

# Uber Trips Analysis: Machine Learning Clustering for Analysing Travel Patterns

<sup>1</sup>B Manohar Prasad, <sup>2</sup>T Subhashini, <sup>3</sup>N Prasanna, <sup>4</sup>O A N Dev Varma, <sup>5</sup>V L Priyanka

<sup>1</sup>Assistant Professor, Srinivasa Institute of Engineering and Technology

<sup>2,3,4,5</sup>Student Scholar, Department of Computer Science and Engineering, Srinivasa Institute of Engineering and Technology

doi.org/10.64643/IJIRTV12I10-194955-459

**Abstract** The rapid growth of ride-hailing services has generated massive amounts of trip data, providing valuable insights into urban mobility and travel behaviour. Analysing such large-scale datasets manually is inefficient and prone to error. This research presents a machine learning-based approach to analyse Uber trip data and identify significant travel patterns using clustering techniques. The study utilizes a publicly available Uber trips dataset containing spatial and temporal attributes such as pickup and drop-off locations, trip time, distance, and fare amount. Data preprocessing techniques including handling missing values, normalization, and feature selection were applied to enhance data quality. The K-Means clustering algorithm was implemented to group trips based on similarity in spatial and temporal features. The performance of the clustering model was evaluated using silhouette score and visual analysis. Results indicate that distinct clusters representing high-demand travel zones and peak travel periods can be effectively identified. The findings demonstrate that machine learning clustering can provide meaningful insights into passenger movement patterns and support better transportation planning, demand forecasting, and resource allocation.

**Keywords:** Travel Pattern Analysis, Machine Learning, Uber Trips, K-Means, Clustering, Transportation Analytics.

## I. INTRODUCTION

Urban transportation systems have undergone significant transformation with the introduction of ride-sharing platforms such as Uber. These platforms generate large volumes of trip data that include pickup and drop-off locations, timestamps, distance travelled, and fare information. This data reflects real-world

passenger movement and demand across urban regions.

Understanding travel patterns from this data is essential for traffic management, driver allocation, and city planning. Traditional statistical methods are often insufficient for analysing such complex and high-dimensional datasets. Moreover, manual analysis becomes impractical due to the scale of data involved.

Recent advancements in machine learning provide powerful tools for extracting hidden patterns from large datasets. In particular, unsupervised learning techniques such as clustering can group similar data points without requiring labelled outputs. By applying clustering techniques to Uber trip data, it is possible to identify high-demand zones, peak travel periods, and frequently travelled routes. This study focuses on using K-Means clustering to analyse Uber trip data and uncover meaningful travel patterns.

## II. LITERATURE SURVEY

Urban transportation analysis has been widely studied in the fields of data mining and intelligent transportation systems. Early research focused on statistical analysis of travel time, trip frequency, and route distribution. These approaches provided descriptive summaries of transportation behaviour but lacked predictive and pattern discovery capabilities.

With the advancement of machine learning, researchers began applying clustering techniques to transportation data. K-Means clustering has been widely used to group trips based on distance, duration, and location coordinates. It is simple and

computationally efficient, making it suitable for large datasets. However, K-Means requires the number of clusters to be predefined and is sensitive to outliers.

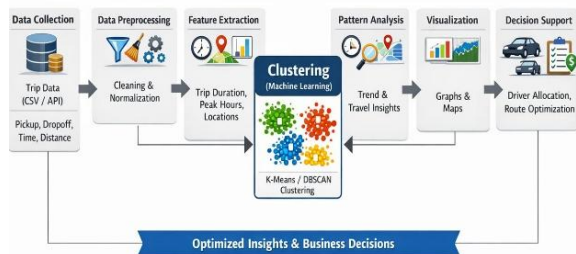
Density-based clustering methods such as DBSCAN have also been used for spatial analysis of pickup and drop-off locations. DBSCAN is effective in detecting high-density regions and identifying noise points that do not belong to any cluster. This makes it useful for identifying popular travel zones and unusual trips.

Recent studies have integrated clustering with visualization techniques such as heatmaps and geographic mapping to analyse urban mobility. Some research has also combined clustering with time-series analysis to study peak and off-peak travel behaviour. However, many existing studies focus on a single clustering algorithm and do not provide a structured framework that includes preprocessing, feature engineering, and systematic evaluation.

This research addresses these limitations by integrating multiple clustering techniques with proper data preprocessing and evaluation metrics to improve the discovery of meaningful travel patterns from Uber trip data.

### III. SYSTEM ARCHITECTURE

The proposed Uber Trip Analysis System follows a layered machine learning architecture to ensure modularity, scalability, and interpretability. The system is organized into five main layers: Data Collection Layer, Data Preprocessing Layer, Feature Engineering Layer, Machine Learning Layer, and Pattern Analysis & Evaluation Layer.



#### Data Collection Layer

This layer is responsible for acquiring historical Uber trip data. The dataset contains attributes such as pickup

latitude, pickup longitude, drop-off latitude, drop-off longitude, trip distance, trip duration, and pickup time. These features serve as the foundation for analysing travel behaviour.

#### Data Preprocessing Layer

The preprocessing layer prepares raw data for clustering. It includes handling missing values, removing invalid coordinates, filtering out extreme outliers, and converting timestamps into useful components such as hour and day. Feature normalization using Standard Scaler is applied to ensure that all attributes contribute equally to distance calculations.

#### Feature Engineering Layer

This layer derives additional features from existing data, such as peak-hour indicator, weekday or weekend classification, and estimated travel speed. These features improve clustering performance by capturing both spatial and temporal aspects of trips.

#### Machine Learning Layer

This layer implements unsupervised learning algorithms. K-Means clustering is used to group trips based on similarity in distance and time, while DBSCAN is applied to detect dense pickup and drop-off regions. Hyperparameters are tuned to obtain optimal clustering performance.

#### Pattern Analysis and Evaluation Layer

The final layer evaluates the clustering results using silhouette score and inertia. Visual tools such as scatter plots and heatmaps are used to interpret cluster distributions and travel patterns.

### IV. METHODOLOGY

This study adopts a machine learning-based analytical approach to examine Uber trip data and discover hidden travel patterns through clustering. A publicly available Uber trips dataset was used, consisting of spatial and temporal attributes such as pickup and drop-off coordinates, trip time, travel distance, and fare amount. These attributes collectively describe passenger movement behaviour across different regions and time periods.

### 1. Data Collection

The dataset used in this study consists of Uber trip records obtained from publicly available open data sources. Each record represents a single trip and contains attributes such as pickup date and time, pickup latitude and longitude, trip distance, and trip duration. These attributes capture both the spatial and temporal characteristics of passenger movement. The dataset includes trips collected over different time periods and geographic locations, enabling the analysis of variations in travel behaviour across time and space. This data serves as the foundation for discovering hidden mobility patterns using machine learning techniques.

### 2. Data Pre-processing

Raw trip data often contains missing values, noise, and inconsistencies that can negatively affect clustering performance. Therefore, data preprocessing is a crucial step in this study. Records with missing or inconsistent values are removed to improve data quality. The pickup date and time attribute is converted into numerical features such as hour of the day and day category (weekday or weekend) in order to capture temporal variations in travel demand. Continuous variables such as trip distance and trip duration are normalized to ensure uniform scaling and to prevent any single feature from dominating the clustering process. Outliers, which represent abnormal or extreme trip values, are identified and removed to improve clustering accuracy and stability.

### 3. Feature Selection

Feature selection is performed to retain only the most relevant attributes for clustering. In this study, key features such as trip distance, pickup hour, and geographical coordinates (latitude and longitude) are selected because they directly influence travel behaviour and demand patterns. These features provide meaningful information about when and where trips occur, as well as how far passengers travel. Redundant or less informative attributes are excluded to reduce dimensionality and computational complexity while improving the interpretability of the clustering results.

### 4. Clustering Algorithm

To group similar trips based on their characteristics, the K-Means clustering algorithm is applied to the pre-processed dataset. K-Means is chosen due to its

simplicity, efficiency, and suitability for large-scale datasets. The algorithm works by partitioning the data into a predefined number of clusters such that trips within the same cluster are more similar to each other than to those in other clusters. The optimal number of clusters is determined using the Elbow Method, which evaluates the within-cluster sum of squares for different cluster values and identifies the point where further increase in clusters does not significantly reduce variance. This ensures that the resulting clusters are compact and well separated.

### 5. Pattern Analysis

After clustering and initial result interpretation, pattern analysis is performed to extract meaningful insights from the clustered Uber trip data. Each cluster is examined to identify common characteristics in terms of spatial and temporal behaviour. Spatial patterns are analysed by observing the concentration of pickup locations within each cluster, which helps in identifying frequently used travel zones and high-demand regions. Temporal patterns are studied by analysing the distribution of trips across different hours of the day and days of the week, enabling the identification of peak and off-peak travel periods.

### 6. Result Analysis

The results obtained from the clustering process demonstrate the effectiveness of the proposed approach in identifying meaningful travel patterns from Uber trip data. After applying the K-Means clustering algorithm, the dataset was divided into distinct clusters based on spatial and temporal similarities among trips. Each cluster represents a group of trips sharing common characteristics in terms of pickup location, travel time, and trip distance.

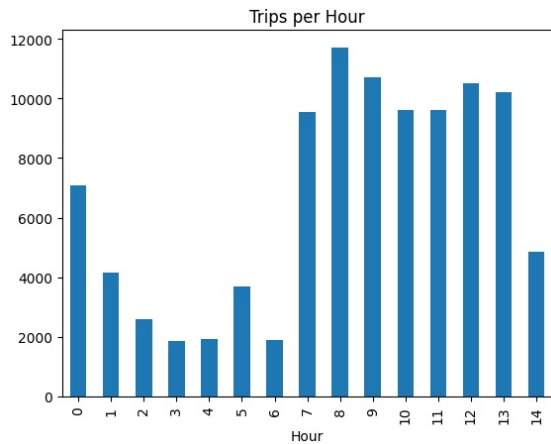
The spatial distribution of clusters reveals that certain regions exhibit a significantly higher concentration of trips compared to others. These high-density clusters correspond to major commercial areas, transportation hubs, and densely populated regions, indicating zones of high travel demand. In contrast, clusters with lower trip density represent suburban or less frequently traveled areas. This spatial segregation confirms that travel activity is not uniformly distributed across the city but is instead concentrated in specific regions.

## V. RESULTS

The Random Forest Regression model was trained on the pre-processed medical insurance dataset to predict insurance charges based on demographic and health-related features. The model achieved a high  $R^2$  score, indicating that it effectively captures the variance in the target variable and provides accurate predictions.

A. Time-Based Travel Pattern Analysis (Trips per Hour)

The time-based analysis shows the variation in Uber trip demand across different hours of the day. The number of trips is low during late-night and early-morning hours, indicating reduced travel activity. A sharp increase in trip volume is observed during morning hours, with the highest demand occurring during typical commuting time. After the morning peak, demand remains relatively high throughout the daytime and increases again during evening hours, corresponding to return journeys from workplaces. This pattern indicates that Uber usage is strongly influenced by daily routines and peak commuting periods.



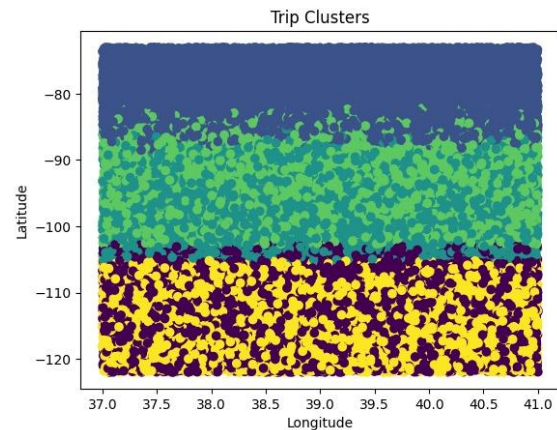
B. Cluster Analysis Based on Trip Distance and Fare Amount

The distance-fare clustering result groups Uber trips into distinct clusters based on similarity in trip distance and fare amount. Most trips belong to clusters representing short to medium distances with low to moderate fares, which correspond to common intra-city travel. Another cluster represents trips with higher distances and fares, indicating longer journeys across different regions. A few outliers with very large distances and high fares represent rare long-distance trips. This clustering highlights different types of travel behaviour based on cost and distance.



C. Spatial Cluster Analysis (Geographical Distribution of Trips)

The spatial clustering result shows that Uber trips are grouped into distinct geographical regions based on pickup location coordinates. High-density clusters indicate areas with frequent pickups and drop-offs, such as commercial zones, residential hubs, and transportation centres. Other clusters with lower density represent regions with less travel activity. The clear separation of clusters confirms that travel demand is concentrated in specific locations rather than being uniformly distributed across the city. This spatial analysis helps identify high-demand travel zones and understand regional mobility patterns.



D. Overall Interpretation of Results

The combined analysis of time-based patterns, distance-fare clustering, and spatial clustering provides a comprehensive understanding of Uber travel behaviour. The results show that Uber trips are influenced by time of day, trip characteristics, and geographical location. High-demand periods

correspond to commuting hours, while spatial clusters reveal concentrated travel activity in specific regions. These findings demonstrate that machine learning clustering techniques are effective in extracting meaningful travel patterns from Uber trip data and can support transportation planning and service optimization.

## VI. DISCUSSION

The results obtained from the clustering-based analysis highlight clear temporal, spatial, and trip-related patterns in Uber travel behaviour. The time-based analysis shows that trip demand varies significantly throughout the day, with higher activity during morning and evening hours. This trend reflects typical commuting behaviour and indicates that Uber usage is strongly influenced by daily work schedules and urban routines. Such temporal patterns are useful for understanding peak demand periods and can support better planning of ride-hailing services.

The distance-fare clustering results reveal that most trips belong to short to medium distance categories with relatively low to moderate fares, representing regular intra-city travel. A smaller number of trips fall into clusters with higher distances and fares, indicating longer journeys between major city regions. This variation suggests that Uber serves both short local travel needs and longer-distance transportation requirements. The presence of a few outliers further indicates that although most trips follow common patterns, occasional special or long-distance trips also occur.

Spatial clustering demonstrates that travel demand is not evenly distributed across geographical regions. Instead, trips are concentrated in specific zones, which likely correspond to commercial districts, residential hubs, and transportation centres. These high-density clusters represent important travel hotspots and reflect areas with consistent passenger demand. Regions with lower cluster density indicate comparatively less travel activity, which may be associated with suburban or less populated areas. This uneven distribution highlights the role of location in influencing ride-hailing demand. The combined findings show that machine learning clustering is effective in extracting meaningful patterns from large-scale Uber trip data without the need for labelled outputs. By grouping trips based on similarity

in time, distance, and location, the approach provides interpretable insights into urban mobility behaviour. These insights can be applied to improve driver allocation strategies, reduce passenger waiting time, and enhance service efficiency. Urban planners can also use such information to identify congestion-prone zones and plan better transportation infrastructure.

However, the study has certain limitations. The clustering model does not consider external factors such as traffic conditions, weather, or special events, which can also influence travel behaviour. In addition, K-Means clustering requires prior selection of the number of clusters and is sensitive to outliers. These factors may affect the accuracy and generalization of the results. Despite these limitations, the analysis demonstrates the potential of unsupervised machine learning techniques in understanding complex travel patterns from real-world transportation data.

Overall, the discussion confirms that clustering-based analysis provides valuable insights into travel demand and passenger movement trends. The observed temporal and spatial patterns validate the suitability of machine learning approaches for transportation analytics and support their application in intelligent urban mobility systems.

## VII. CONCLUSION

This project presented a machine learning-based clustering approach to analyse Uber trip data and identify meaningful travel patterns. By applying clustering techniques to spatial, temporal, and trip-related attributes, the study successfully grouped trips with similar characteristics and revealed important trends in urban mobility behaviour.

The results showed that Uber trip demand varies significantly with time, with higher activity during morning and evening commuting hours. Clustering based on trip distance and fare highlighted different categories of travel behaviour, ranging from short local trips to longer inter-regional journeys. Spatial clustering further demonstrated that travel demand is concentrated in specific geographical regions, indicating high-demand zones such as commercial and residential hubs.

These findings confirm that unsupervised machine learning methods, particularly clustering, are effective in extracting useful information from large-scale

transportation datasets. The proposed approach provides a data-driven way to understand passenger movement patterns and can support better decision-making for ride-hailing services, traffic management, and urban transportation planning.

Overall, this study demonstrates the potential of machine learning clustering techniques in transportation analytics and highlights their usefulness in identifying travel hotspots, peak travel periods, and common movement trends. The outcomes of this work can serve as a foundation for future research and the development of intelligent transportation systems based on real-world trip data.

#### REFERENCE

- [1] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed. San Francisco, CA, USA: Morgan Kaufmann, 2012.
- [2] S. Raschka and V. Mirjalili, *Python Machine Learning*, 3rd ed. Birmingham, U.K.: Packt Publishing, 2019.
- [3] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed. New York, NY, USA: Springer, 2009.
- [4] A. K. Jain, "Data clustering: A review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264–323, Sep. 1999.
- [5] Kaggle, "Uber trips dataset," [Online]. Available: <https://www.kaggle.com>
- [6] Uber Technologies Inc., "Uber Movement: Travel and mobility data," [Online]. Available: <https://movement.uber.com>
- [7] P.-N. Tan, M. Steinbach, and V. Kumar, *Introduction to Data Mining*. Boston, MA, USA: Pearson Education, 2016.
- [8] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.