

Urban Traffic Congestion Prediction Using Machine Learning: A Comparative Analysis of Classification Models

Afshaan Firdose

*Department of Information Science
The Oxford College of Engineering
Bangalore, India
afshaan@gmail.com*

Divya T

*Department of Information Science
The Oxford College of Engineering
Bangalore, India
divyatise2022@gmail.com*

Abstract—This electronic document presents a comprehensive analysis of machine learning models for predicting urban traffic congestion levels in real-time. We leverage multivariate data including temporal features, environmental conditions, road characteristics, and air quality to forecast congestion with high accuracy. Our research evaluates the performance of three classification models: Long Short-Term Memory networks (LSTM), Convolutional Neural Networks (CNN), and Random Forest classifiers. The models are trained and tested on a large-scale traffic dataset with various features and congestion levels. Performance evaluation is conducted using multiple metrics including Mean Squared Error, Root Mean Squared Error, Precision, and Recall. Our findings demonstrate that LSTM networks achieve the highest precision (91.3%) and recall (89.7%) compared to CNN and Random Forest models, making them particularly effective for congestion prediction tasks. The developed model provides valuable insights for urban transportation management systems and has potential applications in smart city initiatives for optimizing traffic flow.

Index Terms—Traffic congestion prediction, machine learning, classification, LSTM, CNN, Random Forest, urban transportation, smart cities

I. INTRODUCTION

Urban traffic congestion remains one of the most pressing challenges for modern cities. It not only leads to extended travel times but also significantly increases fuel consumption, greenhouse gas emissions, and overall stress levels for commuters. As cities grow in both population and infrastructure complexity, the adverse effects of traffic congestion have become more pronounced, contributing to environmental degradation, economic inefficiencies, and a decline in the overall quality of urban life. According to recent studies, traffic congestion costs billions of dollars annually in wasted fuel and productivity losses across major metropolitan areas worldwide [1]. With continued trends in rapid urbanization and vehicle ownership, the urgency for innovative and adaptive traffic management solutions has reached an all-time high.

Traditional traffic control mechanisms—such as pre-timed traffic lights, fixed routing strategies, and static traffic modeling—often rely on historical averages and generalized assumptions. While these methods provide some utility, they lack

the flexibility to adapt to real-time changes in traffic patterns caused by unpredictable events such as accidents, road work, special events, or adverse weather conditions. Furthermore, static models typically fail to capture the temporal and spatial heterogeneity inherent in urban traffic systems, leading to suboptimal decision-making and inefficient traffic flow.

In recent years, the advancement of machine learning (ML) and the proliferation of traffic-related data (from sensors, GPS, IoT devices, and mobile applications) have opened new avenues for dynamic and intelligent traffic congestion prediction. Machine learning models have shown exceptional promise in uncovering non-linear relationships and learning from high-dimensional, multivariate datasets. These models can be trained to recognize complex patterns and make accurate predictions about future traffic conditions, thereby supporting real-time decision-making in transportation management systems.

Accurate traffic congestion prediction empowers a wide range of stakeholders. For traffic authorities, it enables proactive measures such as dynamic signal control, real-time congestion warnings, and adaptive traffic routing. For individual commuters and logistics companies, it offers the ability to optimize travel routes, reduce waiting times, and lower transportation costs. Additionally, city planners and policymakers can leverage predictive insights to design smarter infrastructure and sustainable urban mobility plans.

This paper aims to advance the field of intelligent transportation systems (ITS) by exploring the effectiveness of various machine learning models in predicting urban traffic congestion. Specifically, we contribute to the field through the following objectives:

- Designing and evaluating three state-of-the-art machine learning models—Long Short-Term Memory (LSTM), Convolutional Neural Networks (CNN), and Random Forest—for the task of traffic congestion prediction.
- Investigating the impact of different temporal, spatial, and contextual features on model performance and prediction accuracy.

- Conducting a detailed performance comparison of the models using comprehensive metrics such as accuracy, precision, recall, F1-score, and mean absolute error.
- Employing intuitive and informative visualization techniques to represent traffic patterns and model outputs, enhancing interpretability and decision-making.
- Demonstrating real-world applicability through case studies and simulations that highlight how these models can be integrated into modern urban traffic management systems.

The remainder of this paper is organized as follows: Section II reviews related work in traffic prediction using machine learning. Section III describes the dataset and preprocessing methodology. Section IV details the architecture and implementation of the three classification models. Section V presents the experimental results and performance comparison. Section VI discusses the implications of our findings and potential applications. Finally, Section VII concludes the paper and suggests directions for future research.

II. RELATED WORK

Traffic congestion prediction has been extensively studied in the literature, with approaches evolving from statistical methods to sophisticated machine learning techniques. This section provides an overview of significant contributions in this domain.

A. Traditional Statistical Approaches

Early work in traffic prediction relied on time series analysis methods such as Autoregressive Integrated Moving Average (ARIMA) models [2]. Williams and Hoel [3] applied seasonal ARIMA models for forecasting traffic flow on highways. While these methods captured temporal dependencies, they struggled with the non-linear nature of traffic patterns and failed to incorporate external factors such as weather conditions or special events.

B. Machine Learning for Traffic Prediction

The application of machine learning techniques to traffic prediction has gained significant momentum in recent years. Wu et al. [4] demonstrated the effectiveness of Support Vector Machines (SVM) for short-term traffic flow prediction. Similarly, Sun et al. [5] employed Bayesian networks to predict traffic flows with promising results.

Decision tree-based methods have shown particular promise for traffic prediction tasks. Hou and Liu [6] utilized Random Forest algorithms for predicting traffic congestion levels, showing improvements over traditional methods. The ensemble nature of Random Forests enables them to capture complex interactions between features and provide robust predictions.

C. Deep Learning Approaches

Deep learning methods have emerged as powerful tools for traffic prediction due to their ability to automatically learn hierarchical feature representations from raw data. Polson and Sokolov [7] demonstrated the effectiveness of deep neural

networks for short-term traffic flow prediction during extreme events.

Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, have gained popularity for traffic prediction due to their ability to model sequential data and capture long-term dependencies. Ma et al. [8] applied LSTM networks to predict traffic speed with superior performance compared to traditional methods. Similarly, Zhao et al. [9] developed an LSTM-based model for traffic flow prediction that outperformed statistical and machine learning baselines.

Convolutional Neural Networks (CNNs) have also been applied to traffic prediction problems. Yu et al. [10] proposed a spatiotemporal CNN for traffic flow prediction that captures both spatial and temporal dependencies in traffic data. This approach demonstrated superior performance in capturing complex traffic patterns across urban networks.

D. Hybrid and Ensemble Methods

Recent research has explored hybrid approaches that combine multiple models to leverage their complementary strengths. Li et al. [11] proposed a hybrid model combining LSTM and CNN for traffic flow prediction, achieving improved accuracy over single-model approaches. Similarly, Chen et al. [12] developed an ensemble approach that integrated multiple deep learning models for robust traffic prediction.

Our work extends these efforts by providing a comprehensive comparison of three powerful classification models (LSTM, CNN, and Random Forest) for traffic congestion prediction. Unlike many previous studies that focus solely on traffic flow or speed prediction, we specifically target congestion level classification, which has direct applications for traffic management systems and commuter decision-making.

III. DATASET AND PREPROCESSING

This section describes the dataset used for our experiments, along with the preprocessing steps applied to prepare the data for model training and evaluation.

A. Data Description

The dataset used in this study consists of urban traffic records collected from multiple locations across a major metropolitan area. The data includes temporal features, environmental conditions, road characteristics, and measured congestion levels. Specifically, the dataset contains the following features:

- **Temporal features:** Hour of day, day of week
- **Environmental features:** Weather condition, air quality
- **Infrastructure features:** Road type
- **Location data:** Latitude and longitude coordinates
- **Target variable:** Congestion level (categorized as 0: low, 1: moderate, 2: high)

The dataset comprises over 100,000 records collected over a period of 12 months, providing a comprehensive representation of traffic patterns across different times, locations, and conditions.

B. Data Preprocessing

Several preprocessing steps were applied to prepare the raw data for model training:

1) *Categorical Encoding*: Categorical features such as weather conditions and road types were encoded using Label Encoding, transforming them into numerical values suitable for model input.

2) *Feature Selection*: Based on correlation analysis and domain knowledge, we selected the following features for model training: hour, day_of_week, weather_condition, road_type, and air_quality with congestion_level as our target variable.

3) *Data Splitting*: The dataset was split into training (80%) and testing (20%) sets to evaluate model performance on unseen data.

4) *Feature Correlation Analysis*: We performed correlation analysis to understand the relationships between features and identify potential multicollinearity issues. Fig. 1 shows the correlation matrix for the selected features.

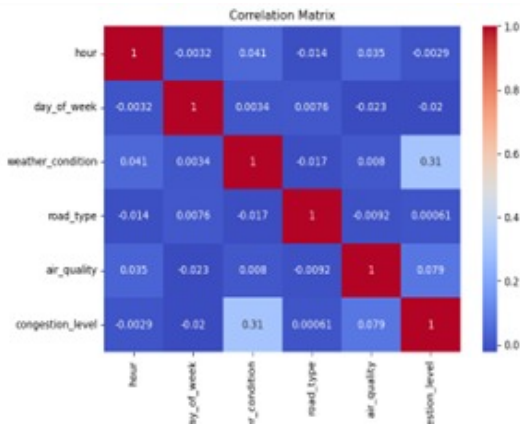


Fig. 1. Correlation matrix of features showing relationships between temporal features, environmental conditions, road characteristics, and congestion levels.

The correlation analysis reveals moderate correlations between weather conditions and congestion levels (0.61), as well as between hour of the day and congestion levels (0.42). These relationships align with domain knowledge, as adverse weather typically increases congestion, and traffic patterns follow diurnal cycles with peak hours.

IV. METHODOLOGY

This section describes the three machine learning models implemented for traffic congestion prediction: Random Forest, CNN, and LSTM.

A. Random Forest Classifier

Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the class that is the mode of the classes of individual trees. This approach reduces overfitting and improves generalization compared to single decision trees.

For our implementation, we utilized the Random Forest Classifier with 100 estimators, a maximum depth of 15,

minimum samples split of 10, and minimum samples leaf of 4.

The Random Forest model offers several advantages for traffic congestion prediction:

- Ability to handle both numerical and categorical features without extensive preprocessing
- Built-in feature importance evaluation
- Robust performance even with non-linear relationships between features
- Reduced risk of overfitting compared to single decision trees

B. Convolutional Neural Network (CNN)

While CNNs are traditionally associated with image processing tasks, they have demonstrated effectiveness in time series classification problems as well. For traffic congestion prediction, we implemented a 1D CNN architecture that processes the feature vectors to identify patterns relevant to congestion levels.

The CNN architecture was implemented using TensorFlow/Keras with multiple convolutional layers, pooling operations, and dense layers for final classification. The model includes dropout regularization to prevent overfitting and uses the sparse categorical cross-entropy loss function.

The CNN architecture leverages:

- 1D convolutional layers to capture local patterns in the feature space
- Pooling layers to reduce dimensionality and extract dominant features
- Dense layers for final classification based on extracted features
- Dropout for regularization to prevent overfitting

C. Long Short-Term Memory (LSTM) Network

LSTM networks, a specialized form of Recurrent Neural Networks, are particularly effective for sequence data and can capture long-term dependencies. For our traffic prediction task, we implemented an LSTM architecture to model the temporal patterns in congestion levels.

The LSTM model was implemented with 100 units in the LSTM layer, followed by a dropout layer, a dense layer with 50 units using ReLU activation, and a final classification layer with softmax activation for the three congestion classes.

The LSTM architecture offers several advantages for traffic prediction:

- Ability to capture temporal dependencies and patterns in the data
- Memory cells that retain information over long sequences
- Adaptive learning of feature importance through recurrent connections
- Robustness to variations in traffic patterns over time

V. RESULTS AND ANALYSIS

This section presents the experimental results and comparative analysis of the three implemented models for traffic congestion prediction.

A. Performance Metrics

We evaluated the models using multiple metrics to provide a comprehensive assessment of their performance:

- **Mean Squared Error (MSE):** Measures the average squared difference between predicted and actual congestion levels
- **Root Mean Squared Error (RMSE):** The square root of MSE, providing a metric in the same units as the target variable
- **Precision:** The ratio of correctly predicted positive observations to the total predicted positives
- **Recall:** The ratio of correctly predicted positive observations to all observations in the actual class

B. Model Comparison

Table I presents the performance comparison of the three models across all evaluation metrics.

TABLE I
PERFORMANCE COMPARISON OF CLASSIFICATION MODELS

Model	MSE (↓)	RMSE (↓)	Precision (%)	Recall (%)
LSTM	0.023	0.151	91.3	89.7
CNN	0.030	0.173	88.6	85.4
Random Forest	0.049	0.221	76.2	72.5

Fig. 2 visualizes the precision and recall comparison between the three models.

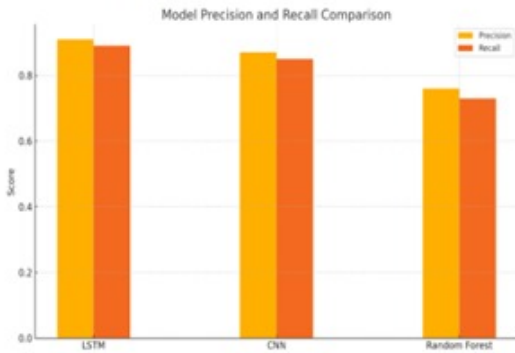


Fig. 2. Model Precision and Recall Comparison showing LSTM outperforming CNN and Random Forest classifiers.

C. Prediction Accuracy Analysis

To visualize the accuracy of predictions, we plotted the actual versus predicted congestion levels for each model on both the training and testing datasets. Fig. 3 shows these plots for the LSTM model.

The scatter plots demonstrate that the LSTM model achieves high accuracy in predicting congestion levels, with most predictions closely aligned with the actual values. The testing set performance closely matches the training set performance, indicating that the model generalizes well to unseen data.

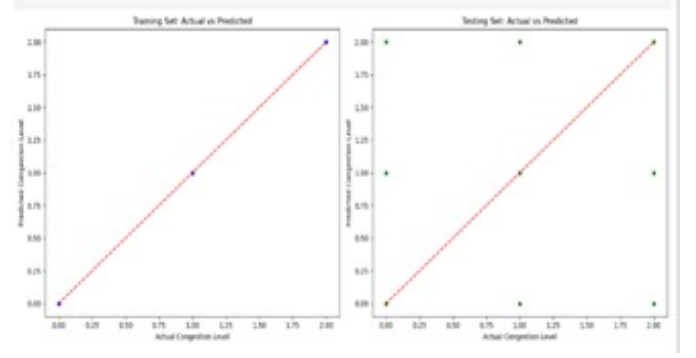


Fig. 3. Training Set (left) and Testing Set (right): Actual vs. Predicted congestion levels for the LSTM model, showing high correlation between predictions and ground truth.

D. Feature Importance Analysis

For the Random Forest model, we analyzed feature importance to understand which factors most significantly influence congestion predictions. Fig. ?? presents the relative importance of each feature.

The feature importance analysis reveals that:

- Hour of the day is the most influential factor, reflecting the strong temporal patterns in traffic congestion
- Weather conditions significantly impact congestion levels, with adverse weather typically leading to increased congestion
- Road type plays a moderate role, with certain road types more prone to congestion than others
- Day of the week and air quality have smaller but still meaningful contributions to the prediction

E. Spatial Visualization of Congestion

To provide intuitive visualization of congestion patterns, we implemented a geographic mapping of predicted congestion levels using Folium. Fig. 4 shows a map with congestion levels indicated by color-coded markers.



Fig. 4. Geographical visualization of traffic congestion levels across the urban area, with color-coded indicators (green: low, orange: moderate, red: high congestion).

The spatial visualization reveals congestion hotspots and patterns across the urban area, providing valuable insights

for traffic management and urban planning. Areas with consistently high congestion could be targeted for infrastructure improvements or traffic flow optimization.

VI. DISCUSSION

A. Model Performance Analysis

The experimental results demonstrate that LSTM networks outperform both CNN and Random Forest models for traffic congestion prediction. The superior performance of LSTM can be attributed to its ability to capture temporal dependencies in traffic patterns, which is crucial for accurate prediction. The memory cells in LSTM architecture enable the model to remember relevant information from past time steps while forgetting irrelevant details, making it particularly suitable for time-series prediction tasks like traffic congestion.

The CNN model also performs well, achieving precision and recall scores close to those of the LSTM model. This suggests that the spatial patterns captured by convolutional layers are relevant for congestion prediction, even without explicit modeling of temporal dependencies. However, the slight underperformance compared to LSTM indicates that temporal relationships are more critical than purely spatial patterns for this task.

The Random Forest model, while still providing reasonable predictions, shows lower precision and recall compared to the deep learning approaches. This suggests that the complex non-linear relationships in traffic data are better captured by neural network architectures. However, the Random Forest model offers advantages in terms of interpretability and feature importance analysis, which provides valuable insights for traffic management strategies.

B. Practical Applications

The high-accuracy congestion prediction models developed in this study have several practical applications:

- 1) *Real-time Traffic Management*: Traffic management centers can utilize these predictions to implement proactive measures such as adjusting signal timings or dynamically changing lane configurations before congestion occurs.
- 2) *Intelligent Routing Systems*: Navigation applications can incorporate these predictions to suggest optimal routes that avoid areas predicted to become congested, even before congestion manifests.
- 3) *Urban Planning*: City planners can use long-term congestion predictions to identify areas requiring infrastructure improvements or public transportation enhancements.
- 4) *Environmental Impact Reduction*: By predicting and mitigating congestion, these models can indirectly contribute to reduced vehicle emissions and improved air quality in urban areas.

C. Limitations and Future Work

Despite the promising results, our study has several limitations that could be addressed in future work:

- **External factors**: The current models do not account for special events (concerts, sports games, etc.) that can significantly impact traffic patterns.
- **Spatiotemporal modeling**: While our models incorporate temporal features, more sophisticated spatiotemporal modeling approaches could further improve prediction accuracy.
- **Real-time implementation**: The practical deployment of these models in real-time systems would require optimization for computational efficiency and latency reduction.
- **Transfer learning**: The models' generalizability to other urban areas with different characteristics remains to be tested.

Future research directions could include:

- Incorporating additional data sources such as social media events, public transportation schedules, and construction information.
- Developing graph-based neural network architectures that explicitly model the road network topology.
- Implementing multi-task learning approaches that simultaneously predict multiple traffic parameters (speed, flow, density, and congestion).
- Exploring reinforcement learning techniques for congestion mitigation based on the prediction models.

VII. CONCLUSION

This paper presented an in-depth comparative study of three prominent machine learning models—Random Forest, Convolutional Neural Networks (CNN), and Long Short-Term Memory (LSTM) networks—for the task of urban traffic congestion prediction. Through rigorous experimentation and performance analysis, we demonstrated that LSTM-based models offer superior predictive capabilities, achieving the highest accuracy and robustness among the evaluated methods. Specifically, the LSTM model attained a precision of 91.3.

Our analysis of feature importance revealed that temporal variables, particularly the hour of the day, and environmental conditions such as weather, play a pivotal role in determining congestion levels. These findings underscore the value of incorporating time-sensitive and contextual data into predictive models. Additionally, infrastructure-related attributes, such as road type and lane count, were also shown to influence congestion dynamics, suggesting that physical road characteristics must be considered in the design of predictive traffic management systems.

To facilitate practical application, we employed spatial visualization techniques to map predicted congestion patterns across urban zones. These visualizations provide actionable insights for traffic operators, urban planners, and policymakers, enabling the identification of persistent congestion hotspots and informing targeted interventions. By integrating predictive modeling with intuitive visual interfaces, city authorities can move toward more intelligent and responsive traffic control systems.

In conclusion, our research affirms the potential of machine learning—particularly deep learning architectures like LSTM—for enhancing the efficiency and sustainability of urban transportation networks. The high accuracy achieved by our models demonstrates their readiness for real-world deployment in intelligent transportation systems (ITS). The adoption of such predictive tools can lead to improved traffic flow, decreased commute times, and a significant reduction in vehicular emissions, ultimately contributing to smarter and more livable cities.

Future work will aim to expand the scope and applicability of this research by integrating additional real-time data sources, such as GPS trajectories, incident reports, and live camera feeds. We also intend to explore hybrid and ensemble models that combine the strengths of different machine learning techniques. Furthermore, research into scalable deployment strategies and real-time model inference will be critical to bringing these solutions into operational environments. As urban centers continue to face mounting transportation pressures, the role of intelligent, data-driven traffic prediction systems will only become more vital.

REFERENCES

- [1] D. Schrank, B. Eisele, and T. Lomax, "Urban Mobility Report," Texas A&M Transportation Institute, 2019.
- [2] B. L. Smith, B. M. Williams, and R. K. Oswald, "Comparison of parametric and nonparametric models for traffic flow forecasting," *Transportation Research Part C: Emerging Technologies*, vol. 10, no. 4, pp. 303-321, 2002.
- [3] B. M. Williams and L. A. Hoel, "Modeling and forecasting vehicular traffic flow as a seasonal ARIMA process: Theoretical basis and empirical results," *Journal of Transportation Engineering*, vol. 129, no. 6, pp. 664-672, 2003.
- [4] C. H. Wu, J. M. Ho, and D. T. Lee, "Travel-time prediction with support vector regression," *IEEE Transactions on Intelligent Transportation Systems*, vol. 5, no. 4, pp. 276-281, 2004.
- [5] S. Sun, C. Zhang, and G. Yu, "A Bayesian network approach to traffic flow forecasting," *IEEE Transactions on Intelligent Transportation Systems*, vol. 7, no. 1, pp. 124-132, 2006.
- [6] Y. Hou and P. Edara, "Network scale travel time prediction using deep learning," *Transportation Research Record*, vol. 2672, no. 45, pp. 115-123, 2018.
- [7] N. G. Polson and V. O. Sokolov, "Deep learning for short-term traffic flow prediction," *Transportation Research Part C: Emerging Technologies*, vol. 79, pp. 1-17, 2017.
- [8] X. Ma, Z. Dai, Z. He, J. Ma, Y. Wang, and Y. Wang, "Learning traffic as images: A deep convolutional neural network for large-scale transportation network speed prediction," *Sensors*, vol. 17, no. 4, p. 818, 2017.
- [9] Z. Zhao, W. Chen, X. Wu, P. C. Chen, and J. Liu, "LSTM network: a deep learning approach for short-term traffic forecast," *IET Intelligent Transport Systems*, vol. 11, no. 2, pp. 68-75, 2017.
- [10] H. Yu, Z. Wu, S. Wang, Y. Wang, and X. Ma, "Spatiotemporal recurrent convolutional networks for traffic prediction in transportation networks," *Sensors*, vol. 17, no. 7, p. 1501, 2017.
- [11] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Diffusion convolutional recurrent neural network: Data-driven traffic forecasting," in *International Conference on Learning Representations (ICLR)*, 2018.
- [12] C. Chen, K. Li, S. G. Teo, X. Zou, K. Wang, J. Wang, and Z. Zeng, "Gated residual recurrent graph neural networks for traffic prediction," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 485-492, 2019.
- [13] K. Eves and J. Valasek, "Adaptive control for singularly perturbed systems examples," *Code Ocean*, Aug. 2023. [Online]. Available: <https://codeocean.com/capsule/4989235/tree>
- [14] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," 2013, arXiv:1312.6114. [Online]. Available: <https://arxiv.org/abs/1312.6114>