



IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 13 Issue: V Month of publication: May 2025

DOI: https://doi.org/10.22214/ijraset.2025.71576

www.ijraset.com

Call: 🛇 08813907089 🕴 E-mail ID: ijraset@gmail.com



# **Exploratory Data Analytics on Uber Transportation Patterns**

Shanthini D

Department of MCA, Paavai Engineering College (Autonomous)

Abstract: This paper presents a comprehensive analysis of Uber trip data in New York City using data science techniques to improve ride-hailing efficiency. The project utilizes historical trip data from the NYC High Volume For-Hire Vehicle (HVFHV) dataset, incorporating clustering (K-Means) and prediction (Random Forest Regression) to forecast demand patterns. It also integrates Power BI for real-time visualization and insights. The aim is to optimize driver allocation, reduce passenger wait times, and enhance urban mobility.

Keywords: Uber, Data Analysis, K-Means Clustering, Random Forest, Power BI, Ride-hailing Optimization

# I. INTRODUCTION

Urban mobility faces increasing challenges in matching ride-hailing demand with efficient driver allocation. This study focuses on analyzing Uber trip data to understand spatial and temporal demand patterns, enabling better decision-making through data science and visualization.

# II. METHODOLOGY

#### A. Data Collection and Preprocessing

Trip data in Parquet format was collected from the NYC TLC HVFHV dataset, merged with taxi zone and weather data. The dataset underwent preprocessing to clean nulls, convert types, and extract features.

### B. Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) employed visualization tools like Matplotlib and Seaborn to uncover patterns such as peak travel hours, demand across zones, and seasonal variations.

### C. Clustering with K-Means

K-Means was applied on features like pickup/drop zones and trip metrics to group locations with similar demand profiles.

### D. Prediction using Random Forest Regression

Random Forest predicted trip volume and estimated wait times. Feature engineering and performance validation using  $R^2$  score ensured accuracy (~84%).

#### E. Visualization with Power BI

Interactive dashboards showcased daily trends, zone distribution, and seasonal demand, enabling decision-makers to act on insights.

#### **III. RESULTS AND DISCUSSION**

The clustering identified high-demand pickup zones, while regression modeling yielded a high R<sup>2</sup> score, confirming strong predictive capability. Dashboards revealed consistent peak demand in zones like Midtown and Downtown during weekdays and evenings. Recommendations include real-time adjustments to vehicle distribution and dynamic pricing.

#### **IV. CONCLUSION**

The project successfully applied data analytics and machine learning to Uber trip data for proactive fleet management. By integrating visualization and prediction, it supports smarter, faster, and data-driven urban mobility solutions.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue V May 2025- Available at www.ijraset.com

		TAE	BLE I		
		TEST CAS	SE REPORT		
TEST	TEST CASE	DESCRIPTION	EXPECTED	ACTUAL	STATUS
CASE ID	TITLE		RESULT	RESULT	
TC_01	LOAD	VERIFY	DATA LOADS	AS	PASS
	PARQUET	CORRECT	WITH REQUIRED	EXPECTED	
	TRIP DATA	LOADING OF	COLUMNS		
		TRIP DATA			
TC_02	MERGE	ENSURE	WEATHER	AS	PASS
	WEATHER	WEATHER DATA	ALIGNED WITH	EXPECTED	
	DATA	INTEGRATES	TRIP DATES		
		WITH TRIP			
		RECORDS			
TC_03	APPLY K-	CLUSTER TRIPS	CLUSTERS	AS	PASS
	MEANS	BY LOCATION	FORMED	EXPECTED	
	CLUSTERING	AND TIME	SHOWING		
			DEMAND		
			HOTSPOTS		
TC_04	TRAIN	PREDICT TRIP	HIGH R <sup>2</sup> SCORE	AS	PASS
	RANDOM	TIME AND	AND EXCEL	EXPECTED	
	FOREST	EXPORT	EXPORT		
	MODEL	RESULTS			
TC_05	VISUALIZE	DISPLAY	PROPER DATE	AS	PASS
	TRENDS IN	DAILY/MONTHLY	FORMAT AND	EXPECTED	
	POWER BI	TRENDS	VISUALS		
		CORRECTLY			
TC_06	DEMAND	IDENTIFY HIGH-	VISUAL DISPLAY	AS	PASS
	ZONE	DEMAND ZONES	OF TOP ZONES	EXPECTED	
	CLUSTERING	USING			
		CLUSTERING			



Fig .1 PowerBICharts



International Journal for Research in Applied Science & Engineering Technology (IJRASET)

ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 13 Issue V May 2025- Available at www.ijraset.com

	by Borough		Pick-up Zone Distribu	tion	
D_Borough Cour	nt of DO_LocationID	1	PU_Zone	Count of PULocationID	
anhattan	7002041		JFK Airport	346138	
ooklyn	4757877	1	LaGuardia Airport	307590	
ueens	3540301		East Village	276059	
onx	2093274 18479031		Times Sq/Theatre District	242050	
tai			Crown Heights North	241050	
			Midtown Center	237437	
			TriBeCa/Civic Center	226617	
			Union Sq	216000	
			East Chelsea	215758	
			Bushwick South	206472	
			Midtown South	202965	
			Total	18479031	
IP PICK-UPS D	y Service Zone				
1M					
(35.53%)	(2.55%)	PU_service_zone Boro Zone Vellow Zone Airports	1M (3.4%)	DO_servit Boro Zo P Yellow 2 Cueens Arports	e_zone ne lone

Fig.2ZonebasedVisualization

# REFERENCES

- [1] Poritigadda, L., et al. (2024). Spatial Data Analysis on On-Demand Cab Services Using Spark.
- [2] Golshanrad, P., et al. (2024). Proposing a model for predicting passenger origin-destination in online taxi-hailing systems. Public Transport.
- [3] Kokkiligadda, M. R., et al. (2023). Spatial Data Analysis on On-Demand Cab Services using Spark. IEEE ICIMI.
- [4] Roy, B., & Rout, D. (2021). Predicting Taxi Travel Time Using ML Techniques. Springer.
- [5] Pradhan, R., et al. (2021). Analysing Uber Trips using PySpark. IOP Conf. Ser.: Mater. Sci. Eng.
- [6] Wang, H., et al. (2021). Applying deep learning to taxi demand forecasting: CNN-LSTM model. Transp. Res. Part C.











45.98



IMPACT FACTOR: 7.129







# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 🕓 (24\*7 Support on Whatsapp)