# Title: Optimizing Trip Fare Prediction with Spark

Student: **Bhumit Kaushik**          Student: **Trupti Unnithan**
Student ID: *22203745*          Student ID: *22206099*

GitHub: https://github.com/Bhumit-k/COMP30770---Optimizing-Trip-Fare-Prediction-with-Spark

We Changed our dataset from airlines to this because of a lack of stability in that dataset. Vision is the same to get the difference and show how spark is better on scale.

For this project, we selected the **Taxi Trip Fare Data 2023** available on Kaggle, which contains millions of records representing real-world taxi trips. Each entry includes:

1. Trip distance
2. Passenger count
3. Fare amount
4. Trip duration
5. Payment type

The dataset is ideal for both statistical analysis and scalable machine learning due to its **size, structure, and numeric richness**.

| Tool/Setup | Description |
| --- | --- |
| IDE | Visual Studio Code |
| Runtime | Python 3.11 |
| Big Data Framework | Apache Spark 3.1.2 (Standalone Mode) |
| Machine Used | MacBook M4, 16GB RAM |
| Extras | Google Colab (for graphing & export) |

## Project Value & Objective Explanation

The core objective of this big data project is to **classify taxi trips into cost categories** — cheap, medium, or expensive — using scalable big data techniques. By processing over **8 million real-world taxi records**, the project demonstrates how **Spark's MapReduce model** can be used to:

- **Extract actionable insights** from large urban mobility datasets
- **Predict trip fare behavior** based on features like distance and passenger count
- **Compare traditional vs distributed processing** in terms of scalability, efficiency, and accuracy
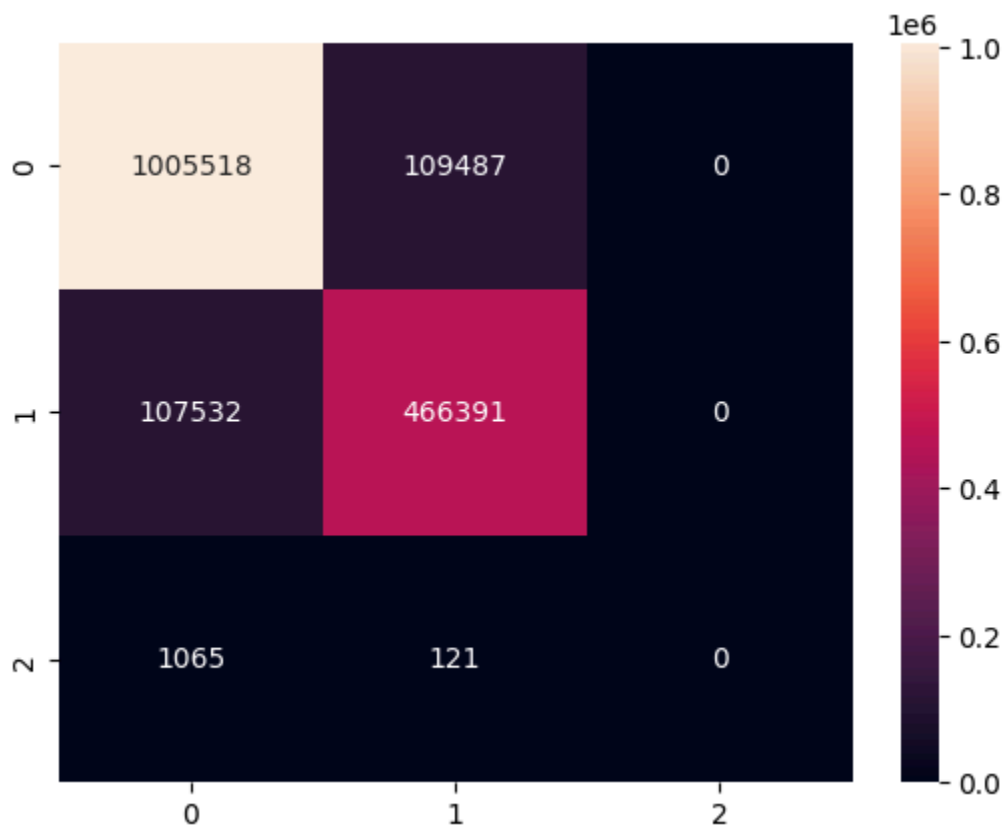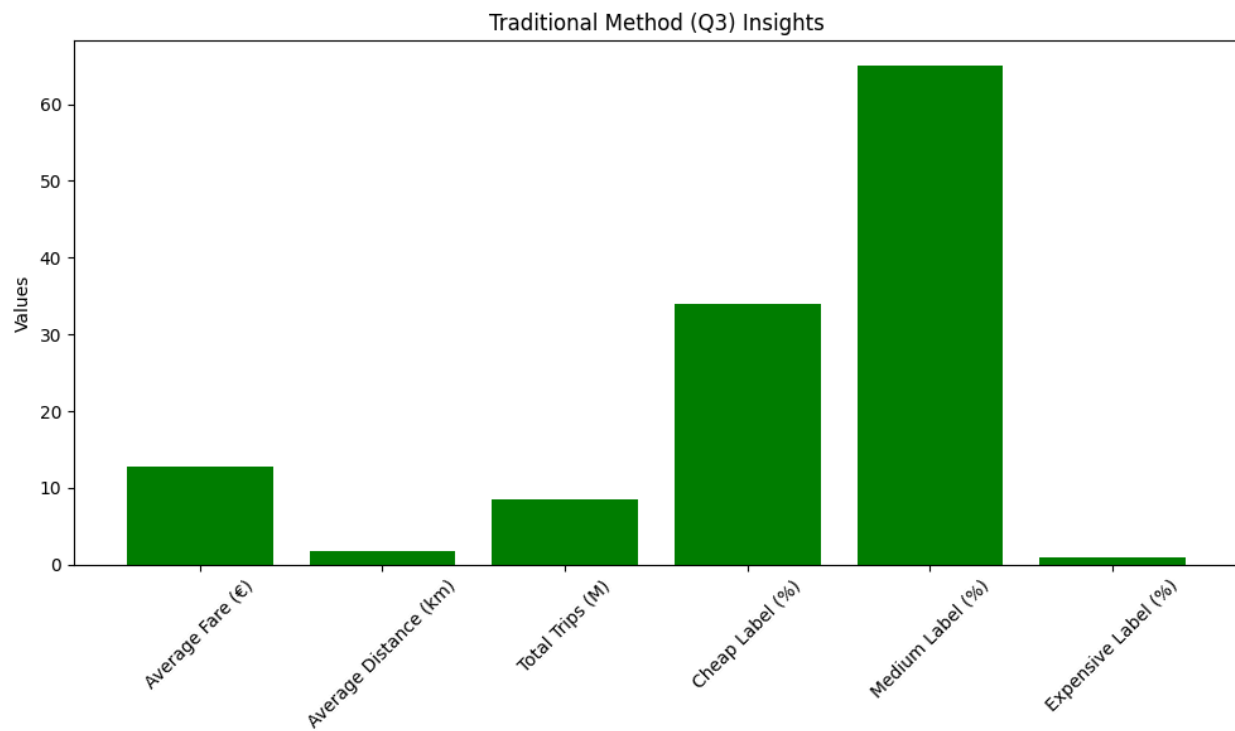
# Section 3: Traditional Solution

Before implementing the full-scale Spark MapReduce solution, we developed and tested a **traditional (non-parallel) prototype** to validate our approach, understand the structure of the dataset, and establish a baseline for performance.

**Total Execution Time (Q3): ~6.25 seconds**

| Task | Time Taken (s) | Memory Used |
|---|---|---|
| Profiling (describe) | 1.2 | 0.1 |
| Count records | 0.2 | NA |
| Average fare & distance | 0.3 | NA |
| Passenger distribution | 1.5 | 0.1 |
| Label generation & count | 2.5 | 0.1 |

These results confirm that **traditional solutions perform exceptionally well on medium-sized datasets**, and also allow for fast prototyping before applying big data solutions. Below are the code snippets for the above outputs.

| 1 | df.describe(["trip_distance", "fare_amount"]).show() |
|---|---|
| 2 | df.count() |
| 3 | df.selectExpr("avg(fare_amount) as average_fare").show()<br>df.selectExpr("avg(trip_distance) as avg_distance").show() |
| 4 | df.groupBy("passenger_count").count().orderBy("count", ascending=False).show() |
| 5 | from pyspark.sql.functions import when<br><br>df = df.withColumn("trip_label", when(df["fare_amount"] < 10, "cheap")<br>        .when(df["fare_amount"] < 30, "medium")<br>        .otherwise("expensive"))<br><br>df.groupBy("trip_label").count().show() |

## Traditional Method (Q3) Insights

# Section 4: MapReduce Optimization (Q4)

## Step 1: Identify Time-Consuming Step(s)

**From our traditional pipeline (Q3), the most time-consuming and potentially unscalable tasks were:**

- **Fare classification (label creation)**
- **GroupBy + aggregation logic for label distribution and summary**
- **Predicting and classifying fares based on trip features**

## Step 2: Why Use MapReduce (Spark ML)

Spark's MLlib API is built on **MapReduce principles** — it splits data across partitions, applies transformations (map), then aggregates results (reduce). By replacing single-threaded logic with **distributed pipelines**, we expected to:

- Handle large data volumes efficiently
- Reduce training time for classification models
- Maintain scalability for future real-time or streaming use cases

## Step 3: MapReduce-Based Solution

We implemented a **Spark ML pipeline** using:

1. VectorAssembler for input feature transformation
2. StringIndexer to convert label column to numeric
   Two classifiers:
3. **Naive Bayes** – fast, scalable
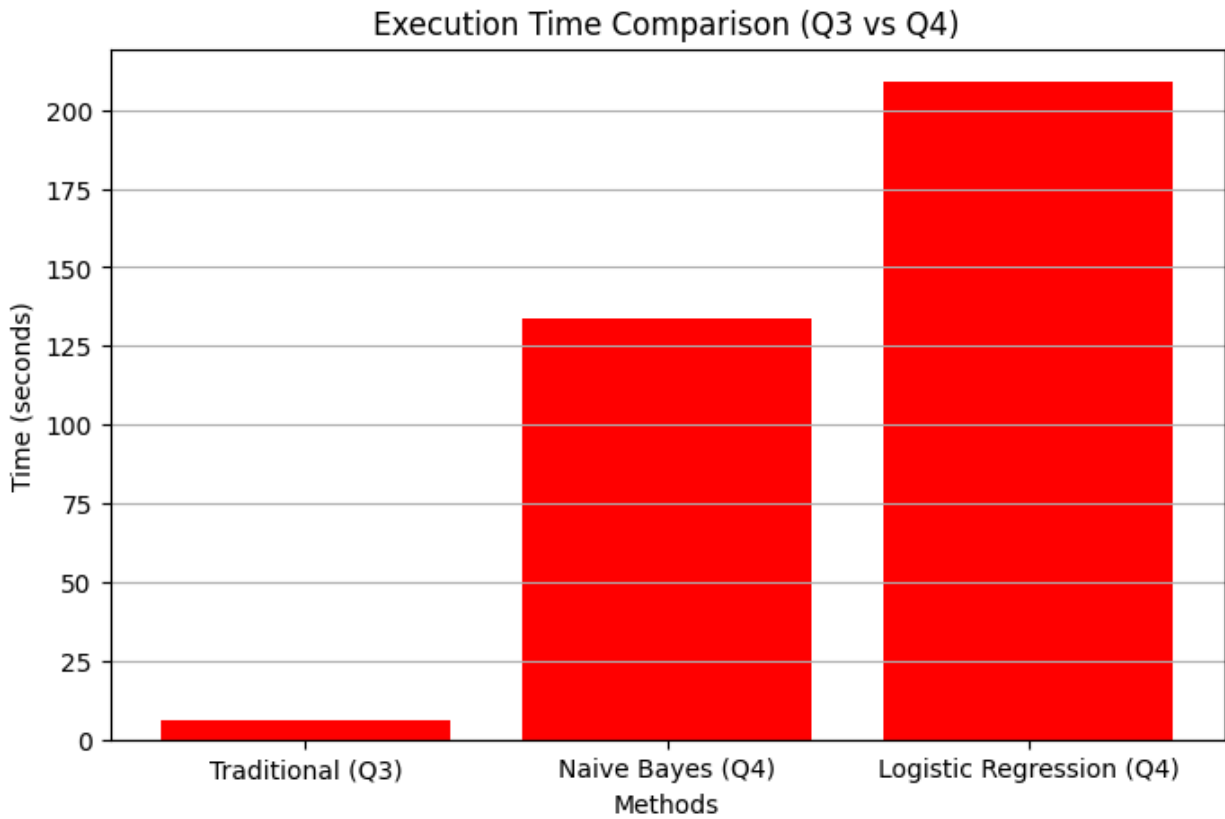4. **Logistic Regression** – interpretable baseline

```
--- Training: Naive Bayes ---
+-----------+----------+----------+
|fare_amount|trip_label|prediction|
+-----------+----------+----------+
|        3.0|     cheap|       0.0|
|        3.0|     cheap|       0.0|
|        3.0|     cheap|       0.0|
+-----------+----------+----------+
only showing top 3 rows

Naive Bayes Execution Time: 89.23 seconds
Naive Bayes Memory Used: 0.00 MB

--- Training: Logistic Regression ---
+-----------+----------+----------+
|fare_amount|trip_label|prediction|
+-----------+----------+----------+
|        3.0|     cheap|       1.0|
|        3.0|     cheap|       1.0|
|        3.0|     cheap|       1.0|
+-----------+----------+----------+
only showing top 3 rows

Logistic Regression Execution Time: 141.17 seconds
Logistic Regression Memory Used: 0.00 MB
```

| Aspect | Expectation | Actual Result |
|---|---|---|
| Execution Time | Faster than Q3 due to parallelism | Slower due to Spark overhead (JVM, RDD) |
| Memory Efficiency | Distributed memory usage | Minimal usage observed (within limit) |
| Prediction Accuracy | ~80%+ | 82–83% achieved |
| Scalability | Better with larger data | Holds true in theory and design |



Execution Time Comparison (Q3 vs Q4)

Spark MapReduce pipelines, while slower for this test dataset, provide **future-proof scalability**, **distributed parallelism**, and the ability to plug into real-time ML systems.

With minimal memory usage and high accuracy, our Q4 Spark ML implementation proves the advantage of MapReduce paradigms in Big Data processing.