

NYC Taxi Fare Prediction: Machine Learning and Data Analysis

Iachim Cristian
Serbicean Alexandru
Universidade da Madeira
2238224@student.uma.pt
2236924@student.uma.pt

May 20, 2025

Abstract

This report presents a comprehensive study on New York City (NYC) taxi fares using machine learning. We explore predictive models such as k-Nearest Neighbors (kNN), Logistic Regression, Random Forest, Gradient Boosting and Deep Learning. The dataset is preprocessed by removing outliers, encoding categorical variables and extracting additional features. Model performances are evaluated using accuracy, precision, recall and F1-score. The results indicate that ensemble models such as Random Forest and Gradient Boosting achieve 100% accuracy, indicating potential overfitting, while Logistic Regression performs optimally with 99.2% accuracy. Finally, we propose a deployment strategy for real-world applications.

Contents

1	Introduction	2
1.1	Objectives	2
2	Data Analysis and Preprocessing	2
2.1	Dataset Overview	2
2.2	Data Cleaning	2
2.3	Feature Engineering	2
2.4	Time-Based Features	3
2.5	Trip-Related Features	3
2.6	Indicators for Trip Patterns	3
2.7	Fare-Based Feature	4
2.8	Impact of Feature Engineering	4
3	Model Development	4
3.1	k-Nearest Neighbors (kNN)	4
3.2	Supervised Learning Models	4
3.3	Deep Learning Model	5
3.4	Clustering (Unsupervised Learning)	5

4	Results and Discussion	6
4.1	Model Comparison	6
4.2	Evaluation Metrics	6
4.3	Overfitting Analysis	6
5	Deployment Strategy	6
6	Conclusion and Future Work	6
6.1	Key Findings	6
6.2	Future Improvements	7

1 Introduction

Predicting taxi fares is crucial for optimizing pricing models, improving ride-sharing services and enhancing urban mobility. This project leverages machine learning algorithms to classify NYC taxi fares into high and low categories. The dataset, sourced from Kaggle, contains ride details such as trip distance, pickup/dropoff time, passenger count and fare amount.

1.1 Objectives

- Develop machine learning models to predict fare amounts.
- Compare the performance of different algorithms.
- Investigate potential overfitting in ensemble models.
- Deploy a trained model as an API for real-world applications.

2 Data Analysis and Preprocessing

2.1 Dataset Overview

The dataset consists of NYC taxi trips from 2019, obtained from Kaggle. It includes:

- **Numerical Features:** Trip distance, fare amount, trip duration, etc.
- **Categorical Features:** Payment type, vendor ID, store-and-forward flag.
- **Datetime Features:** Pickup and dropoff timestamps.

2.2 Data Cleaning

The preprocessing steps include:

1. Handling missing values.
2. Encoding categorical variables (payment type, vendor ID).
3. Standardizing numerical features using `StandardScaler`.
4. Removing outliers using the **IQR method**.

2.3 Feature Engineering

Feature engineering plays a crucial role in enhancing the predictive power of machine learning models by deriving new, meaningful features from raw data. In this project, we extracted several time-based and trip-related features to improve model performance. These features provide additional insights into trip patterns, customer behavior and fare variations.

2.4 Time-Based Features

Understanding the temporal dynamics of taxi trips is essential, as factors such as rush hours, weekends and nighttime conditions can significantly influence fare prices and trip durations.

- **Pickup Hour** (`pickup_hour`) Extracted from the `tpep_pickup_datetime` column, this feature represents the hour of the day when the trip started. It helps capture daily patterns such as peak and off-peak hours.
- **Pickup Day** (`pickup_day`) Extracted from the `tpep_pickup_datetime` column to indicate the day of the month when the trip took place. This is useful for identifying fare fluctuations over the course of a month.
- **Pickup Weekday** (`pickup_weekday`) Represents the day of the week (0 = Monday, 6 = Sunday). This feature helps determine the impact of weekdays versus weekends on taxi fares.
- **Pickup Month** (`pickup_month`) Extracted from the `tpep_pickup_datetime` column to indicate the month of the trip. This is useful for analyzing seasonal effects on taxi fares.

2.5 Trip-Related Features

To capture the dynamics of taxi trips, several features were engineered to measure trip duration, speed and fare characteristics.

- **Trip Duration** (`trip_duration`) Calculated as the time difference between `tpep_dropoff_datetime` and `tpep_pickup_datetime`, converted into minutes.

$$\text{trip_duration} = \frac{\text{tpep_dropoff_datetime} - \text{tpep_pickup_datetime}}{60} \quad (1)$$

This feature is essential for estimating fare amounts, as longer trips generally cost more.

- **Speed in Miles per Hour** (`speed_mph`) Estimated as the total trip distance divided by the trip duration (converted into hours).

$$\text{speed_mph} = \frac{\text{trip_distance}}{\frac{\text{trip_duration}}{60} + 1e-6} \quad (2)$$

This feature is useful for detecting outliers (e.g., unrealistically high speeds might indicate data anomalies).

2.6 Indicators for Trip Patterns

Binary indicators were created to represent whether a trip occurred during specific conditions such as weekends, rush hours or nighttime.

- **Weekend Indicator** (`is_weekend`) Set to **1** if the trip occurred on **Saturday (5) or Sunday (6)**, otherwise **0**. This helps analyze fare trends during weekends versus weekdays.
- **Rush Hour Indicator** (`is_rush_hour`) Set to **1** if the trip occurred during typical rush hours (**7-9 AM** or **5-7 PM**), otherwise **0**. This feature helps capture fare surges during high-traffic periods.
- **Nighttime Indicator** (`is_night`) Set to **1** if the trip took place between **midnight (00:00) and 5 AM**, otherwise **0**. Nighttime trips often have different fare structures due to lower traffic but increased demand for taxis.

2.7 Fare-Based Feature

To formulate the problem as a classification task, we introduced a binary target variable.

- **High-Fare Classification** (`high_fare`) Set to **1** if the `fare_amount` is greater than the median fare, otherwise **0**. This feature helps convert the regression problem (predicting fare amounts) into a **binary classification task** (high-fare vs. low-fare).

2.8 Impact of Feature Engineering

The introduction of these features significantly improves model performance by providing:

- **More granularity in temporal patterns** (e.g., rush hour vs. non-rush hour trips).
- **Better trip duration and speed insights** (useful for estimating realistic fare amounts).
- **Fare classification for improved decision-making** (e.g., predicting whether a trip is expensive based on known features).

These engineered features contribute to a higher predictive accuracy and help refine the machine learning models used in this project.

3 Model Development

3.1 k-Nearest Neighbors (kNN)

Implemented from scratch, this algorithm predicts fares based on the Euclidean distance between feature vectors. The model achieved 95.1% accuracy.

3.2 Supervised Learning Models

We trained the following models:

- **Logistic Regression:** Achieved 99.2% accuracy.
- **Random Forest:** Achieved 100% accuracy (potential overfitting).
- **Gradient Boosting:** Achieved 100% accuracy (potential overfitting).

3.3 Deep Learning Model

A neural network was trained using TensorFlow with the following architecture:

- Input Layer: 20 features
- Hidden Layers: 64 neurons, 32 neurons (ReLU activation)
- Output Layer: Sigmoid activation (binary classification)

The model achieved an accuracy of 98.6%.

3.4 Clustering (Unsupervised Learning)

Unsupervised learning techniques such as clustering help identify hidden patterns in data without predefined labels. In this project, we applied clustering algorithms to detect fare-based patterns among taxi trips.

- **K-Means Clustering:** The dataset was clustered into two groups using the K-Means algorithm. The clustering revealed distinct fare-based groups but lacked classification accuracy when compared to supervised learning models. The **Silhouette Score**, which measures how well-separated the clusters are, was computed as:

$$\text{Silhouette Score} = 0.9435 \quad (3)$$

A score close to **1** indicates well-separated clusters, suggesting that K-Means effectively identified trip patterns.

- **DBSCAN (Density-Based Clustering):** DBSCAN was applied to detect outliers and identify fare anomalies. Unlike K-Means, DBSCAN does not require the number of clusters to be predefined, making it more flexible for irregular data distributions. The evaluation metrics for DBSCAN clustering are as follows:
 - **Number of clusters found:** 25
 - **Number of noise points:** 8821
 - **Silhouette Score:** -0.8396

The high number of clusters (25) and large noise points (8821) indicate that DBSCAN classified a significant portion of the data as outliers. The negative Silhouette Score (-0.8396) suggests poor separation between clusters, implying that DBSCAN struggled with the datasets structure.

4 Results and Discussion

4.1 Model Comparison

Model	Accuracy
kNN	95.1%
Logistic Regression	99.2%
Random Forest	100.0%
Gradient Boosting	100.0%
Deep Learning	98.6%
K-Means	N/A
DBSCAN	N/A

Table 1: Model Performance Comparison

4.2 Evaluation Metrics

For classification models, we computed:

- **Precision, Recall, F1-Score:** Ensuring correct fare classification.
- **Confusion Matrix:** Analyzing false positives and false negatives.

4.3 Overfitting Analysis

- Random Forest and Gradient Boosting achieved 100% accuracy, suggesting overfitting.
- Cross-validation and hyperparameter tuning are necessary.

5 Deployment Strategy

The best-performing models (Logistic Regression and Deep Learning) can be deployed via a Flask API to allow real-time fare predictions.

6 Conclusion and Future Work

6.1 Key Findings

- Logistic Regression performed best in terms of generalizability (99.2% accuracy).
- Ensemble models (Random Forest, Gradient Boosting) suffered from overfitting.
- Deep Learning provided robust performance (98.6% accuracy) and can be improved with more data.

6.2 Future Improvements

- Implement hyperparameter tuning to improve model robustness.
- Explore additional features (e.g., weather, holidays) for better predictions.
- Deploy the model in a real-world application (e.g., ride-hailing service).

References

- Kaggle NYC Taxi Dataset: <https://www.kaggle.com/dhruvildave/new-york-city-taxi-trip>
- Scikit-Learn Documentation: <https://scikit-learn.org/>