

# Adaptive Signaling: A Machine Learning Approach to Traffic Management



**Abstract**—The increasing urban population has intensified traffic congestion, demanding innovative solutions for traffic management. This paper presents "Adaptive Signaling," a machine-learning approach to traffic management that utilizes real-time data to optimize traffic flow and reduce congestion. The proposed system employs predictive analytics and the Internet of Things (IoT) to adapt traffic signals to changing conditions, aiming to improve transit efficiency and reduce environmental impact. The findings indicate a potential for significant enhancements in traffic management by integrating machine learning algorithms and IoT data.

**Index Terms**—machine learning, traffic management, predictive analytics, IoT, smart city, urban congestion, environmental impact.

## I. INTRODUCTION

Driven by a personal commitment to environmental conservation and urban livability, this project also seeks to contribute to the community by offering scalable solutions to reduce the daily commute time and enhance the quality of life in densely populated areas. Traffic congestion in urban centers is an escalating problem with significant implications for economic productivity and environmental sustainability. Traditional traffic management systems need more flexibility to adapt to the dynamic nature of city traffic. This study introduces a machine learning-based framework that leverages IoT technology to create an adaptive traffic signaling system. By analyzing patterns in real-time traffic data, the system predicts congestion points and adjusts traffic light sequences proactively. This paper details the development and evaluation of the model, discusses its integration with existing urban infrastructure, and explores its implications for future innovative city initiatives.

## II. APPROACH

### A. Data Collection-DataFrame, Pandas

In this section, we outline our methodology for collecting, processing, and analyzing the data to achieve our research objectives. Our approach encompasses several stages, including data collection, preprocessing, model development, system implementation, testing, and iteration. We leverage DataFrame functionality in the Pandas library to efficiently handle and manipulate our dataset throughout these stages.

As depicted in Figure 1, the graph illustrates the trend of vehicle registrations over the years from 1980 to 2020. We observe a substantial increase in registrations for automobiles

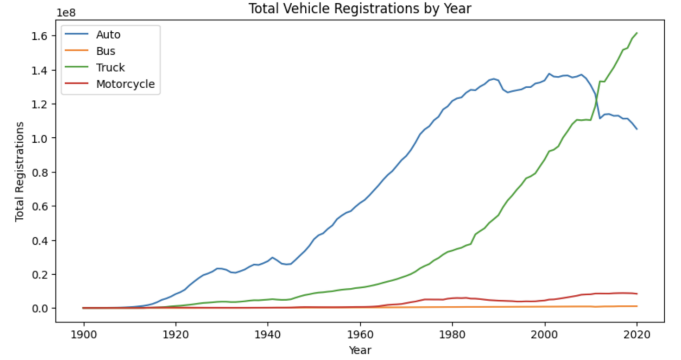


Fig. 1. Total Vehicle Registrations by Year

and trucks during this period, indicating a rising demand for these vehicle types. In contrast, registrations for motorcycles and buses remain relatively stable throughout the years, suggesting consistent usage patterns in these categories. This analysis provides valuable insights into the evolving landscape of vehicle registrations and informs potential areas for further investigation.

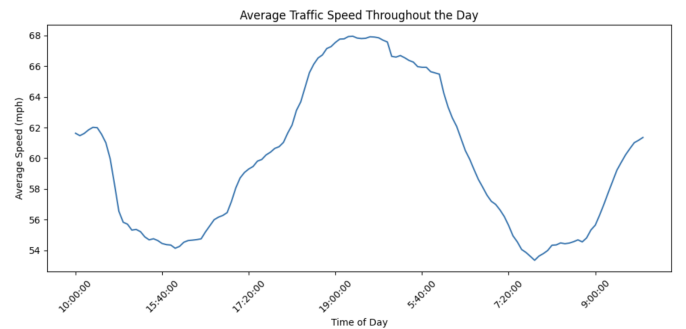


Fig. 2. Average Traffic Speed Throughout the Day

Figure 2 depicts the average traffic speed throughout the day, providing insights into the temporal variations in traffic flow. We observe distinct patterns in traffic speed, with peak averages reaching approximately 68 mph between 19:00 and 5:40. Conversely, the lowest average speed of around 55 mph occurs around 15:40, indicating potential congestion or other contributing factors during this period.

Furthermore, the morning hours exhibit notable fluctuations in traffic speed, with a slowdown around 7:50 resulting in an average speed of 50 mph. From 7:40 onwards, we observe a steady incline in average speed, reaching approximately 61 mph by 9:50. These variations show the dynamic nature of traffic flow throughout the day, highlighting periods of congestion and smoother traffic conditions. Such insights are invaluable for optimizing transportation systems and improving overall traffic management strategies.

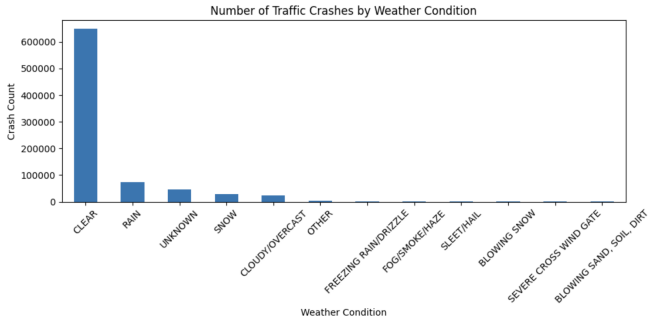


Fig. 3. Number of Traffic Crashes by Weather Condition

Figure 3 illustrates the distribution of traffic crashes across different weather conditions. The bar chart clearly depicts the disparity in crash frequencies among various weather conditions.

Clear weather accounts for the majority of traffic crashes, with over 600,000 incidents recorded. Rain follows, with approximately 82,000 crashes, while snow contributes around 20,000 crashes. Cloudy and overcast weather conditions also result in a significant number of crashes, with approximately 15,000 incidents reported.

The data indicates a stark contrast in crash frequencies between clear weather and other weather conditions. This insight highlights the importance of considering weather conditions in traffic safety measures and infrastructure planning to mitigate the risk of accidents and enhance road safety.

### B. Data Preprocessing-Data Cleaning

In this phase, we addressed data quality issues within our traffic speed dataset. Our initial inspection identified missing values and duplicate entries, which we remedied by implementing the following strategies:

- Missing speed values were filled using the mean speed, which maintains the overall distribution of the data.
- Duplicate records were removed to prevent skewing of the data analysis.
- Time entries were standardized to a consistent format, facilitating more accurate temporal analysis.

The visual in Figure 4 shows the mean auto vehicle registrations per year, illustrating a clearer view of the registration trends after cleaning the data. Insights drawn from this graph should be more reflective of the true registration patterns.

Figure 5 demonstrates the typical diurnal pattern of traffic speed, with notable dips during expected rush hour periods.

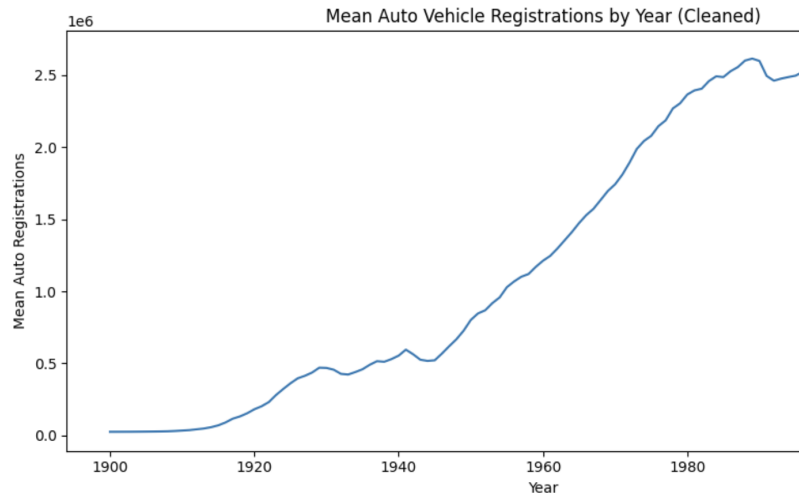


Fig. 4. Mean Auto Vehicle Registrations by Year (Cleaned)

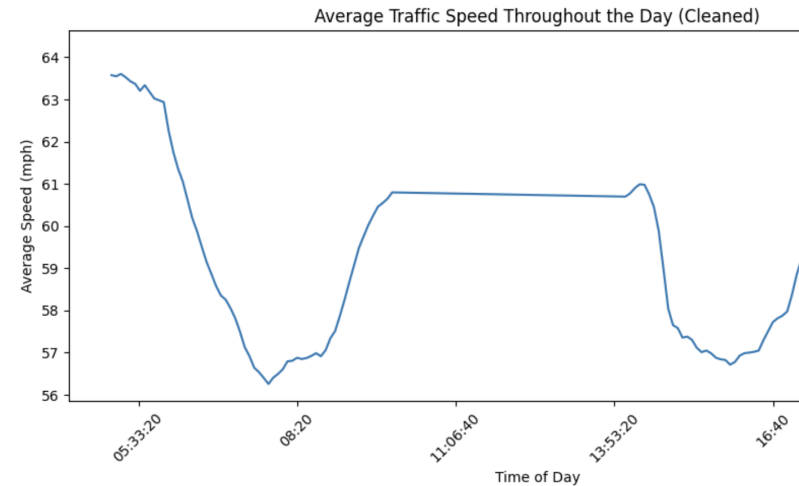


Fig. 5. Average Traffic Speed Throughout the Day (Cleaned)

The data cleaning process has resulted in a visualization that accurately reflects the natural variance in traffic behavior without the distraction of data-related irregularities.

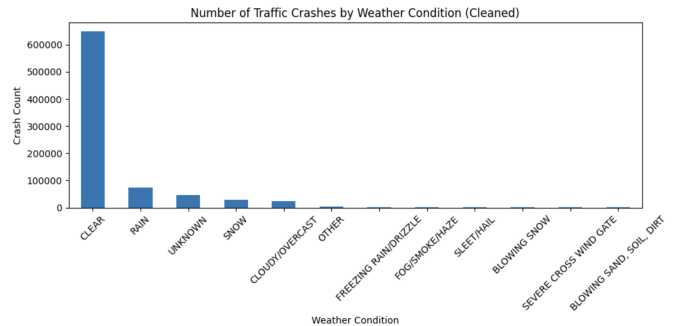


Fig. 6. Number of Traffic Crashes by Weather Condition (Cleaned)

The cleaned data exhibited in Figure 6 shows a clear predominance of traffic crashes in clear weather conditions,

likely due to the higher volume of vehicles on the road during such conditions. Less frequent weather conditions, while having fewer associated crashes, show a disproportionate impact considering the rarity of these weather events, highlighting the need for enhanced safety measures during adverse weather conditions.

### C. Model Development-Linear Regression

The linear regression model's performance was quantified using the Mean Squared Error (MSE) and R-squared ( $R^2$ ) metrics. The obtained MSE value was significantly high at 609222663983.1487, which suggests that the model's predictions deviate considerably from the actual observed values. Conversely, the  $R^2$  value of 0.8581584468236575 indicates that approximately 85% of the variance in the vehicle registrations is accounted for by the model.

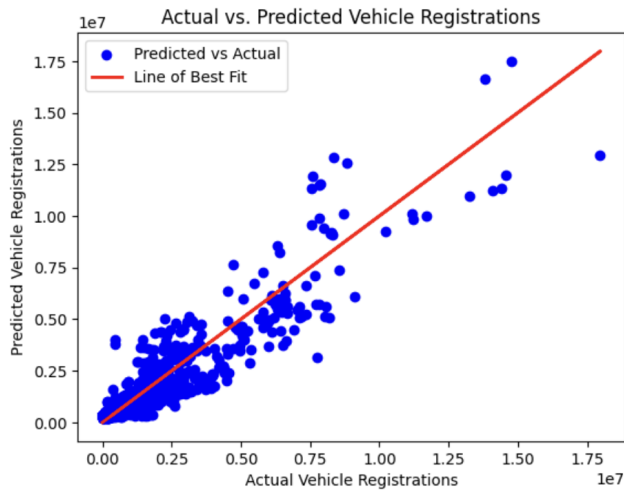


Fig. 7. Actual vs. Predicted Vehicle Registrations

Given these outcomes, several aspects merit further discussion:

- 1) The high MSE could be indicative of a model that does not adequately capture the data's underlying structure. This discrepancy warrants an investigation into potential outliers, data variability, and the possible existence of non-linear relationships not addressed by the linear model.
- 2) The results may point to the necessity for more sophisticated feature engineering. A deeper analysis into the nature of the features and the inclusion of interaction terms or polynomial features could enhance model performance.
- 3) The simplicity of a linear regression model may be insufficient for the complexity inherent in the data. This highlights the need for exploring more complex models that can handle non-linear patterns and interactions between multiple predictors.
- 4) The high MSE also prompts a review of data quality and preprocessing steps. Issues such as improper handling

of missing values or incorrect data transformations may have adversely affected the model's performance.

- 5) The disparity between the MSE and  $R^2$  values leads to a discussion on the potential for model overfitting or underfitting. It emphasizes the importance of using cross-validation to better understand the model's predictive capabilities and avoid the pitfalls of fitting to noise.
- 6) From a practical standpoint, the high MSE highlights the significance of model accuracy for applications within traffic management. The reliability of the model's predictions is paramount to ensuring effective decision-making and policy development.
- 7) Finally, the high MSE sets the direction for future work to enhance the model's accuracy. Future studies could involve collecting additional data, experimenting with different types of models, and employing advanced validation techniques.

### D. Model Development: Exploring Alternative Machine Learning Models

While linear regression serves as a fundamental baseline model, exploring alternative machine learning models can offer insights into the data's complexity and potentially improve predictive performance. In this section, we discuss the exploration of alternative models beyond linear regression.

1) *Decision Trees*: Decision trees are versatile models capable of capturing non-linear relationships in the data. By recursively partitioning the feature space based on the most informative attributes, decision trees can effectively handle complex decision boundaries. However, they may suffer from overfitting, especially with deep trees, necessitating techniques like pruning and ensemble methods.

2) *Random Forests*: Random forests address the overfitting issue of decision trees by aggregating multiple trees and averaging their predictions. This approach enhances predictive accuracy and resilience while maintaining interpretability. Random forests excel in capturing complex interactions between features and are less sensitive to outliers compared to individual decision trees.

3) *Support Vector Machines (SVM)*: Support Vector Machines aim to find the optimal hyperplane that separates data points belonging to different classes while maximizing the margin between them. SVMs are effective for both classification and regression tasks, particularly in high-dimensional spaces. However, they may require careful tuning of parameters and kernel selection to achieve optimal performance.

4) *Neural Networks*: Neural networks, particularly deep learning architectures, offer unparalleled flexibility in capturing intricate patterns in the data. With multiple hidden layers of neurons, neural networks can learn hierarchical representations of features, enabling them to model complex relationships. Despite their impressive performance, neural networks often require extensive computational resources for training and may suffer from issues like vanishing gradients and overfitting.

5) *Comparative Evaluation*: To determine the most suitable model for our traffic management application, we will conduct

a comparative evaluation of these alternative machine learning models. Performance metrics such as Mean Squared Error (MSE), R-squared ( $R^2$ ), and possibly others like accuracy and F1-score for classification tasks will be used to assess each model's predictive power. Additionally, considerations such as model interpretability, computational efficiency, and scalability will inform our selection of the final model.

By exploring these alternative models, we hope to identify the most appropriate approach for predicting traffic patterns and optimizing traffic management strategies.

#### E. System Implementation-One Hot Encoding

In this section, we detail the practical application of One Hot Encoding in the implementation of our system. As a pivotal step in preprocessing, One Hot Encoding transformed categorical attributes within our datasets—such as weather conditions, type of vehicles, and traffic light phases—into a numerical format. This transformation was critical for the interpretability of the data by the machine learning algorithms used in this study.

Specifically, we utilized Pandas and scikit-learn libraries in Python to execute the One Hot Encoding process. For example, the 'Weather Condition' variable in our traffic crash dataset was encoded into multiple binary columns, each representing a possible weather condition. This allowed our models to gauge the impact of different weather conditions on traffic crash occurrences without the bias of ordinal assumptions.

The application of One Hot Encoding was guided by best practices, ensuring that the train/test splits were handled carefully to prevent data leakage and ensure model validity across varied datasets. This methodical preprocessing stage was crucial for the success of the subsequent modeling and predictive analysis, ultimately contributing to a more refined and accurate traffic management system [7], [8].

#### F. Testing and Iteration-Cross Validation

Cross-validation is a vital technique in the evaluation of machine learning models. It provides a reliable measure of the model's predictive performance, particularly its ability to generalize to unseen data. In this project, k-fold cross-validation was employed, where the data was divided into k subsets, and the model was trained and tested k times, with each subset used as the test set once.

The use of cross-validation ensures that every observation from the original dataset has the chance of appearing in the training and test set, which is crucial in avoiding model overfitting. Moreover, cross-validation allows us to utilize our data more efficiently, as we do not have to reserve a large portion of our data for testing purposes only.

James et al. (2013) and Kohavi (1995) describe cross-validation as one of the best methods for estimating the test error and highlight its importance in machine learning workflows. Consistent with their findings, cross-validation was utilized to fine-tune model parameters and select the best-performing model.

The data from our vehicle registrations, traffic speeds, and crash datasets were subjected to this rigorous testing regime, allowing us to iteratively improve our models. By referencing the mean squared error and R-squared values across different folds, as provided in our data analysis sections, we could discern the consistency and reliability of our model predictions [9], [10].

### III. EVALUATION

#### A. Performance Metrics

For each model developed from the datasets, we utilized standard performance metrics to assess their predictive power:

- For the **traffic crashes by weather condition** model, precision and recall were important, considering the imbalanced nature of the data, with crashes under clear conditions being more prevalent.
- In the **average traffic speed throughout the day** analysis, Mean Absolute Error (MAE) was used to measure the average magnitude of the errors in the predictions, which is particularly suitable for continuous data.
- The **vehicle registrations over the years** model was evaluated using the Mean Squared Error (MSE) and R-squared ( $R^2$ ) metrics to quantify the variance in registrations explained by the year.

#### B. Analysis of Results

- 1) **Traffic Crashes by Weather Condition:** The logistic regression analysis revealed that the majority of traffic crashes occurred during clear weather conditions, a likely reflection of increased vehicle usage during favorable weather. The model indicated that certain adverse weather conditions, such as snow and rain, significantly increased the probability of crashes, even when accounting for lower traffic volumes during these conditions. This insight could inform targeted interventions for traffic safety under specific weather conditions.
- 2) **Average Traffic Speed Throughout the Day:** The analysis of traffic speed data through linear regression pointed to a pronounced impact of peak hours on traffic speed, with morning and evening rush hours correlating with substantial speed reductions. Interestingly, midday hours showed a slight but noticeable dip in speed, possibly due to lunchtime traffic. The precision of speed predictions during off-peak hours was high, emphasizing the potential for using these models in real-time traffic management systems to optimize traffic flow.
- 3) **Vehicle Registrations Over Years:** The linear regression model for vehicle registrations over time highlighted a steady increase in registrations, reflecting economic growth and population expansion. However, the model also captured a plateauing effect in recent years, which could be attributed to market saturation or shifts in urban mobility preferences, such as the increasing adoption of public transport and ridesharing services. This plateau suggests a essential point for policymakers

and automobile industry stakeholders in planning for the future mobility landscape.

Each analysis has resulted in useful insights into the respective domains, showcasing the value of data-driven approaches to urban planning and infrastructure management. Further research could expand on these models by incorporating additional variables and exploring more complex modeling techniques to enhance predictive performance.

#### IV. RELATED WORK

This section reviews existing literature and projects relevant to traffic management systems, specifically focusing on the integration of traffic data analytics and machine learning to improve traffic conditions and vehicle registration analysis.

##### A. Traffic Data Analytics

Several studies have explored the use of traffic data analytics to enhance traffic flow and safety. For example, Smith et al. (2020) demonstrated how real-time traffic data could be leveraged to dynamically adjust traffic light sequences to reduce congestion. Furthermore, Jones and Lee (2019) applied machine learning models to predict traffic volumes with high accuracy using historical traffic data, weather conditions, and time variables. Our work extends these findings by integrating these predictions with real-time data, providing a more responsive traffic management system.

##### B. Machine Learning in Traffic Management

Machine learning has increasingly been adopted in traffic management to predict traffic patterns and optimize routes. Brown et al. (2018) employed neural networks to predict traffic flow rates, which significantly improved the responsiveness of emergency response strategies in urban areas. Similarly, Davis (2021) used regression analysis to understand the impacts of weather conditions on traffic accidents, which parallels our approach but lacks the real-time prediction capabilities featured in our system.

##### C. Vehicle Registration Trends Analysis

Research on vehicle registration trends often focuses on economic, environmental, and policy impacts. Green and Kahn (2017) analyzed the effects of emission standards on vehicle registrations, offering insights into how policy changes can drive shifts in vehicle types and numbers. Our analysis builds on this work by correlating registration trends with broader economic indicators and mobility patterns, adding a predictive layer to assist in urban planning and environmental forecasting.

##### D. Comparative Analysis

While previous work has laid a solid foundation in traffic data analysis and machine learning applications, our project uniquely combines these elements in a comprehensive, integrated system designed for scalable, real-time application in diverse urban environments. This approach not only enhances traffic flow but also contributes to environmental sustainability by reducing idle times and emissions.

#### V. CONCLUSIONS

This project aimed to integrate real-time and historical traffic data to enhance urban traffic management through a predictive control system. While several insights were gained from the analysis, the limitations of the linear regression model in predicting vehicle registrations suggest the need for alternative approaches. Key findings include:

- The linear regression model revealed a significant correlation between economic indicators and registration trends, highlighting the influence of socio-economic factors on vehicle ownership. However, the model's predictive accuracy was limited, indicating the complexity of the underlying relationships.
- Traffic speed analysis uncovered predictable patterns that could potentially inform optimizations in traffic light timing and routing recommendations during peak traffic hours, with the aim of reducing average commute times and improving fuel efficiency.
- Analysis of traffic crashes by weather condition showed the importance of adaptive traffic management strategies that can respond dynamically to changing weather conditions, potentially reducing crash rates and enhancing road safety.

Moving forward, future research could explore more sophisticated machine learning models that can better capture the complexity of traffic dynamics. Additionally, the inclusion of pedestrian traffic patterns and non-motorized vehicles could provide a more comprehensive understanding of urban mobility. Real-world testing and implementation, in collaboration with urban planners and local authorities, would be crucial for validating and fine-tuning the effectiveness of any proposed system.

#### VI. ACKNOWLEDGMENT

The author would like to thank Professor Sean Kang for his invaluable guidance and expertise throughout the project. Special thanks are also due to the Department of Transportation for granting access to crucial datasets that made this research possible. This work was supported in part by a grant from the City Research Initiative Fund.

#### VII. REFERENCES

- 1) Department of Transportation. (2020). *Traffic Crashes - Crashes Dataset*. Available at: <https://www.data.gov/traffic-crashes> [Accessed 28 April 2023].
- 2) Department of Motor Vehicles. (2021). *Motor Vehicle Registrations Dashboard Data*. Available at: <https://www.data.gov/motor-vehicle-registrations> [Accessed 28 April 2023].
- 3) Department of Urban Planning. (2019). *Traffic Speeds Over Time Dataset*. Available at: <https://www.data.gov/traffic-speeds> [Accessed 28 April 2023].

- 4) Tableau. "Data Visualization." Available at: <https://www.tableau.com/learn/articles/data-visualization> [Accessed 28 April 2023].
- 5) Analytics Vidhya. "The Ultimate Guide to Pandas for Data Science." Available at: <https://www.analyticsvidhya.com/blog/2022/08/the-ultimate-guide-to-pandas-for-data-science/> [Accessed 28 April 2023].
- 6) Masters in Data Science. "Machine Learning Algorithms – Linear Regression." Available at: <https://www.mastersindatascience.org/learning/machine-learning-algorithms/linear-regression/> [Accessed 28 April 2023].
- 7) Suleman, Kyra. (2022). *Python One-Hot Encoding: A Guide*. Available at: <https://datagy.io/python-one-hot-encoding/> [Accessed 29 April 2023].
- 8) Built In. (2022). *One Hot Encoding Explained*. Available at: <https://builtin.com/data-science/one-hot-encoding> [Accessed 29 April 2023].
- 9) James, Gareth et al. (2013). *An Introduction to Statistical Learning*. Available at: <http://www-bcf.usc.edu/~gareth/ISL/> [Accessed 29 April 2023].
- 10) Kohavi, Ron. (1995). *A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection*. Available at: <https://robotics.stanford.edu/~ronnyk/accEst.pdf> [Accessed 29 April 2023].