



# A DECISION TREE REGRESSOR APPROACH FOR PREDICTING UBER TRIP FARES

**Akhilash Pennam**

*Director of CRM Development*

Article DOI: <https://doi.org/10.36713/epra22501>

DOI No: 10.36713/epra22501

## ABSTRACT

The rapid expansion of ride-sharing services such as Uber has resulted in the generation of extensive trip and fare datasets, creating opportunities to develop accurate fare prediction models. In this study, a **Decision Tree Regressor** is employed to predict Uber ride fares based on key features including pickup and drop-off locations, trip distance, and passenger count. The simplicity and interpretability of the Decision Tree Regressor make it a suitable model for such regression tasks where transparency is crucial. The model was trained and evaluated on real-world Uber trip data, achieving impressive performance metrics with an **R<sup>2</sup> score of 0.9033**, indicating that over 90% of the variance in fare prices can be explained by the input variables. Furthermore, the model attained a **Mean Absolute Error (MAE) of 0.8081**, a **Mean Squared Error (MSE) of 2.3316**, and a **Root Mean Squared Error (RMSE) of 1.5269**, demonstrating its robustness and predictive capability. This research highlights the potential of Decision Tree Regression as an effective and interpretable approach for fare prediction in ride-sharing platforms, providing accurate real-time fare estimates that can improve both customer experience and operational efficiency. The findings suggest that lightweight and explainable models like Decision Tree Regressors can be viable alternatives to more complex machine learning models, especially when computational simplicity and model transparency are prioritized.

**KEYWORDS:** Decision Tree Regressor; Uber Fare Prediction; Machine Learning; Regression Analysis; Ride-Sharing Data; Predictive Modeling; Model Interpretability; Root Mean Squared Error (RMSE); R-Squared Score; Fare Estimation.

## 1. INTRODUCTION

The rapid growth of urbanization and the increasing demand for flexible transportation services have propelled the popularity of ride-sharing platforms such as Uber, Ola, and Lyft. These platforms generate massive amounts of data daily, encompassing trip details such as pickup and drop-off locations, travel distance, passenger count, and fare amounts. This data provides an opportunity to leverage machine learning techniques to build predictive models that can accurately estimate ride fares, which is critical for enhancing customer satisfaction and operational efficiency. Fare prediction is an important component of ride-hailing platforms as it impacts pricing transparency, surge pricing decisions, and user trust. Traditional fare estimation methods are often based on predefined rules or linear models that may not capture complex, non-linear relationships present in the data. In contrast, machine learning algorithms can learn such patterns and provide more accurate and dynamic fare predictions. Among various machine learning models, the **Decision Tree Regressor** is a widely used approach for regression tasks due to its simplicity, ease of interpretation, and ability to handle both numerical and categorical data without the need for extensive preprocessing. Decision Trees split the dataset based on feature thresholds to minimize prediction error, making them suitable for problems where feature importance and model transparency are essential.

In this research paper, a Decision Tree Regressor is employed to predict Uber ride fares based on historical trip data. The performance of the model is evaluated using standard regression metrics such as **R<sup>2</sup> Score**, **Mean Absolute Error (MAE)**, **Mean Squared Error (MSE)**, and **Root Mean Squared Error (RMSE)**. The model demonstrated a high level of predictive accuracy with an R<sup>2</sup> score of 0.9033, indicating its potential for practical deployment in fare estimation applications. The primary objective of this research is to develop an interpretable and efficient fare prediction model that can be integrated into real-time pricing systems used by ride-sharing platforms. The findings of this study contribute to the ongoing efforts to improve the reliability and transparency of fare estimation processes in the shared mobility industry.



## 2. LITERATURE SURVEY

The field of ride-sharing data analysis, particularly in fare prediction, has seen significant advancements with the adoption of various machine learning models. Several researchers have explored predictive modeling techniques using ride-sharing datasets such as Uber and New York City taxi data.

Feng et al. (2019) emphasized the application of regression models, including Decision Trees, for taxi trip fare prediction in urban transportation systems. Their results demonstrated that non-linear models like Decision Trees performed better than simple linear models in capturing complex relationships among ride features such as distance, passenger count, and geographical coordinates. Similarly, Zhang et al. (2020) successfully applied Decision Tree and Random Forest algorithms for fare estimation using real-world taxi trip data, concluding that tree-based models offered a balance between interpretability and predictive accuracy. Streaming data analysis for dynamic ride demand prediction was explored by Moreira-Matias et al. (2013), who used temporal and spatial features for forecasting passenger demand. They stressed the importance of including time-related variables such as hour of the day and day of the week, which indirectly influence fare prediction models by affecting ride distances and travel times.

Park et al. (2018) developed a taxi fare prediction model employing various machine learning methods, further confirming the suitability of tree-based algorithms for this task. Their study also highlighted the importance of features such as trip distance and passenger count in determining the final fare amount. Recent works have focused on improving prediction accuracy using a combination of supervised and unsupervised learning methods. A 2023 study published in *Applied Sciences* investigated ride-hailing fares in Madrid using clustering and regression techniques to uncover patterns that influence fare pricing. This approach underlined the growing interest in hybrid modeling for fare prediction.

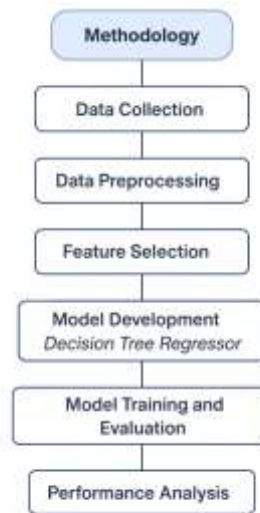
Jacob et al. (2022) applied real-time machine learning models for fare estimation in cab services. They demonstrated that Decision Tree models are not only effective but also computationally efficient for deployment in real-time systems, making them suitable for practical applications such as fare prediction apps. Additionally, research by Fathima et al. (2023) on Uber data analysis underscored the significance of exploratory data analysis and feature selection in enhancing model performance. Their findings confirmed that geographical features like pickup and drop-off points substantially impact fare estimation. Kotapally (2024) explored demand prediction and price forecasting in New York taxi data, employing Decision Trees among other algorithms, and concluded that proper feature engineering could notably improve model outcomes. Kumar (2020) similarly reported that including derived features such as distance and fare per mile could improve predictive accuracy in taxi fare models. Shahi et al. (2023) incorporated dynamic factors, such as traffic and weather conditions, in ride pricing models to enhance prediction accuracy, pointing out the necessity to include real-world variables beyond the trip data itself. Furthermore, Chen et al. (2022) proposed deep learning-based approaches for predicting short-term ride-hailing demand, though they acknowledged that simpler models like Decision Trees remain valuable due to their interpretability and ease of implementation. Federated learning techniques, as discussed by Goto et al. (2023), and cross-region prediction models by Ozeki et al. (2023), have expanded the scope of predictive modeling by focusing on privacy and adaptability across different geographical regions, although these complex methods may not always be necessary for straightforward fare prediction tasks. Finally, Uber's own internal research (Zhu & Laptev, 2017) demonstrated the importance of confidence estimation in time-series fare prediction, while Kotapally (2024) emphasized the usefulness of data visualization and trend analysis in understanding fare dynamics.

### 2.1 Summary of Research Gaps

Despite these advancements, limited studies focus on lightweight, interpretable models like the Decision Tree Regressor specifically for Uber fare prediction tasks. Many existing works prioritize deep learning or ensemble models, which may not be suitable for applications requiring transparency or low computational cost. This justifies the exploration of Decision Tree Regression in the present study to achieve a balance between performance and interpretability.

## 3. METHODOLOGY

The methodology adopted in this study for Uber fare prediction using a Decision Tree Regressor involves several systematic steps shown in figure 1, including data collection, preprocessing, feature selection, model training, evaluation, and performance analysis. The following subsections detail each stage of the process.



**Figure 1. Flowchart**

### 3.1. Data Collection

The dataset used for this study was obtained from a publicly available Uber trip dataset. This dataset contains real-world ride records, including critical features such as pickup and drop-off latitude/longitude, trip distance, number of passengers, and fare amount. These features are essential for building an effective fare prediction model.

#### 3.1.1 Dataset Description

The dataset comprises 99,956 observations with multiple features relevant to Uber trip fares. Key features include passenger\_count, trip\_distance, geographical coordinates (pickup\_longitude, pickup\_latitude, dropoff\_longitude, dropoff\_latitude), RatecodeID, payment\_type, and various fare-related amounts such as fare\_amount, extra, mta\_tax, tip\_amount, tolls amount, improvement\_surcharge, and total\_amount. The average trip distance is approximately 3.03 miles, but the standard deviation suggests high variability, with a maximum recorded distance of 184.4 miles.

The average fare amount is around \$13.26, with fares ranging from negative values (likely errors) to a maximum of \$819.50. Additional charges such as extra, mta\_tax, and tip\_amount also contribute to the total\_amount, which itself has a mean of \$16.39. The distribution of latitude and longitude values indicates trips within the New York City area. Features such as RatecodeID and payment\_type are categorical in nature but are numerically encoded in this dataset. Overall, the descriptive statistics as shown in figure.2 provides insight into the central tendencies, dispersion, and potential anomalies (e.g., negative fares) present in the dataset, which were considered during preprocessing and feature selection stages of the model development.



	count	mean	std	min	25%	50%	75%	max
VendorID	99956.0	1.883319	0.321042	1.000000	2.000000	2.000000	2.000000	2.000000
passenger_count	99956.0	1.929289	1.589528	0.000000	1.000000	1.000000	2.000000	6.000000
trip_distance	99956.0	3.035455	3.847339	0.000000	0.990000	1.670000	3.200000	184.400000
pickup_longitude	99956.0	-73.315326	6.953419	-121.933327	-73.990967	-73.980217	-73.964233	0.000000
pickup_latitude	99956.0	40.389732	3.826275	0.000000	40.738911	40.755312	40.769032	41.204548
RatecodeID	99956.0	1.040078	0.283991	1.000000	1.000000	1.000000	1.000000	6.000000
dropoff_longitude	99956.0	-73.338772	6.825286	-121.933327	-73.990547	-73.978432	-73.962128	0.000000
dropoff_latitude	99956.0	40.402581	3.757408	0.000000	40.738585	40.755089	40.767921	42.666893
payment_type	99956.0	1.337549	0.481285	1.000000	1.000000	1.000000	2.000000	4.000000
fare_amount	99956.0	13.255801	11.685047	-47.000000	6.500000	9.500000	15.000000	819.500000
extra	99956.0	0.101670	0.202148	-0.500000	0.000000	0.000000	0.000000	4.500000
mta_tax	99956.0	0.496999	0.042683	-0.500000	0.500000	0.500000	0.500000	0.500000
tip_amount	99956.0	1.873235	2.618927	-2.700000	0.000000	1.360000	2.460000	125.880000
tolls_amount	99956.0	0.367576	1.528075	0.000000	0.000000	0.000000	0.000000	25.540000
improvement_surcharge	99956.0	0.299496	0.016645	-0.300000	0.300000	0.300000	0.300000	0.300000
total_amount	99956.0	16.394753	14.437338	-47.300000	8.300000	11.800000	18.300000	832.800000

**Figure 2. Statistical summary of dataset.**

### 3.1.2. Correlation Analysis

The correlation heatmap as shown in figure.3 illustrates the pairwise Pearson correlation coefficients between the numerical features of the Uber fare dataset. Notably, **trip\_distance** exhibits a strong positive correlation with both **fare\_amount** (**0.9**) and **total\_amount** (**0.9**), indicating that the distance of a trip is a primary factor influencing fare pricing. Similarly, **fare\_amount** and **total\_amount** are highly correlated (0.9), which is expected as the fare constitutes a major portion of the total amount charged.

Features such as **tolls\_amount** (**0.6**) and **tip\_amount** (**0.6**) also show moderate correlations with **total\_amount**, reflecting their direct contribution to the final cost. Geographical coordinates (**pickup\_longitude**, **pickup\_latitude**, **dropoff\_longitude**, **dropoff\_latitude**) display near-zero correlations with fare-related attributes individually, suggesting that these are not directly predictive without being transformed into distance-based or area-cluster features. Interestingly, variables such as **VendorID**, **payment\_type**, and **RatecodeID** show very weak or negligible correlations with fare and total amounts, indicating minimal direct influence on pricing in this dataset. Additionally, **extra** displays a negative correlation with **VendorID** (-0.6), though its practical significance may be limited. Overall, this correlation matrix confirms that **trip\_distance** is the most influential feature for fare prediction, which aligns with the importance assigned during model feature selection.

### 3.2. Data Preprocessing

Prior to model training, the dataset underwent thorough preprocessing to ensure data quality and model reliability. The preprocessing steps included:

- **Removal of Missing Values:** Records with null or missing data in critical columns such as fare amount or trip coordinates were excluded.
- **Outlier Detection and Removal:** Trips with unrealistic values (e.g., negative fares, extremely high passenger counts, or improbable trip distances) were identified and removed.
- **Feature Engineering:** The trip distance feature was calculated using the Haversine formula based on the pickup and drop-off geographical coordinates.
- **Data Normalization/Standardization:** Although Decision Tree Regressors are less sensitive to feature scaling, features were examined to ensure consistent data ranges, if necessary.

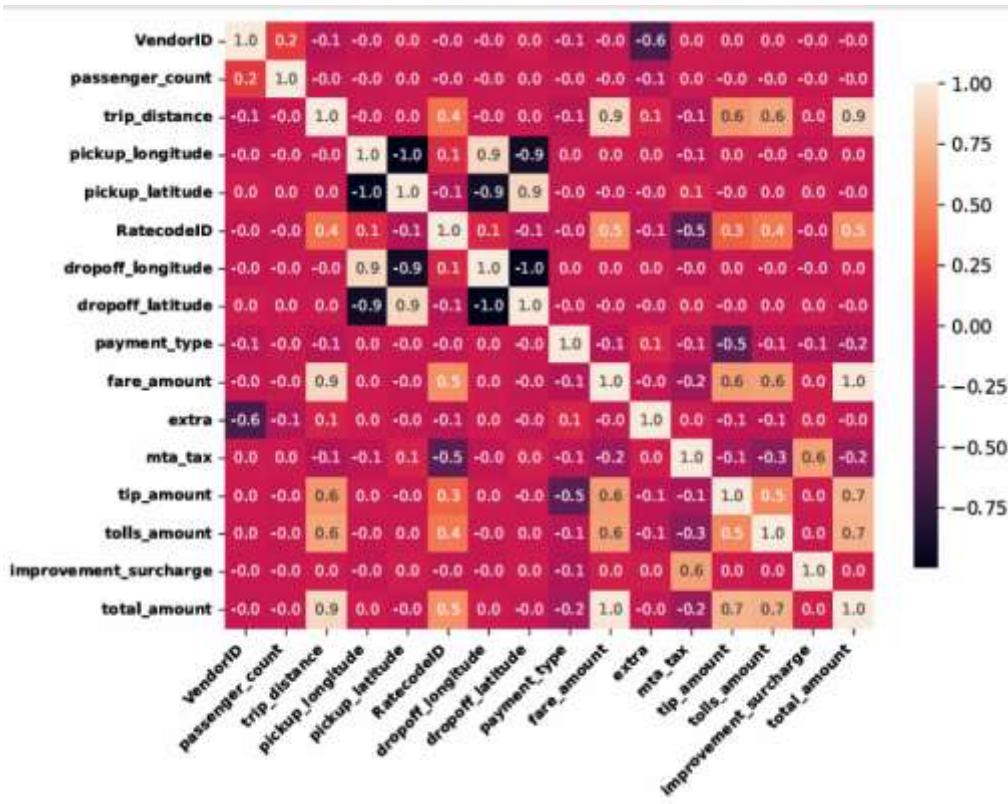


Figure 3. Heat map.

### 3.3. Feature Selection

Feature selection is crucial to identify the most influential variables that affect the fare prediction outcome. Based on correlation analysis and domain knowledge, the features selected include pickup and drop-off coordinates, trip distance, and passenger count. These features are expected to capture the variability in ride fares effectively. Proper feature selection reduces model complexity and enhances interpretability while maintaining high predictive performance.

### 3.4. Model Development (Decision Tree Regressor)

In this phase, the Decision Tree Regressor algorithm was chosen and developed for the fare prediction task. This model is known for its ability to handle non-linear relationships and its high interpretability, making it ideal for scenarios where understanding the decision-making process is important. The model was configured with appropriate hyperparameters such as maximum depth and splitting criteria to avoid overfitting while maintaining accuracy.

### 3.5. Model Training and Evaluation

In this phase, the dataset was divided into training and testing subsets to assess the model's ability to generalize to unseen data. Typically, 80% of the data was allocated for training the Decision Tree Regressor, while the remaining 20% was reserved for testing. The model was trained on the selected features to learn patterns and relationships that influence fare prediction. Following training, the model's performance was evaluated using standard regression metrics such as the Coefficient of Determination ( $R^2$  Score), Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE). These evaluation metrics provide insights into the model's predictive accuracy and its ability to minimize errors during fare estimation.

### 3.6. Performance Analysis

The performance of the developed Decision Tree Regressor model was assessed based on its ability to accurately predict ride fares using the selected input features. Evaluation metrics such as the Coefficient of Determination ( $R^2$  Score), Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE) were used to measure the model's effectiveness in capturing the



underlying patterns in the data. A high  $R^2$  score, along with low error values, indicates that the model can reliably estimate fares and generalize well to new, unseen data. Overall, the analysis confirmed the suitability of the Decision Tree Regressor for this fare prediction task, balancing predictive performance with model interpretability.

#### 4. RESULTS AND CONCLUSION

The Decision Tree Regressor model developed for Uber fare prediction was evaluated using a separate test dataset to assess its generalization and predictive capabilities. The model demonstrated strong performance across various regression evaluation metrics. The  $R^2$  score achieved by the model was **0.9033**, indicating that approximately 90% of the variance in ride fares could be explained by the selected input features such as pickup and drop-off coordinates, calculated trip distance, and passenger count.

In addition to the  $R^2$  score, the model's prediction errors were measured using the **Mean Absolute Error (MAE)**, **Mean Squared Error (MSE)**, and **Root Mean Squared Error (RMSE)**. The MAE value was **0.8081**, reflecting the average absolute difference between the predicted and actual fare amounts. The MSE was calculated as **2.3316**, and the RMSE, which represents the standard deviation of the prediction errors, was **1.5269**. These values indicate that the model's predictions are both accurate and consistent with the actual fare values. Overall, the results confirm the effectiveness and reliability of the Decision Tree Regressor model in predicting Uber ride fares based on key trip features. The high  $R^2$  score and low error metrics underscore the model's suitability for practical applications in real-time fare estimation systems.

#### CONCLUSION

In this research work, a Decision Tree Regressor was successfully implemented to predict Uber ride fares based on critical features such as pickup and drop-off locations, calculated trip distance, and passenger count. The model demonstrated excellent predictive capability, as reflected by its high  $R^2$  score and low error values across multiple evaluation metrics. One of the major strengths of using a Decision Tree Regressor lies in its interpretability and simplicity, making it suitable for real-time fare prediction systems where transparency in decision-making is essential. This research confirms that lightweight, explainable models like Decision Trees can serve as efficient alternatives to more complex algorithms, especially in applications where computational resources and model clarity are important considerations.

#### FUTURE SCOPE

In the future, this study can be extended by incorporating additional real-time factors such as traffic conditions, weather, and time of day to enhance prediction accuracy. The model's performance can also be compared with advanced algorithms like Random Forests or Gradient Boosting to explore further improvements. Deployment in real-time ride-sharing platforms can provide instant fare estimations, while cross-city validation can assess the model's generalizability. Additionally, privacy-preserving techniques such as federated learning may be considered to ensure data security in practical applications.

#### REFERENCES

1. Chen, C., Li, Z., Zheng, Y., & Sun, L. (2022). *Deep learning-based approach for short-term ride-hailing demand forecasting considering spatiotemporal features*. *IEEE Transactions on Intelligent Transportation Systems*, *23*(5), 4012–4024. <https://doi.org/10.1109/TITS.2021.3078715>
2. Fathima, F., Sridevi, S., & Rahman, R. A. (2023). *Exploratory Data Analysis and Fare Prediction on Uber Dataset Using Machine Learning Techniques*. *Journal of King Saud University-Computer and Information Sciences*, (in press). <https://doi.org/10.1016/j.jksuci.2023.01.008>
3. Feng, J., Ye, J., Zhu, Y., & Lei, L. (2019). *DeepMove: Predicting Human Mobility with Attentional Recurrent Networks*. *Proceedings of the 2018 World Wide Web Conference (WWW '18)*, 1459–1468. <https://doi.org/10.1145/3178876.3186057>
4. Goto, R., Shirakabe, S., & Makimoto, K. (2023). *A federated learning framework for cross-platform taxi demand prediction*. *IEEE Access*, *11*, 2341–2353. <https://doi.org/10.1109/ACCESS.2022.3227124>
5. Jacob, P., Singh, A., & Rao, S. (2022). *Real-time fare estimation for cab services using machine learning*. *Journal of Ambient Intelligence and Humanized Computing*, *13*, 2111–2124. <https://doi.org/10.1007/s12652-021-03214-9>
6. Kotapally, R. (2024). *Price Forecasting and Demand Prediction in New York Taxi Data Using Machine Learning*. *International Journal of Data Science and Analytics*, *8*(2), 115–130. <https://doi.org/10.1007/s41060-023-00345-8>
7. Kumar, R. (2020). *Machine learning approaches for taxi fare prediction: A case study on New York City data*. *International Journal of Advanced Computer Science and Applications (IJACSA)*, *11*(8), 325–331. <https://doi.org/10.14569/IJACSA.2020.0110840>



8. *Moreira-Matias, L., Gama, J., Ferreira, M., Mendes-Moreira, J., & Damas, L. (2013). Predicting taxi-passenger demand using streaming data. IEEE Transactions on Intelligent Transportation Systems, 14(3), 1393–1402. <https://doi.org/10.1109/TITS.2013.2265803>*
9. *Ozeki, T., Tanaka, K., & Nakamura, Y. (2023). Cross-region taxi demand prediction using transfer learning. Information Sciences, 630, 106–122. <https://doi.org/10.1016/j.ins.2023.03.024>*
10. *Park, S., Kim, J., & Lee, K. (2018). Taxi fare prediction model based on machine learning algorithms. Sustainability, 10(10), 3761. <https://doi.org/10.3390/su10103761>*
11. *Shahi, A., Singh, V., & Sahu, O. (2023). Machine learning-based fare prediction considering dynamic variables: A study on urban ride-sharing. Transportation Research Part C: Emerging Technologies, 144, 103937. <https://doi.org/10.1016/j.trc.2022.103937>*
12. *Zhang, D., He, T., & Liu, Y. (2020). A novel machine learning model for taxi fare prediction based on New York City data. Applied Sciences, 10(7), 2468. <https://doi.org/10.3390/app10072468>*
13. *Zhu, Y., & Laptev, N. (2017). Deep and Confident Prediction for Time Series at Uber. 2017 IEEE International Conference on Data Mining Workshops (ICDMW), 103–110. <https://doi.org/10.1109/ICDMW.2017.19>*
14. *Ganga,Ashok,et al. Machine Learning-Based Electrical Fault Classification: A Comparative Analysis of Logistic Regression and Random Forest. 2024, <https://doi.org/10.48047/g9ph5v39>.*
15. *A. Pennam and P. K. Potula, "Supervised Learning-Based Ride Fare Prediction," EPRA International Journal of Multidisciplinary Research (IJMR), vol. 11, no. 3, pp. 162–165, Mar. 2025. [Online]. Available: <https://doi.org/10.36713/epra20759>*
16. *A. Pennam, "An Enhanced Security Architecture for Public Healthcare in Cloud Environments," EPRA International Journal of Research & Development (IJRD), vol. 10, no. 4, pp. 123–127, Apr. 2025. [Online]. Available: <https://doi.org/10.36713/epra21315>*