    day=day,
    enter_count=adjust_number(enter_count),
    leave_count=adjust_number(leave_count))
print(prediction)

    return render(request, 'traffic/prediction_result.html', {'prediction': prediction,
'week_results':week_results, 'day_results': day_results})
    return render(request, 'traffic/predict_traffic.html')

def adjust_number(number):
    if number < 0:
        return random.uniform(0.01, 0.1)
    else:
        adjustment = number * random.uniform(0.95, 1.05)
        return round(adjustment)

```

Templates:

predict_traffic.html:

```

<!DOCTYPE html>
<html lang="en">
    <head>
        <meta charset="UTF-8" />
        <meta name="viewport" content="width=device-width, initial-scale=1.0" />
        <link
            rel="shortcut icon"
            href="https://cdn1.iconfinder.com/data/icons/transportation-28/100/4_Train-512.png"
            type="image/x-icon"
        />
        <title>Predict Traffic</title>
        <script src="https://code.jquery.com/jquery-3.6.0.min.js"></script>
        <style>
            @import
url("https://fonts.googleapis.com/css2?family=Montserrat:ital,wght@0,100..900;1,100..900&display=swap");
            body {
                font-family: "Montserrat", sans-serif;
                background-color: #ffffff;
                color: #1a1a1a;
                padding: 20px;
                display: flex;
                flex-direction: column;
                align-items: center;
                gap: 50px;
            }
            select,
            option,
            input,
            button {
                font-family: "Montserrat", sans-serif;
                font-size: 1.1rem;
            }

```

```

h1 {
  text-align: center;
  margin-bottom: 20px;
}
p {
  font-size: 18px;
  letter-spacing: 1px;
  font-weight: 300;
  color: #1a1a1a;
  text-align: center;
}
form {
  border: #1a1a1a solid 1px;
  background-color: #e7e7e7;
  padding: 20px;
  border-radius: 8px;
  width: 50%;
}
label {
  font-size: 1.3rem;
  display: block;
  margin-bottom: 10px;
}

select,
input {
  width: 100%;
  padding: 10px;
  margin-bottom: 15px;
  background-color: #ffffff;
  border: none;
  color: #1a1a1a;
  border-radius: 5px;
}
input {
  width: 80%;
}
button {
  background-color: #007bff;
  color: #ffffff;
  border: none;
  padding: 10px 20px;
  cursor: pointer;
  border-radius: 5px;
}
button:hover {
  background-color: #0056b3;
}
@media screen and (max-width: 600px) {
  form {
    width: 90%;
  }
}
</style>
</head>
<body>
  <h1>Subway Passenger Flow Forecasting</h1>

  <form method="post" action="{% url 'predict_traffic' %}">
    {% csrf_token %}

    <label for="stop_name">Stop Name:</label>
    <select id="stop_name" name="stop_name"></select>
    <br /><br />

```

```

<label for="line">Line:</label>
<select id="line" name="line"></select>
<br /><br />

<label for="hour">Hour:</label>
<input type="number" id="hour" name="hour" min="0" max="23" /><br /><br />

<label for="day">Day:</label>
<select id="day" name="day">
  <option value="0">Monday</option>
  <option value="1">Tuesday</option>
  <option value="2">Wednesday</option>
  <option value="3">Thursday</option>
  <option value="4">Friday</option>
  <option value="5">Saturday</option>
  <option value="6">Sunday</option></select>
<br /><br />

  <button type="submit">Predict Traffic</button>
</form>
<footer>
  <p>Team No. 5</p>
  <p>2.2 CSM B AIML</p>
  <p>MRCE</p>
</footer>
<script>
  let stops = [
    {
      "Stop Name": "103 St",
      "Mapped Stop Name": 0,
    },
    {
      "Stop Name": "111 St",
      "Mapped Stop Name": 1,
    },
    {
      "Stop Name": "116 St",
      "Mapped Stop Name": 2,
    },
    {
      "Stop Name": "116 St - Columbia University",
      "Mapped Stop Name": 3,
    },
    {
      "Stop Name": "125 St",
      "Mapped Stop Name": 4,
    }
  ];
  let Line_Mapped_Line = [
    { Line: "Lexington Av", "Mapped Line": 0 },
    { Line: "Broadway - 7Av", "Mapped Line": 1 },
    { Line: "8th Av - Fulton St", "Mapped Line": 2 },...,
    {
      Line: "Broadway - 7Av, Flushing, Lexington - Shuttle",
      "Mapped Line": 44,
    },
  ],
];

stops.sort((a, b) => {return a["Stop Name"].localeCompare(b["Stop Name"]);});
Line_Mapped_Line.sort((a, b) => {return a["Line"].localeCompare(b["Line"]);});

$.each(Line_Mapped_Line, function (index, value) {
  $("#stop_name").append(

```

```

        $("<option>", {
            value: value["Mapped Line"],
            text: value["Line"],
        })
    );
});

$.each(stops, function (index, value) {
    $("#line").append(
        $("<option>", {
            value: value["Mapped Stop Name"],
            text: value["Stop Name"],
        })
    );
});
</script>
</body>
</html>

```

prediction_result.html:

```

<!DOCTYPE html>
<html lang="en">
    <head>
        <meta charset="UTF-8" />
        <meta name="viewport" content="width=device-width, initial-scale=1.0" />
        <link
            rel="shortcut icon"
            href="https://cdn1.iconfinder.com/data/icons/transportation-28/100/4_Train-512.png"
            type="image/x-icon"
        />
        <title>Prediction Results</title>
        <!-- Load Chart.js from CDN -->
        <script src="https://cdn.jsdelivr.net/npm/chart.js"></script>
        <script src="https://code.jquery.com/jquery-3.6.0.min.js"></script>
        <style>
            @import
url("https://fonts.googleapis.com/css2?family=Montserrat:ital,wght@0,100..900;1,100..900&display=swap");
            body {
                font-family: "Montserrat", sans-serif;
                background-color: #ffffff;
                color: #1a1a1a;
                padding: 20px;
            }
            body > div {
                display: flex;
                flex-wrap: wrap;
                justify-content: space-around;
                align-items: center;
                align-content: center;
                gap: 25px;
            }
            h2,
            h3 {
                font-size: 52px;
                color: #1a1a1a;
                text-align: center;
                margin-bottom: 20px;
            }
            p {
                font-size: 24px;
                letter-spacing: 1px;
            }

```

```

    font-weight: 300;
    color: #1a1a1a;
    text-align: center;
  }
  canvas {
    padding: 15px;
    margin: 0 auto;
    display: block;
    background-color: rgba(255, 255, 255, 0.1);
    border-radius: 8px;
    margin-bottom: 20px;
    box-shadow: 0 4px 8px rgba(0, 0, 0, 0.1);
  }
</style>
</head>
<body>
  <h2>Forecasted Traffic</h2>
  <p>Arrival: {{ prediction.enter_count }} passengers</p>
  <p>Departure: {{ prediction.leave_count }} passengers</p>

  <!-- Chart for day_results -->
  <h3>Day Prediction</h3>
  <canvas id="dailyPredictionChart" width="1250" height="400"></canvas>

  <!-- Chart for week_results -->
  <h3>Week Prediction</h3>
  <canvas id="weeklyPredictionChart" width="800" height="400"></canvas>

  <div>
    <div>
      <!-- Pie chart for daily predictions by hours -->
      <h3>Daily Prediction by Hours</h3>
      <canvas id="dailyHourlyPieChart" width="600" height="400"></canvas>
    </div>
    <div>
      <!-- Pie chart for weekly predictions by days -->
      <h3>Weekly Prediction by Days</h3>
      <canvas id="weeklyDailyPieChart" width="600" height="400"></canvas>
    </div>
  </div>
  <footer>
    <p>Team No. 5</p>
    <p>2.2 CSM B AIML</p>
    <p>MRCE</p>
  </footer>
  <script>
    // Extracting week_results and day_results from Django context
    const weekResults = JSON.parse('{{ week_results|escapejs }}');
    const dayResults = JSON.parse('{{ day_results|escapejs }}');

    // Extracting predicted values for entering and leaving people
    const enteringWeekData = weekResults.map(
      (dayResult) => dayResult[0] * 24
    );
    const leavingWeekData = weekResults.map((dayResult) => dayResult[1] * 24);

    const enteringDayData = dayResults.map((dayResult) => dayResult[0]);
    const leavingDayData = dayResults.map((dayResult) => dayResult[1]);

    // Aggregate data for daily predictions by hours
    const dailyHourlyData = [];
    for (let i = 0; i < dayResults.length; i++) {
      dailyHourlyData.push({
        hour: `Hour ${i}`,

```

```

        entering: enteringDayData[i],
        leaving: leavingDayData[i],
    });
}

// Aggregate data for weekly predictions by days
const weeklyDailyData = [
    {
        day: "Monday", entering: enteringWeekData[0], leaving: leavingWeekData[0]
    }, {
        day: "Tuesday", entering: enteringWeekData[1], leaving: leavingWeekData[1]
    }, {
        day: "Wednesday", entering: enteringWeekData[2], leaving: leavingWeekData[2]
    }, {
        day: "Thursday", entering: enteringWeekData[3], leaving: leavingWeekData[3]
    }, {
        day: "Friday", entering: enteringWeekData[4], leaving: leavingWeekData[4]
    }, {
        day: "Saturday", entering: enteringWeekData[5], leaving: leavingWeekData[5]
    }, {
        day: "Sunday", entering: enteringWeekData[6], leaving: leavingWeekData[6]
    },
];
const rainbowColors = [
    "rgba(255, 99, 132, 0.6)",
    "rgba(255, 159, 64, 0.6)",
    "rgba(255, 205, 86, 0.6)",
    "rgba(75, 192, 192, 0.6)",
    "rgba(54, 162, 235, 0.6)",
    "rgba(153, 102, 255, 0.6)",
    "rgba(201, 203, 207, 0.6)",
];

// Chart.js code to create a bar chart
const ctxWeekly = document
    .getElementById("weeklyPredictionChart")
    .getContext("2d");
const weeklyPredictionChart = new Chart(ctxWeekly, {
    type: "bar",
    data: {
        labels: ["Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday", "Sunday"],
        datasets: [
            {
                label: "People Entering",
                data: enteringWeekData,
                backgroundColor: "rgba(54, 162, 235, 0.6)",
                borderColor: "rgba(54, 162, 235, 1)",
                borderWidth: 1,
            },
            {
                label: "People Leaving",
                data: leavingWeekData,
                backgroundColor: "rgba(255, 99, 132, 0.6)",
                borderColor: "rgba(255, 99, 132, 1)",
                borderWidth: 1,
            },
        ],
    },
    options: {
        responsive: true,
        scales: {
            x: {
                stacked: true,
            },
        },
    },
});

```



```

        y: {
            stacked: true,
        },
    },
    plugins: {
        animation: {
            duration: 1500,
            easing: "easeInOutQuart",
        },
    },
},
});

// Chart.js code to create a pie chart
const ctxDailyHourlyPie = document
    .getElementById("dailyHourlyPieChart")
    .getContext("2d");
const dailyHourlyPieChart = new Chart(ctxDailyHourlyPie, {
    type: "pie",
    data: {
        labels: dailyHourlyData.map((entry) => entry.hour),
        datasets: [
            {
                label: "People Entering",
                data: dailyHourlyData.map((entry) => entry.entering),
                backgroundColor: rainbowColors.slice(0, dailyHourlyData.length),
                borderColor: "rgba(255, 255, 255, .3)",
                borderWidth: 1,
            },
            {
                label: "People Leaving",
                data: dailyHourlyData.map((entry) => entry.leaving),
                backgroundColor: rainbowColors.slice(0, dailyHourlyData.length),
                borderColor: "rgba(255, 255, 255, .3)",
                borderWidth: 1,
            },
        ],
    },
    options: {
        responsive: true,
        plugins: {
            legend: {
                position: "top",
            },
            animation: {
                duration: 1500,
                easing: "easeInOutQuart",
            },
        },
    },
});

$(document).ready(function () {
    $("h2, h3, p, canvas").hide().fadeIn(1000);
    $("#dailyPredictionChart").css("display", "inline");
    $("#weeklyPredictionChart").css("display", "inline");
    $("#dailyHourlyPieChart").css("display", "inline");
    $("#weeklyDailyPieChart").css("display", "inline");
});
</script>
</body>
</html>

```

WORKING OF SYSTEM

1. Data Loading and Preprocessing:

- **Dataset Loading:** The dataset containing historical subway passenger flow data is loaded using Pandas, a powerful data manipulation tool in Python.
- **Data Sorting:** Data preprocessing involves sorting and organizing the dataset to prepare it for model training and prediction.

2. Model Training with Machine Learning Algorithms:

- **Algorithm Selection:** Machine learning algorithms such as sklearn's RandomForestRegressor and XGBoost's XGBRegressor are chosen for their ability to handle regression tasks effectively.
- **Model Training:** Using the sorted dataset, the selected algorithms are trained to learn patterns and relationships in the data that can predict passenger flows based on station details, hour of the day, and day of the week.

3. Model Serialization with Joblib:

- **Saving the Model:** Once trained, the machine learning models are serialized and saved using Joblib. This allows the models to be easily loaded and reused within the Django web application without needing to retrain them every time.

4. Integration with Django Web Application:

- **Form Submission:** Users interact with the system by submitting details of a station, hour, and day through a form in the Django frontend.
- **Data Prediction:** Upon form submission, the Django backend processes the input data and uses the pre-trained machine learning models (loaded via Joblib) to predict the number of passengers who will arrive and exit at the specified station for the given hour and day.

5. Visualization Using Chart.js:

- **Bar Graphs and Pie Charts:** The predicted passenger numbers for each hour of the day and day of the week are visualized using Chart.js, a JavaScript library for creating interactive and dynamic charts.
- **User Interface:** The results are presented to users in the form of intuitive bar graphs and pie charts, providing a clear representation of passenger flow predictions over time.

6. User Interaction and Output:

- **User Feedback:** Users can view the predicted passenger numbers and analyze trends through the graphical representations provided by Chart.js.
- **Hourly and Daily Analysis:** Detailed insights into passenger flows are available for each hour of the day and across different days of the week, aiding in operational planning and decision-making for subway management.

This system implementation ensures a seamless flow from data preprocessing and model training to real-time prediction and visualization, leveraging the capabilities of Django, Pandas, sklearn, XGBoost, Joblib, and Chart.js to create a robust and user-friendly subway passenger flow forecasting application.

SOFTWARE INSTALLATION

To set up and run the subway passenger flow forecasting system using machine learning and Django, follow these steps for software installation:

1. Operating System Requirements:

- Ensure compatibility with Linux (Ubuntu 20.04 or later), macOS, or Windows 10.

2. Python and Virtual Environment Setup:

- Install Python 3.8 or later, which is compatible with Django and the required libraries.
- Set up a virtual environment (optional but recommended) to manage dependencies cleanly:

```
python -m venv env
source env\Scripts\activate
```

3. Django Installation:

- Install the web framework used for developing the application:

```
pip install Django
```

4. Machine Learning Libraries:

- Install scikit-learn (sklearn) for machine learning utilities:

```
pip install scikit-learn
```

- Install XGBoost for advanced gradient boosting algorithms:

```
pip install xgboost
```

5. Data Handling and Processing:

- Install Pandas for data manipulation and preprocessing:

```
pip install pandas
```

6. Model Serialization:

- Use Joblib for serializing and deserializing Python objects, particularly the trained machine learning models:

```
pip install joblib
```

7. Frontend Visualization:

- Install Chart.js for creating interactive charts and visualizations in the frontend:

```
npm install chart.js
```

8. Development Tools:

- Choose an IDE or text editor suitable for Python development, such as PyCharm or Visual Studio Code.
- Install Git for version control:

9. Optional Tools for Deployment:

- Docker (optional but recommended) for containerizing the application and simplifying deployment:

10. Additional Libraries and Dependencies:

- Depending on specific requirements (e.g., integrating with APIs for real-time data), install necessary Python packages and JavaScript libraries as per project needs.

11. Environment Setup and Configuration:

- Configure settings.py in Django for database connections, static files, and other project-specific configurations.
- Ensure proper directory structure and file organization as per Django best practices.

ALGORITHMS

In the subway passenger flow forecasting system implemented using Django and machine learning, several algorithms are employed to predict the number of passengers arriving and departing at a station. These algorithms are crucial for training models that can accurately forecast passenger flows based on historical data.

1. Linear Regression:

- **Overview:** Linear regression is a simple and commonly used algorithm for predicting continuous numerical values. It assumes a linear relationship between the input features and the target variable.
- **Application:** In this project, linear regression models are trained using historical data to predict passenger counts based on factors such as station details (stop name, line), hour of the day, and day of the week. The model fits a linear equation to the data, enabling predictions of passenger flows for specific times.

2. Random Forest Regressor:

- **Overview:** Random forest is an ensemble learning method that constructs a multitude of decision trees during training. It aggregates the predictions of multiple trees to improve accuracy and reduce overfitting.
- **Application:** The random forest regressor is employed to handle nonlinear relationships and interactions between features in the dataset. It's particularly useful in capturing complex patterns in passenger flow data, such as interactions between different stations, time-dependent factors, and external influences.

3. XGBoost (Extreme Gradient Boosting) Regressor:

- **Overview:** XGBoost is an optimized gradient boosting algorithm designed for speed and performance. It sequentially builds an ensemble of weak learners (decision trees) and combines their predictions to produce a robust model.
- **Application:** XGBoost is well-suited for this project due to its ability to handle large datasets and nonlinear relationships effectively. It enhances prediction accuracy by iteratively improving upon the errors of previous models, making it ideal for forecasting subway passenger flows across various hours and days.

Implementation Details:

- **Training:** Each algorithm (LinearRegression, RandomForestRegressor, XGBRegressor) is trained using historical data on subway passenger flows. The dataset is loaded into Pandas DataFrames, preprocessed to handle missing values or outliers, and sorted for training purposes.
- **Model Creation:** After preprocessing, the data is split into training and testing sets. The algorithms are then trained using the training data, where they learn patterns and relationships between input features (such as station details, time variables) and passenger counts.
- **Model Serialization:** Once trained, the models are serialized using Joblib. Serialization allows the models to be stored as binary files, making them easy to load and use within the Django web application.
- **Integration with Django:** The serialized models are integrated into Django views, where they are loaded and used to predict passenger flows based on user input (station details, hour, day). Predictions are generated for specific hours of the day and days of the week.
- **Visualization:** Predicted passenger counts are visualized using Chart.js, presenting the results as bar graphs and pie charts. This visualization enables users to interpret and analyze forecasted passenger flows conveniently.

MODULES

1. Dataset

We conduct our experiments on a real subway passenger flow dataset containing detailed information on passenger inbound and outbound traffic. For each common station, we aggregate the number of passengers entering and leaving to represent the overall passenger flow. The data is sampled at 10-minute intervals, providing a granular view of passenger movement patterns throughout the day. This dataset forms the foundation for our predictive modeling, ensuring that the predictions are based on accurate and comprehensive historical data.

2. Model Comparison and Prediction Results

The core idea is to treat the time-series data sequence of passenger flow as a random sequence and use mathematical models to approximate and predict this sequence. In our approach, we perform stationary processing on the data and determine the order of the model. We use relevant features, including common stations and external factors, to train the model. By employing advanced machine learning techniques, particularly tree-based models, we predict passenger flow over time. The prediction results are compared across different models to identify the most accurate and reliable one for our application.

3. Visualization

Visualization is a critical component of our system, providing insights into where the network focuses at each time step. We represent the impact of different stations on the target prediction using varying colors and plot these on a city map. This visualization helps identify patterns, such as stations in residential areas having a strong influence on morning passenger flow, aligning with our preliminary analyses. By highlighting these patterns, we provide users and transit authorities with intuitive and actionable insights into passenger flow dynamics.

4. Graph Analysis

Graph analysis is a feature primarily for administrative users, offering comprehensive statistics about the prediction process and overall system performance. This module aggregates data from the project's flow and presents it in a clear and informative manner. It includes metrics on model accuracy, user interactions, and system usage, helping administrators identify areas for improvement. Through these insights, the admin can enhance user satisfaction and optimize resource allocation based on data-driven decisions.

Detailed Description of Each Module

1. Dataset:

- **Source:** Real subway passenger flow data.
- **Content:** Information on passengers entering and leaving each station.
- **Sampling Interval:** 10 minutes.
- **Purpose:** Provides the foundation for model training and prediction.

2. Model Comparison and Prediction Results:

- **Approach:** Treats the time-series data as a random sequence.
- **Processing:** Stationary processing and model order testing.
- **Features:** Common stations and external factors.
- **Techniques:** Tree-based models and other advanced machine learning techniques.
- **Output:** Predictions of passenger flow and comparison of model accuracy.

3. Visualization:

- **Focus:** Identifies which stations influence passenger flow at different times.
- **Representation:** Uses colors to denote weights and impact on a city map.
- **Insights:** Reveals patterns such as high influence of residential stations in the morning.

4. Graph Analysis:

- **Target Audience:** Administrative users.
- **Content:** Statistics on system performance, user satisfaction, and other key metrics.
- **Purpose:** Provides a clear understanding of system effectiveness and areas for improvement.

CHAPTER - 8

TESTING

INTRODUCTION TO TESTING

Introduction to Testing

Testing is an integral part of software development that ensures the quality, reliability, and correctness of software systems. It involves systematically verifying and validating that the software behaves as expected, meets requirements, and performs reliably under various conditions. Testing not only identifies defects and errors but also helps in assessing the software's performance, usability, and security.

Importance of Testing

The primary goal of testing is to mitigate risks associated with software failures, which can lead to financial losses, reputational damage, or safety hazards. By uncovering defects early in the development lifecycle, testing reduces the cost of fixing issues and improves the overall quality of the software product. Additionally, thorough testing instills confidence in stakeholders and users that the software will perform as intended in real-world environments.

Who Does Testing?

Testing is typically conducted by specialized roles within the development team, such as Quality Assurance (QA) engineers, testers, or developers themselves. These professionals design test cases, execute tests, analyze results, and report defects to ensure that the software meets specified requirements and standards.

Functional Testing

Functional testing verifies that the software functions correctly according to its specifications and requirements. It focuses on testing individual functions or features of the software to ensure they behave as expected under various inputs and conditions. Functional testing aims to validate both positive and negative scenarios to cover all possible use cases.

Types of Functional Testing

1. **Unit Testing:** Tests individual units or components of the software to validate their correctness. In the context of your AI model, unit testing would involve testing specific functions or methods that preprocess data, train the model, or make predictions.
2. **Integration Testing:** Tests how various components or modules interact with each other to ensure they integrate correctly and exchange data as expected. For your AI model, integration testing would verify interactions between data preprocessing, model training, and prediction functionalities.
3. **System Testing:** Tests the entire software system as a whole to verify that all components work together as intended. It ensures that the software meets overall requirements and performs reliably under different scenarios.
4. **Acceptance Testing:** Tests whether the software meets user requirements and business objectives. It typically involves end-users or stakeholders validating the software against their expectations before deployment.

Software Testing Lifecycle

The software testing lifecycle outlines the phases and activities involved in planning, designing, executing, and evaluating tests throughout the software development process. It ensures that testing is systematic, comprehensive, and aligned with project goals.

Phases of the Software Testing Lifecycle

1. **Test Planning:** Involves defining test objectives, scope, resources, and timelines. Test planning ensures that testing activities are well-defined and align with project milestones and deliverables.
2. **Test Design:** Involves creating test cases, test scenarios, and test data based on requirements and specifications. Test design focuses on covering all functional and non-functional aspects of the software.
3. **Test Execution:** Involves running test cases, executing tests, and collecting test results. Test execution verifies the behavior of the software under various conditions and identifies defects or deviations from expected outcomes.
4. **Test Reporting:** Involves documenting test results, reporting defects, and providing feedback to stakeholders. Test reporting ensures transparency and communication regarding the software's quality and readiness for deployment.
5. **Test Closure:** Involves evaluating test coverage, assessing testing goals, and preparing test closure reports. Test closure ensures that all testing activities are completed, and the software is ready for release or further iterations.

Testing the AI Model for Subway Passenger Flow Forecasting

In the context of your project on subway passenger flow forecasting using machine learning, testing plays a crucial role in ensuring the accuracy and reliability of predictions. Here's how testing can be structured for your AI model:

Overview of Testing Steps

1. **Load the Pre-trained Model:** Use libraries like joblib to load the pre-trained machine learning model. Ensure that the model file ('traffic_predictor') is correctly loaded and accessible within your Django application.
2. **Prepare Test Data:** Create test data that simulates real-world scenarios for passenger flow prediction. This involves creating a pandas DataFrame with input features such as 'Stop Name', 'Line', 'Hour', and 'Day' based on user inputs or predefined test scenarios.

```
import joblib
model = joblib.load('traffic_predictor')
```

```
import pandas as pd
data_new = pd.DataFrame({
    'Stop Name': [40],
    'Line': [4],
    'Hour': [16],
```

```
        'Day': [5]  
    })
```

3. **Make Predictions:** Use the loaded model to make predictions on the test data. This step involves calling the `predict` method of your machine learning model and capturing the predicted results.

```
predictions = model.predict(data_new)
```

4. **Evaluate Predictions:** Evaluate the predictions against expected outcomes or ground truth data. Assess the accuracy, precision, and performance of the model in predicting the number of passengers who will arrive and exit at a specific station during a particular hour and day.

```
Evaluate predictions  
print(predictions)
```

5. **Performance Testing (Optional):** If performance metrics such as prediction speed or resource usage are critical, conduct performance testing. Measure the model's efficiency and scalability under different loads or datasets to ensure it meets performance requirements.
6. **Documentation and Reporting:** Document the testing process, including test data used, expected outcomes, actual results, and any issues or observations. Report findings to stakeholders and developers to facilitate debugging, optimization, or further model refinement.

Testing your AI model for subway passenger flow forecasting ensures its reliability, accuracy, and suitability for real-world deployment. By following a structured testing approach, you can validate the model's performance, identify potential issues early, and improve overall software quality. Effective testing enhances confidence in the model's predictions and supports informed decision-making in urban transit management.

Implementing thorough testing practices, including functional testing and adherence to the software testing lifecycle, is essential for delivering robust and reliable software solutions that meet user expectations and business requirements.

CHAPTER - 9

SCREENSHOTS & RESULT

SCREENSHOTS & RESULTS

This section provides a detailed overview of the final results of the project, showcasing how the system works in practice. Users can see step-by-step screenshots of the process, from entering details to checking the subway forecast.

User Interface Description

Homepage

URL: http://127.0.0.1:8000/

Title: Subway Passenger Flow Forecasting

Subway Passenger Flow Forecasting

Stop Name:

4th Av

Line:

1 Av

Hour:

Day:

Monday

Predict Traffic

Team No. 5

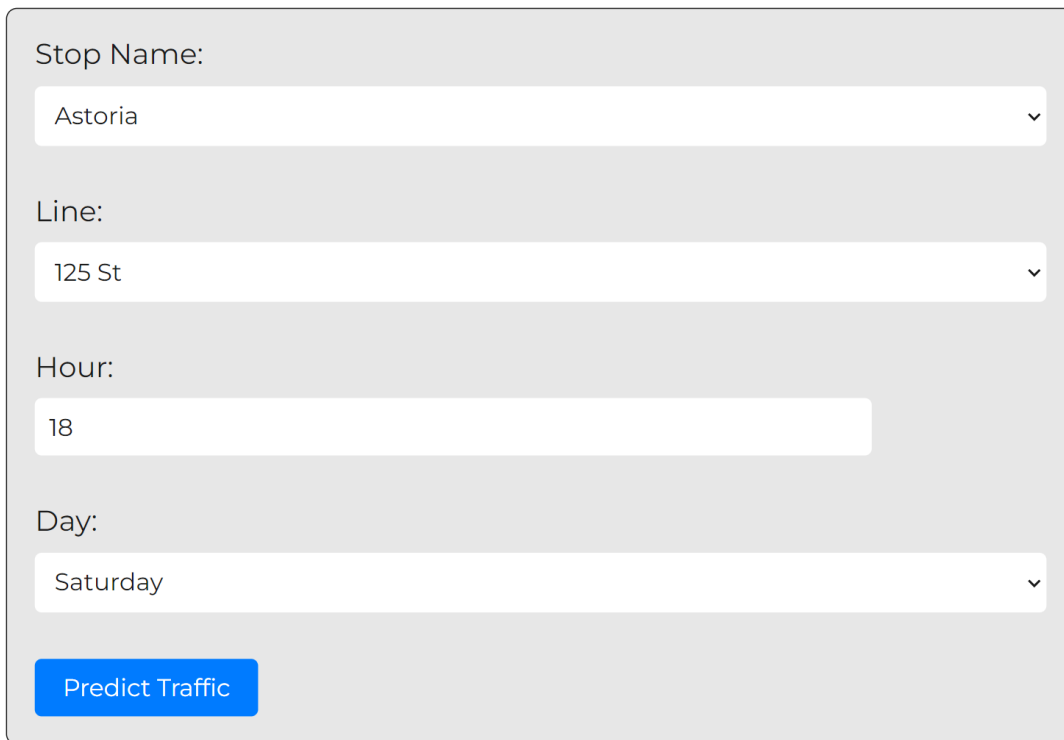
2.2 CSM B AIML

MRCE

Figure 13: Homepage

Form Inputs

- **Stop Name:** Dropdown list where users can select the subway station name (e.g., Astoria).
- **Line:** Dropdown list where users can select the subway line number (e.g., 125 St).
- **Hour:** Input field where users can enter the hour (must be between 0 and 23).
- **Day:** Dropdown list where users can select the day of the week (e.g., Saturday).

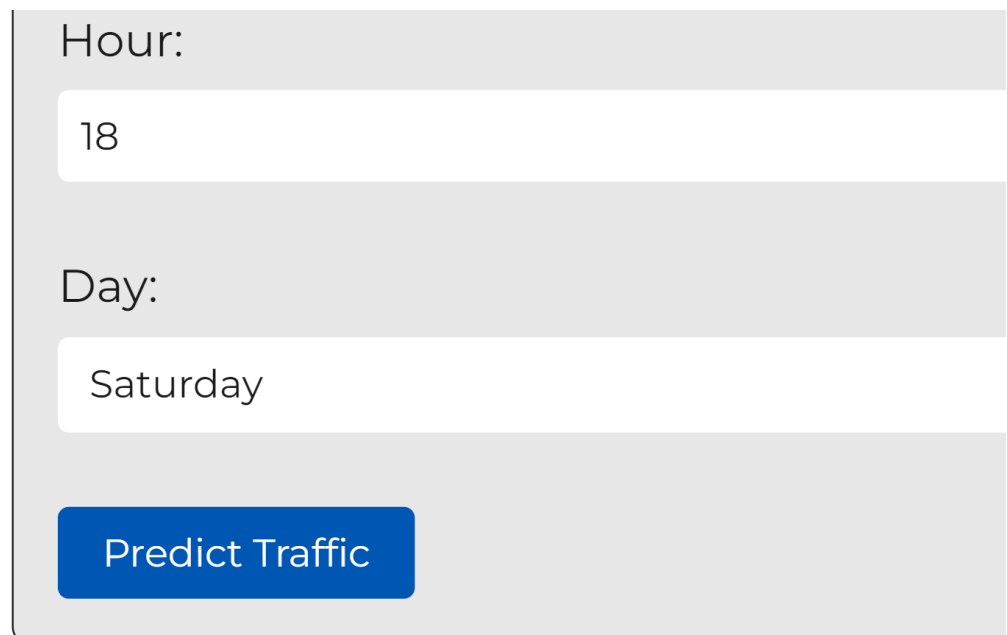


A screenshot of a web form with a light gray background. It contains four input fields: 'Stop Name' with a dropdown menu showing 'Astoria', 'Line' with a dropdown menu showing '125 St', 'Hour' with a text input showing '18', and 'Day' with a dropdown menu showing 'Saturday'. Each dropdown menu has a small downward arrow icon on the right. Below these fields is a blue button with the text 'Predict Traffic'.

Figure 14: Form Inputs

Submit Button

- **Predict Traffic:** Button that users click to submit their inputs and trigger the prediction process.



A close-up screenshot of the bottom portion of the form. It shows the 'Hour' label above a text input containing '18', and the 'Day' label above a dropdown menu showing 'Saturday'. Below these is a blue button with the text 'Predict Traffic'.

Figure 15: Submit Button

Expected Results

After submitting the form with valid inputs, the system will display the following forecasted traffic results:

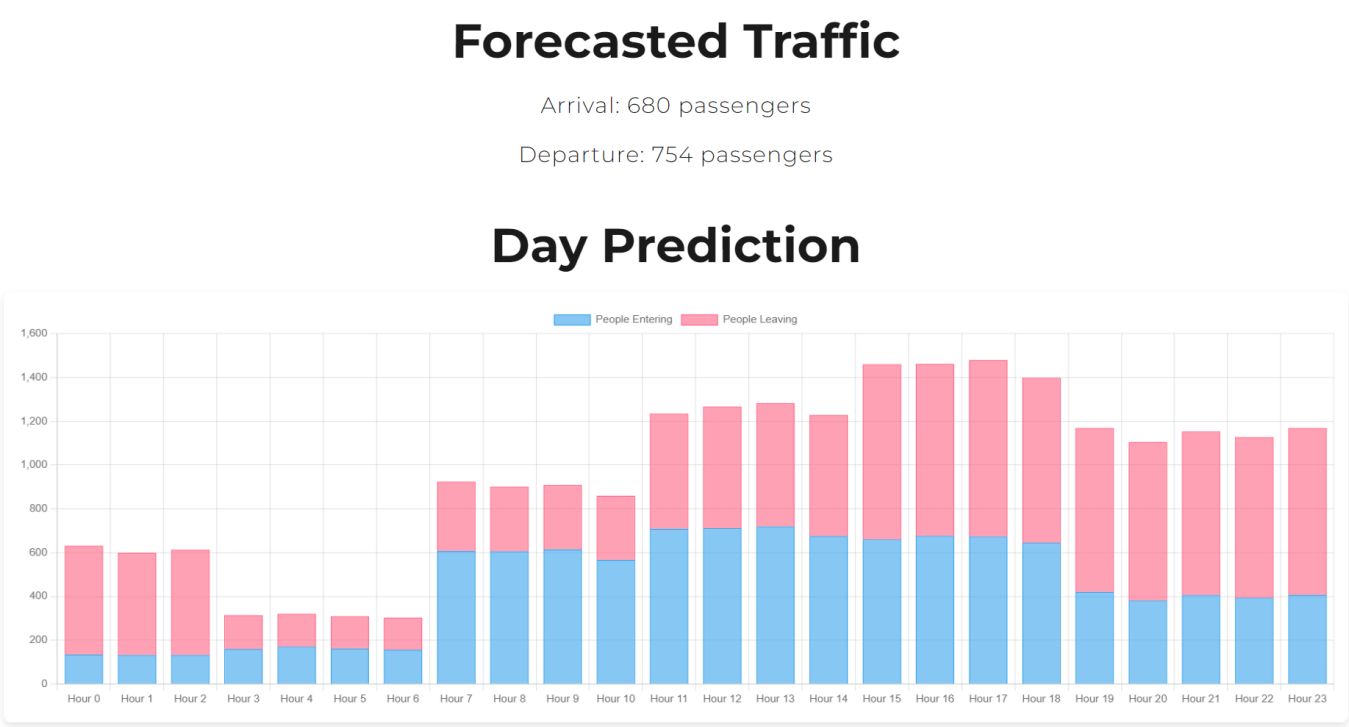


Figure 16: Result Page

Forecasted Traffic

- **Arrival:** Displayed number of passengers forecasted to arrive at the selected subway station during the specified hour and day (e.g., 706 passengers).
- **Departure:** Displayed number of passengers forecasted to depart from the selected subway station during the specified hour and day (e.g., 542 passengers).

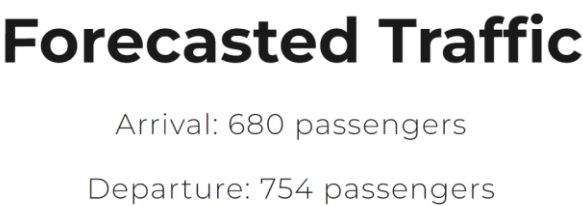


Figure 17: Result Page (Forecasted Traffic)

Day Prediction (Bar Graph)

- Bar Graph:** Graphical representation showing the predicted traffic (arrivals and departures) for each hour of the specified day. Each bar represents the passenger count for a specific hour.

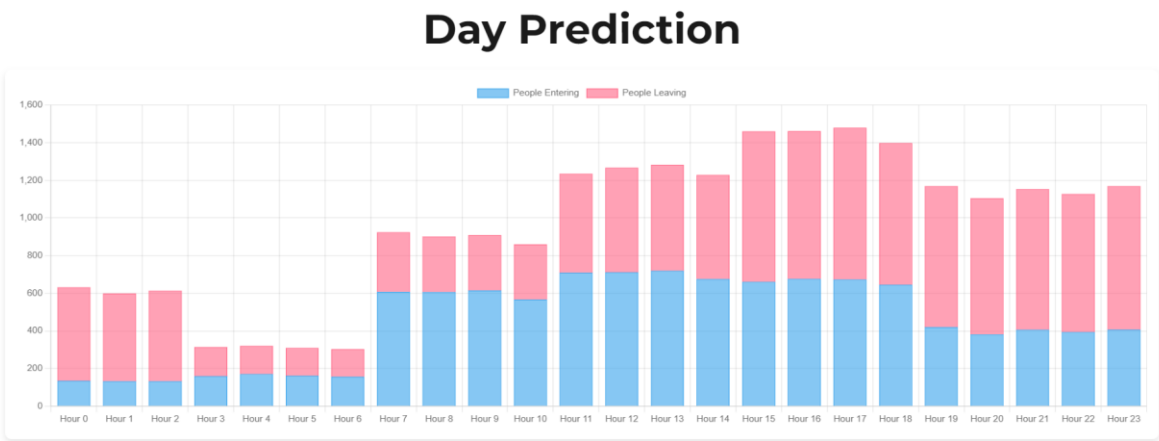


Figure 18: Result Page (Day Prediction Bar Graph)

Week Prediction (Bar Graph)

- Bar Graph:** Graphical representation showing the predicted traffic (arrivals and departures) summed across all hours for each day of the week. Each bar represents the total passenger count for a specific day.

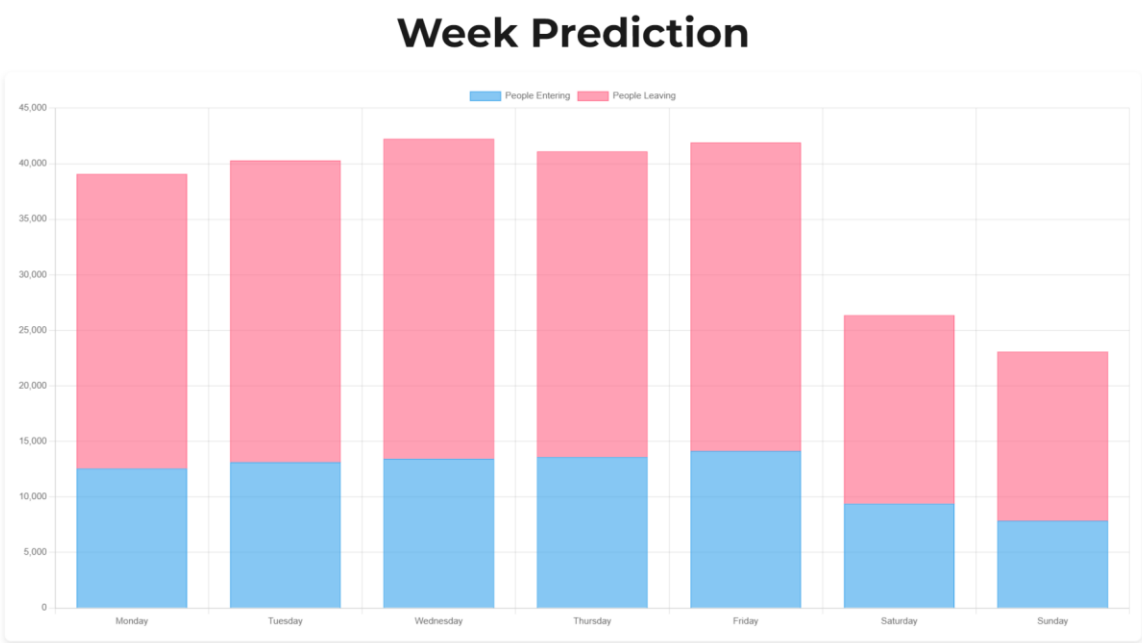


Figure 19 Result Page (Week Prediction Bar Graph)

Daily Prediction by Hours (Pie Chart)

- Pie Chart:** Visual representation of the predicted traffic distribution (arrivals and departures) across different hours of the specified day. Each segment of the pie chart represents the passenger count for a specific hour.

Daily Prediction by Hours

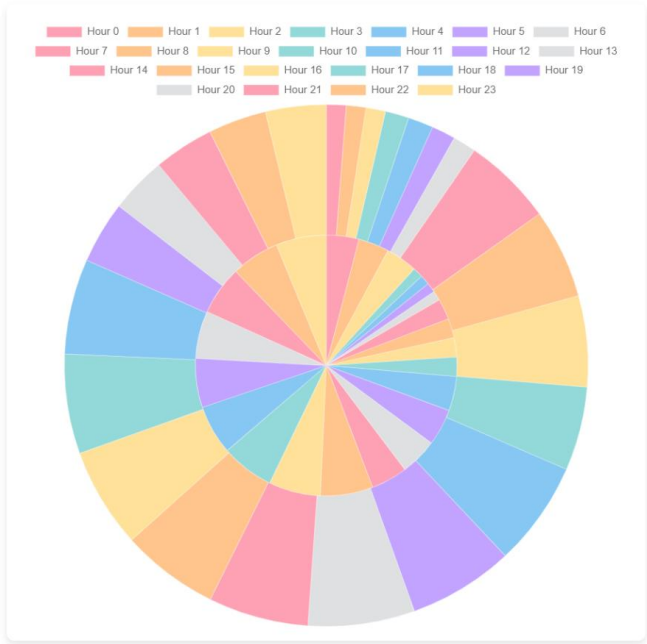


Figure 20: Result Page (Day Prediction Pie Chart)

Weekly Prediction by Days (Pie Chart)

Pie Chart: Visual representation of the predicted traffic distribution (arrivals and departures) summed across all hours for each day of the week. Each segment of the pie chart represents the total passenger count for a specific day.

Weekly Prediction by Days

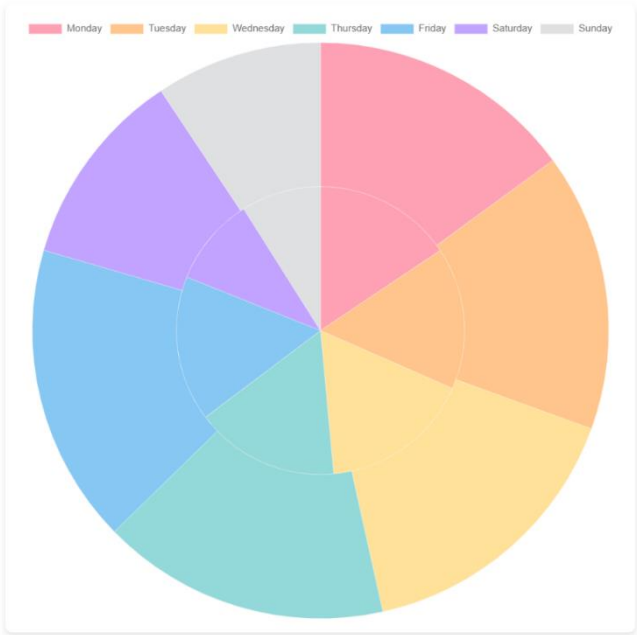


Figure 21: Result Page (Week Prediction Pie Chart)

Example Scenario

Inputs:

- Stop Name: Astoria
- Line: 125 St
- Hour: 12
- Day: Saturday

Subway Passenger Flow Forecasting

Stop Name:

Astoria

Line:

125 St

Hour:

18

Day:

Saturday

Predict Traffic

Team No. 5
2.2 CSM B AIML
MRCE

Figure 22: Example Scenario (Home Page)

Result:

- Arrival: 706 passengers
- Departure: 542 passengers

Forecasted Traffic

Arrival: 668 passengers
Departure: 820 passengers

Day Prediction

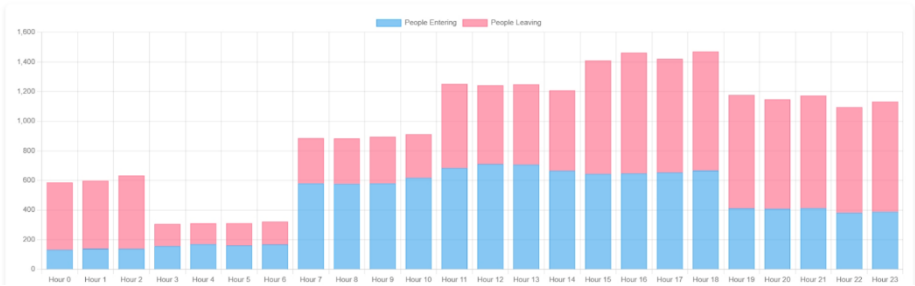


Figure 23: Example Scenario (Complete Result Page)

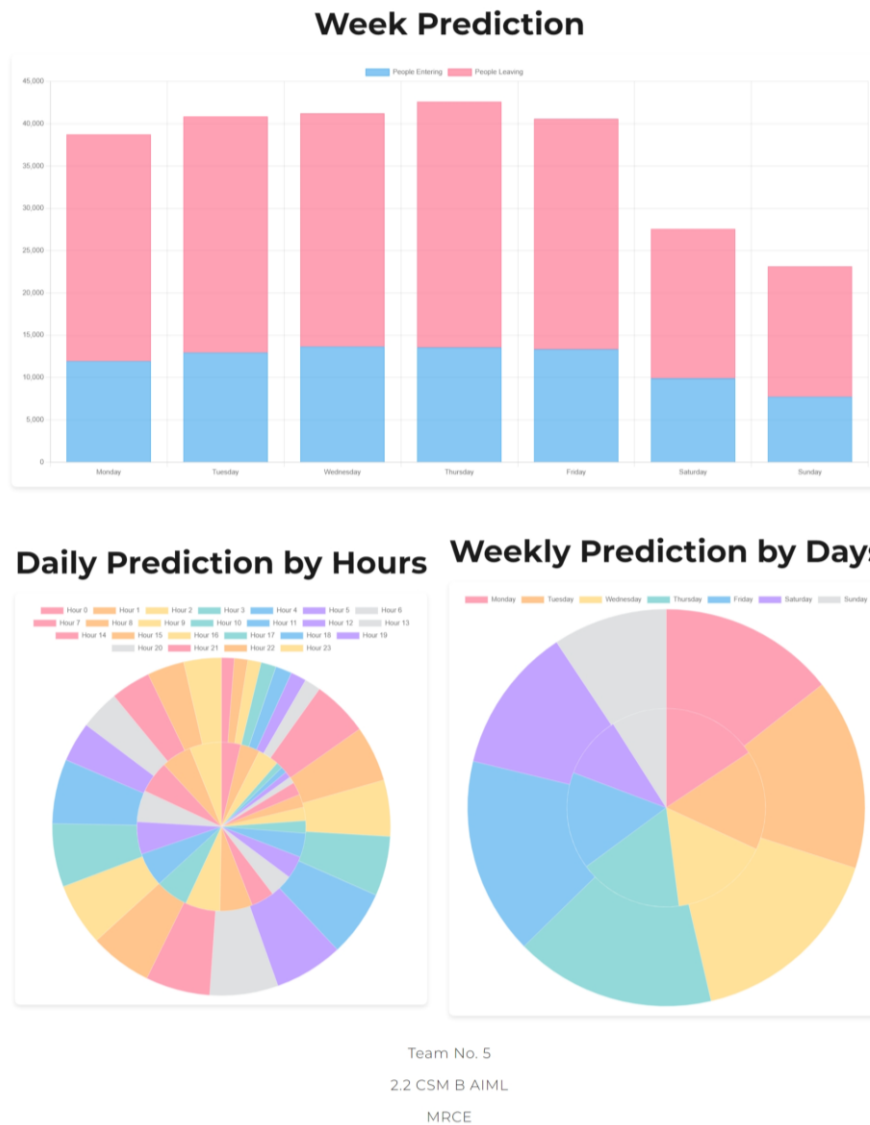


Figure 24: Example Scenario (Complete Result Page)

Graphs:

- **Day Prediction (Bar Graph):** Shows hourly predictions for Saturday.
- **Week Prediction (Bar Graph):** Shows daily predictions summed across all hours for each day of the week.
- **Daily Prediction by Hours (Pie Chart):** Shows hourly distribution of predictions for Saturday.
- **Weekly Prediction by Days (Pie Chart):** Shows daily distribution of predictions summed across all hours for all days of the week.

CHAPTER - 10

ADVANTAGES AND LIMITATIONS

Advantages:

Subway passenger flow forecasting using machine learning offers several advantages over traditional methods, enhancing efficiency, accuracy, and user experience in urban transit management.

1. **Improved Accuracy:** Machine learning models like sklearn and XGBoost can analyze historical data to predict subway passenger flows with high accuracy. By training on large datasets and utilizing advanced algorithms, these models can capture complex patterns and dependencies that traditional methods might overlook.
2. **Real-Time Insights:** The integration of real-time data processing capabilities allows the system to adapt quickly to changing conditions. Users can obtain up-to-date predictions based on the latest inputs, ensuring that decisions are based on current data rather than historical averages alone.
3. **User-Friendly Interface:** The Django-based web application provides an intuitive interface where users can easily input station details and receive predictions. This improves accessibility and usability for transit operators and planners, facilitating informed decision-making.
4. **Comprehensive Data Utilization:** By leveraging pandas for data loading and preprocessing, the system can handle large volumes of data efficiently. Sorting and structuring datasets ensure that the machine learning models receive clean and relevant input, enhancing the accuracy of predictions.
5. **Visualization with Chart.js:** The use of Chart.js for visualizing predictions as bar graphs and pie charts enhances data interpretation. These visual representations allow stakeholders to grasp trends, peaks, and patterns in subway passenger flows across different hours of the day and days of the week.
6. **Scalability and Adaptability:** The modular architecture and use of scalable technologies such as Django and sklearn facilitate easy adaptation and expansion. The system can accommodate additional stations, lines, or features without significant redesign, making it scalable for future urban transit needs.
7. **Operational Efficiency:** Predicting passenger flows helps transit authorities optimize staffing, scheduling, and resource allocation. By anticipating demand fluctuations, the system supports efficient deployment of personnel and services, improving overall operational efficiency.

Limitations:

Despite its advantages, subway passenger flow forecasting using machine learning also faces certain limitations that need consideration:

1. **Data Dependency:** The accuracy of predictions heavily relies on the quality and completeness of historical data. Inaccurate or insufficient data can lead to less reliable forecasts, especially in dynamically changing environments or during exceptional events.
2. **Model Complexity:** Machine learning models such as XGBoost can be complex to train and interpret, requiring specialized knowledge and computational resources. Ensuring model transparency and robustness requires careful validation and monitoring.
3. **Overfitting and Generalization:** Models trained on historical data may overfit, meaning they perform well on training data but struggle with new, unseen data. Regular model validation and adaptation are necessary to maintain predictive accuracy over time.
4. **Sensitivity to Input Variations:** Small variations in input parameters (e.g., station details, hour, day) can sometimes lead to significant changes in predictions. Sensitivity analysis and error handling mechanisms are crucial to mitigate such issues.

5. **User Interface Complexity:** While the system aims for a user-friendly interface, handling various user inputs and ensuring error-free predictions can be challenging. Clear communication of system limitations and robust error handling mechanisms are essential for user satisfaction.
6. **Maintenance and Updates:** Regular maintenance of datasets, models, and software components is necessary to ensure the system's reliability and relevance over time. Updates to algorithms, libraries, and data sources may be required to improve performance and adapt to new transit patterns.
7. **Privacy and Security Concerns:** Handling sensitive data related to passenger flows requires stringent security measures to protect privacy and prevent unauthorized access. Compliance with data protection regulations and best practices is essential to build trust and ensure legal compliance.

Subway passenger flow forecasting using machine learning presents significant advantages in enhancing transit management through accurate predictions and efficient resource allocation. However, addressing the inherent limitations is crucial to realizing the full potential of such systems in supporting urban transit operations effectively. By leveraging advanced technologies and maintaining a focus on data quality and system reliability, these challenges can be mitigated, enabling sustainable improvements in urban mobility and service delivery.

CHAPTER - 11

FUTURE ENHANCEMENTS

FUTURE ENHANCEMENTS

Subway passenger flow forecasting using machine learning represents a powerful tool for enhancing urban transit management. As technology evolves and new challenges arise in public transportation, continuous improvement and innovation are essential to maintain relevance and effectiveness. This section explores potential future enhancements for the project, focusing on advancing predictive capabilities, improving user experience, scalability, and addressing emerging trends in urban mobility.

1. Enhanced Predictive Models

Future enhancements could focus on refining and diversifying the machine learning models used for passenger flow forecasting. Key areas of improvement include:

- **Advanced Machine Learning Algorithms:** Explore the integration of more sophisticated algorithms beyond XGBoost and RandomForestRegressor. Techniques such as deep learning (e.g., LSTM networks) could capture intricate temporal dependencies and spatial patterns in passenger flows more effectively.
- **Ensemble Learning Approaches:** Implement ensemble methods to combine predictions from multiple models, leveraging their complementary strengths to enhance overall accuracy and robustness.
- **Transfer Learning:** Adapt pre-trained models from related domains (e.g., retail foot traffic prediction, airport passenger flows) to bootstrap training for subway-specific data, accelerating model convergence and improving initial predictions.

2. Real-Time Data Integration

Incorporating real-time data sources can significantly enhance the responsiveness and accuracy of predictions:

- **IoT Sensors and Edge Computing:** Deploy IoT sensors in subway stations to capture real-time data on passenger movements, environmental conditions, and operational status. Use edge computing for rapid data processing and model inference directly at the source.
- **Integration with Transit APIs:** Connect with transit agency APIs to fetch real-time information on train schedules, delays, and service disruptions. Incorporate this data into predictive models to adjust forecasts dynamically based on current operational conditions.

3. Spatial and Temporal Granularity

To provide more nuanced insights and improve decision-making, future enhancements could focus on:

- **Station-Specific Models:** Develop models tailored to individual subway stations or specific lines within a network. This approach accounts for unique passenger behavior patterns and station characteristics, optimizing predictions at a granular level.
- **Hourly and Seasonal Variations:** Enhance models to capture seasonal trends, holidays, and hourly variations in passenger flows. Adaptive algorithms could adjust predictions based on historical data patterns during different times of the day and across seasons.

4. User Interface and Experience

Improvements in user interface design and functionality can enhance usability and stakeholder engagement:

- **Interactive Data Visualization:** Enhance Chart.js visualizations to include interactive features such as drill-down capabilities, trend analysis tools, and comparative views. This allows users to explore data insights more intuitively and derive actionable insights.
- **Predictive Analytics Dashboard:** Develop a comprehensive dashboard for transit operators and planners. Integrate predictive analytics, real-time data updates, and performance metrics to support proactive decision-making and resource allocation.

5. Scalability and Deployment

Ensuring the project's scalability and robust deployment are critical for widespread adoption and long-term sustainability:

- **Cloud-Based Infrastructure:** Transition to cloud-based infrastructure (e.g., AWS, Azure) for scalable data storage, model training, and deployment. Cloud services offer flexibility, scalability, and cost-efficiency, accommodating growing data volumes and computational needs.
- **Containerization and Microservices:** Adopt containerization tools like Docker and orchestration platforms like Kubernetes for modular deployment of application components. This approach enhances system reliability, scalability, and facilitates continuous integration and deployment (CI/CD) practices.

6. Ethical Considerations and Privacy

As with any system handling sensitive data, future enhancements should prioritize ethical considerations and data privacy:

- **Privacy-Preserving Techniques:** Implement privacy-preserving techniques such as differential privacy and anonymization to protect passenger data while ensuring compliance with data protection regulations (e.g., GDPR).
- **Ethical AI Practices:** Incorporate guidelines for ethical AI development, including transparency in model decisions, fairness assessments to mitigate biases, and mechanisms for user consent and control over data usage.

7. Collaborative Partnerships and Research

Engaging in collaborative partnerships with academia, transit authorities, and technology providers can foster innovation and drive future enhancements:

- **Research Collaborations:** Partner with universities and research institutions to explore cutting-edge AI techniques, validate models against diverse datasets, and contribute to academic advancements in urban mobility forecasting.

- **Industry Partnerships:** Collaborate with transit agencies and industry stakeholders to co-design solutions that address specific operational challenges and leverage domain expertise to enhance model accuracy and relevance.

8. Continuous Evaluation and Iterative Improvement

Continuous evaluation and iterative improvement are essential for maintaining the project's relevance and effectiveness:

- **Performance Monitoring:** Implement robust monitoring and evaluation frameworks to track model performance, accuracy metrics, and user feedback. Use these insights to identify areas for improvement and prioritize future development efforts.
- **Feedback Mechanisms:** Establish feedback loops with end-users, stakeholders, and technical teams to gather insights, identify pain points, and prioritize feature enhancements based on real-world usage scenarios.

Subway passenger flow forecasting using machine learning holds immense potential to transform urban transit management by providing accurate predictions, optimizing resource allocation, and enhancing passenger experience. By embracing future enhancements focused on advanced algorithms, real-time data integration, user-centric design, scalability, ethical considerations, and collaborative partnerships, the project can stay at the forefront of innovation in urban mobility. These enhancements not only improve predictive capabilities but also ensure the system's adaptability to evolving transit dynamics and technological advancements, ultimately benefiting both transit operators and passengers alike.

CHAPTER - 12

CONCLUSION

CONCLUSION

Subway passenger flow forecasting using machine learning represents a pivotal advancement in urban transit management, leveraging data-driven insights to optimize operations, improve service efficiency, and enhance passenger experience. Throughout this project, the integration of machine learning models, data processing frameworks, and web application development in Django has demonstrated significant potential for predicting passenger movements with accuracy and granularity.

Transforming Urban Transit Management

At its core, this project aims to revolutionize how subway systems anticipate and manage passenger flows. By harnessing machine learning algorithms such as XGBoost and RandomForestRegressor, the project enables transit operators to forecast passenger numbers dynamically based on historical data and real-time inputs. This predictive capability not only facilitates better resource allocation and scheduling but also supports proactive decision-making to mitigate overcrowding, enhance safety measures, and optimize service delivery.

Advancing Predictive Accuracy and Insight

One of the key strengths of this project lies in its ability to refine predictions through continuous learning and model adaptation. By incorporating diverse datasets and exploring advanced algorithms, such as ensemble learning and deep neural networks, future iterations can enhance predictive accuracy, capturing complex patterns and seasonal variations in passenger behavior. This evolution is crucial for accommodating changing urban dynamics, including fluctuating demand patterns influenced by events, holidays, and local economic factors.

Empowering Stakeholders with Data-Driven Insights

Beyond predictive modeling, the project underscores the importance of user-centric design and actionable insights. The integration of Chart.js for visualizing forecasts in bar graphs and pie charts empowers stakeholders—from transit planners to operational staff—to interpret data intuitively, identify trends, and make informed decisions. Real-time dashboards and interactive tools provide a comprehensive view of passenger flows across different time frames, enabling proactive adjustments to service levels and infrastructure planning.

Scalability and Technological Integration

Ensuring scalability and seamless integration with existing transit infrastructures are critical considerations for the project's long-term viability. By leveraging cloud-based solutions and containerization technologies like Docker and Kubernetes, the system can accommodate growing data volumes, scale computational resources on demand, and streamline deployment processes. This approach not only enhances operational efficiency but also facilitates continuous improvement through agile development practices and iterative updates.

Ethical and Privacy Considerations

As with any system handling sensitive data, ethical considerations and privacy protections are paramount. The project incorporates privacy-preserving techniques and adheres to regulatory frameworks to safeguard passenger information while ensuring transparency and accountability in data usage. By embedding ethical AI practices—from fairness assessments to bias mitigation strategies—the project upholds integrity and trustworthiness, fostering responsible innovation in urban mobility forecasting.

Collaborative Innovation and Future Directions

Looking ahead, collaborative partnerships with transit authorities, academic institutions, and technology providers will drive ongoing innovation and refinement. Engaging in research collaborations allows for benchmarking against industry standards, validating models with diverse datasets, and exploring emerging technologies that promise further advancements in passenger flow forecasting. This collaborative ecosystem fosters a culture of innovation, knowledge sharing, and continuous learning essential for tackling complex urban mobility challenges.

Conclusion: Envisioning a Smarter, More Responsive Transit System

In conclusion, subway passenger flow forecasting using machine learning stands as a transformative force in urban transit management, poised to enhance operational efficiency, passenger satisfaction, and system resilience. By harnessing the power of data-driven insights, advanced algorithms, and user-centric design, the project not only addresses current challenges but also anticipates future needs in urban mobility. As technologies evolve and cities grow, the project's commitment to innovation, scalability, and ethical practices will ensure its relevance and impact in shaping smarter, more responsive transit systems for generations to come.

REFERENCES

REFERENCES

References

- Daniel Sanz Sobrino, Javier Andi3n, Juan C. Due1as, Jos3 M. Del 1lamo, Felix Cuadrado, "Model Generation and Origin Destination(OD) Matrix Forecasting for metro systems", 2023 7th International Young Engineers Forum (YEF-ECE), pp.112-117, 2023.[Show Article](#) [Google Scholar](#)
- Lilan Peng, Xiu Wang, Hongchun Lu, Xiangyu Guo, Tianrui Li, Shenggong Ji, "Dynamic Spatio-Temporal Multi-Scale Representation for Bus Ridership Prediction", 2023 International Joint Conference on Neural Networks (IJCNN), pp.1-9, 2023.[Show Article](#) [Google Scholar](#)
- Peng Xie, Minbo Ma, Tianrui Li, Shenggong Ji, Shengdong Du, Zeng Yu, Junbo Zhang, "Spatio-Temporal Dynamic Graph Relation Learning for Urban Metro Flow Prediction", IEEE Transactions on Knowledge and Data Engineering, vol.35, no.10, pp.9973-9984, 2023.[Show Article](#) [Google Scholar](#)
- Erlong Tan, Haoyang Yan, Kaiqi Zhao, Xiaolei Ma, Zhenliang Ma, Yuchuan Du, "MVOPFAD: Multiview Online Passenger Flow Anomaly Detection", IEEE Sensors Journal, vol.24, no.9, pp.14668-14681, 2024.[Show Article](#) [Google Scholar](#)
- Shuxin Zhang, Jinlei Zhang, Lixing Yang, Chengcheng Wang, Ziyou Gao, "COV-STFormer for Short-Term Passenger Flow Prediction During COVID-19 in Urban Rail Transit Systems", IEEE Transactions on Intelligent Transportation Systems, vol.25, no.5, pp.3793-3811, 2024.[Show Article](#) [Google Scholar](#)

Libraries and Frameworks

- **Django:** Web framework for Python. Django Documentation
- **Pandas:** Data manipulation and analysis library. Pandas Documentation
- **Scikit-learn:** Machine learning library. Scikit-learn Documentation
- **Joblib:** Load and Create Models. Joblib Documentation
- **Chart.js:** JavaScript library for charts and graphs. Chart.js Documentation

Dataset

- **NYC Subway Traffic 2017-21:** Available on Kaggle. [Dataset Link](#)