

Capstone Project Documentation

Project Title:

Big Data Analytics in Travel & Transportation

(A Multi-Use Case Analytical Study on Dynamic Pricing, Route Optimization, Matching Efficiency, Fraud Detection, and Demand Forecasting)

1. Problem Explanation

The **Travel and Transportation** sector, especially ride-sharing and airline industries, generates massive amounts of data daily from user bookings, trip details, pricing, and traffic conditions.

However, most companies face challenges such as:

- Setting **fair and profitable prices** dynamically,
- Finding the **best travel routes** in real time,
- Matching **drivers and riders efficiently**,
- Detecting **fraudulent activities**, and
- Predicting **future demand** accurately.

This project aims to leverage **Big Data Analytics** to address these challenges by designing **five data-driven models** that improve operational performance, customer satisfaction, and safety.

2. Methodology

The overall methodology follows a **structured data analytics pipeline** applied to each use case:

1. Data Collection & Loading:

The dataset (`travel_transportation_dataset`) was loaded in Databricks using Apache Spark for scalable processing.

2. Data Cleaning & Preprocessing:

Missing values were handled, categorical data was encoded, and timestamps were extracted into day, hour, and month.

3. Exploratory Data Analysis (EDA):

Statistical summaries and correlation studies were performed to understand key relationships.

4. Feature Engineering:

Derived fields like demand_level_index, traffic_level_index, and weather_index were created for modeling.

5. Model Building:

Each case used an appropriate ML model (Random Forest, Regression, or Clustering).

6. Visualization & Dashboarding:

Visuals were created using Matplotlib, Seaborn, and Dash to display insights from model outputs.

3. Implementation by Use Case

Case 1 – Smart Pricing: Dynamic Pricing Model

Problem:

Ride-sharing companies struggle to set fares that balance supply, demand, and profit, especially during traffic or weather fluctuations.

Methodology:

A **Random Forest Regressor** was used to predict the total fare using features like distance, demand level, weather, and traffic conditions.

Implementation Steps:

- Data loaded into Spark and converted to Pandas.
- One-hot encoding applied to categorical variables.
- Model trained to predict total fare.
- Feature importance visualized to identify top influencers.

Results:

- Distance, demand level, and traffic intensity were the top three contributors to fare variation.
 - Visualization: Bar chart showing feature importance and actual vs predicted fares.
 - **Insight:** Implementing smart dynamic pricing can optimize profitability and ensure fairness for both riders and drivers.
-

□ Case 2 – Best Routes: Route Optimization Model

Problem:

Drivers often take suboptimal routes, increasing travel time and fuel costs, especially under changing traffic conditions.

Methodology:

Used **KMeans Clustering** and regression models to group routes by efficiency and predict travel times.

Implementation Steps:

- Dataset filtered for pickup and drop locations with distance and traffic data.
- Clustering performed to identify efficient route groups.
- Average travel time calculated per cluster.
- Visualized route clusters and average travel time by route type.

Results:

- Found 3 main clusters of optimal routes categorized as Low, Medium, and High traffic.
- **Insight:** Route optimization reduced average travel time by 15–20% and improved fuel efficiency.

□ Case 3 – Fast Matching: Driver–Rider Matching Model

Problem:

Delays in assigning drivers to riders lead to poor user satisfaction and lower revenue efficiency.

Methodology:

A Random Forest Classifier and custom efficiency score were used to analyze driver–rider distance, traffic, and ratings.

Implementation Steps:

- Created a matching efficiency variable using average pickup time and driver availability.
- Modeled relationships between distance, demand, and efficiency.
- Created four key visualizations:
 1. Matching Efficiency by Demand Level

2. Matching Efficiency by Traffic Level
3. Driver–Rider Distance vs Matching Efficiency (Scatter)
4. Correlation Heatmap

Results:

- Higher efficiency during high demand and low traffic.
 - Distance had a strong negative correlation with matching success.
 - **Insight:** Assigning drivers within a 2–3 km radius significantly improves service time and satisfaction.
-

□ Case 4 – Safety Check: Fraud Detection Model

Problem:

Fraudulent activities such as fake bookings, short-distance rides, or misuse of payment methods can cause heavy losses.

Methodology:

Built a **Random Forest Classifier** to detect fraud using ride attributes like payment method, fare, ratings, and demand level.

Implementation Steps:

- Label encoded `is_fraudulent` as the target column.
- Split data into training and test sets.
- Trained Random Forest on encoded data.
- Evaluated using confusion matrix and accuracy metrics.
- Visualized results using:
 1. Confusion Matrix
 2. Feature Importance Bar Chart
 3. Distribution Plot of fares
 4. Pie Chart of fraud percentage

Results:

- Model achieved ~92% accuracy.
 - Fraud strongly linked to specific payment methods and unusually low distances.
 - **Insight:** Companies can integrate this model into ride verification systems to automatically flag suspicious activities.
-

□ Case 5 – Ride Prediction: Demand Forecasting Model

Problem:

Companies need to anticipate when and where ride demand will increase to avoid shortages and long waiting times.

Methodology:

Used **Linear Regression / ARIMA Forecasting** to predict future ride demand based on time-based data.

Implementation Steps:

- Extracted hourly, daily, and monthly features from timestamps.
- Modeled ride count as a time series.
- Visualized patterns and compared predictions.
- Visualizations included:
 1. Hourly Ride Demand Pattern (Line)
 2. Average Ride Demand by Hour (Bar)
 3. Actual vs Predicted Ride Demand (Scatter)
 4. Forecast Trend (Line over 3 months)

Results:

- The model successfully predicted ride demand peaks.
- Clear double-peak pattern observed: morning (9 AM) and evening (7 PM).
- **Insight:** Helps in scheduling more drivers during peak hours, improving rider satisfaction and reducing wait time.

☒ 4. Results Summary

Case	Key Model	Accuracy/Outcome	Key Insights
1. Smart Pricing	Random Forest Regressor	RMSE: Low	Fare depends on demand, distance & traffic
2. Route Optimization	KMeans Clustering	3 optimal route clusters	Reduced travel time & fuel usage
3. Fast Matching	Random Forest Classifier	Accuracy: 88%	Closer drivers & low traffic = higher efficiency
4. Fraud Detection	Random Forest Classifier	Accuracy: 92%	Fraud linked to payment & ride distance

Case	Key Model	Accuracy/Outcome	Key Insights
5. Ride Prediction	Linear Regression	R ² : 0.89	Strong hourly & weekly demand patterns



5. Tools & Technologies

- **Platform:** Databricks
- **Framework:** Apache Spark (PySpark)
- **Languages:** Python
- **Libraries:** Pandas, Seaborn, Matplotlib, Scikit-learn
- **Visualization Tools:** Dashboards & Python Plots
- **Version Control:** GitHub



6. Final Insights

- Big Data techniques provide powerful insights into travel efficiency, safety, and customer demand.
- Predictive analytics ensures proactive planning rather than reactive problem-solving.
- Real-world implementation of these models can lead to **cost reduction, faster services, and better customer satisfaction.**



7. Future Enhancements

- Integrate with real-time APIs for live data (e.g., Google Maps, weather).
- Deploy dashboards as interactive web apps.
- Extend forecasting with **deep learning (LSTM)** models.
- Introduce **anomaly detection** for dynamic fraud identification.