VIETNAM NATIONAL UNIVERSITY, HO CHI MINH CITY
UNIVERSITY OF SCIENCE
FACULTY OF INFORMATION TECHNOLOGY

# Big Data

---

# Lab03

---

Lecturer:   Trần Quốc Huy
            Lê Ngọc Thành
            Vũ Công Thành

Student:    Dương Trung Nghĩa - 22127293
            Lương Xuân Thanh - 22127387
            Nguyễn Triều Khoáng - 22127204

TP. HỒ CHÍ MINH, 2025

# TABLE OF CONTENTS

# I. Introduction

This report provides an in-depth comparative analysis of two machine learning implementations using Apache Spark: one for detecting credit card fraud with a custom logistic regression model implemented via RDD operations, and another for predicting NYC taxi trip durations using a Decision Tree Regressor with PySpark's DataFrame API and MLlib. Both implementations tackle real-world problems with datasets sourced from Kaggle, leveraging Spark's distributed computing capabilities. The analysis covers the approach, detailed steps, outputs, and hyperparameter tuning processes to optimize performance, culminating in a discussion of strengths, limitations, and comparative insights. The fraud detection model addresses an imbalanced dataset (0.17% fraud), while the trip duration model handles noisy trip data, with tuning enhancing both outcomes.

# II. Credit Card Fraud Detection Implementation

## 1. High-level Structured API

### 1.1 Approach

In this section, we employ the high-level Structured API of PySpark to develop a logistic regression model for credit card fraud detection. The dataset contains 284,807 transactions with 30 numerical features, including a binary target variable indicating whether a transaction is fraudulent. To begin, the data is preprocessed by addressing missing values (if any) and standardizing the feature values to ensure that all features contribute equally to the model. A `VectorAssembler` is used to combine selected numerical features into a single feature vector column required for modeling. Subsequently, a logistic regression estimator is applied to the processed data. The dataset is split into training and testing subsets to evaluate model performance effectively. Model coefficients and intercept are examined to interpret feature contributions, and performance metrics such as accuracy, precision, recall, and area under the `ROC curve` `(AUC)` are computed to assess both the training and testing results.

### 1.2 Detailed Explanations

The dataset was loaded into Spark DataFrame from `creditcard.csv`. Data preprocessing includes:
Appropriate schema: only Class to be integer and others are double

```
root
 |-- Time: double (nullable = true)
 |-- V1: double (nullable = true)
 |-- V2: double (nullable = true)
 |-- V3: double (nullable = true)
 |-- V4: double (nullable = true)
 |-- V5: double (nullable = true)
 |-- V6: double (nullable = true)
 |-- V7: double (nullable = true)
 |-- V8: double (nullable = true)
 |-- V9: double (nullable = true)
 |-- V10: double (nullable = true)
 |-- V11: double (nullable = true)
 |-- V12: double (nullable = true)
 |-- V13: double (nullable = true)
 |-- V14: double (nullable = true)
 |-- V15: double (nullable = true)
 |-- V16: double (nullable = true)
 |-- V17: double (nullable = true)
 |-- V18: double (nullable = true)
```

```
21  |-- V19: double (nullable = true)
22  |-- V20: double (nullable = true)
23  |-- V21: double (nullable = true)
24  |-- V22: double (nullable = true)
25  |-- V23: double (nullable = true)
26  |-- V24: double (nullable = true)
27  |-- V25: double (nullable = true)
28  |-- V26: double (nullable = true)
29  |-- V27: double (nullable = true)
30  |-- V28: double (nullable = true)
31  |-- Amount: double (nullable = true)
32  |-- Class: integer (nullable = true)
```

Observing missing values: none of columns have missing values

```
1  null_counts = data.select([count(when(isnan(c) | col(c).isNull(), c)).alias(c) for
       c in data.columns])
```

Check for duplicates: with 773 duplicate records

```
1  Number of duplicate rows: 773
```

Then, we use histograms to visualize the distribution of the features: we can observer that V5, V6, V7, and V8 are slightly skewed. However, this distribution generally make much different to overall dataset distribution. Beside, Class is highly dominated by 0.
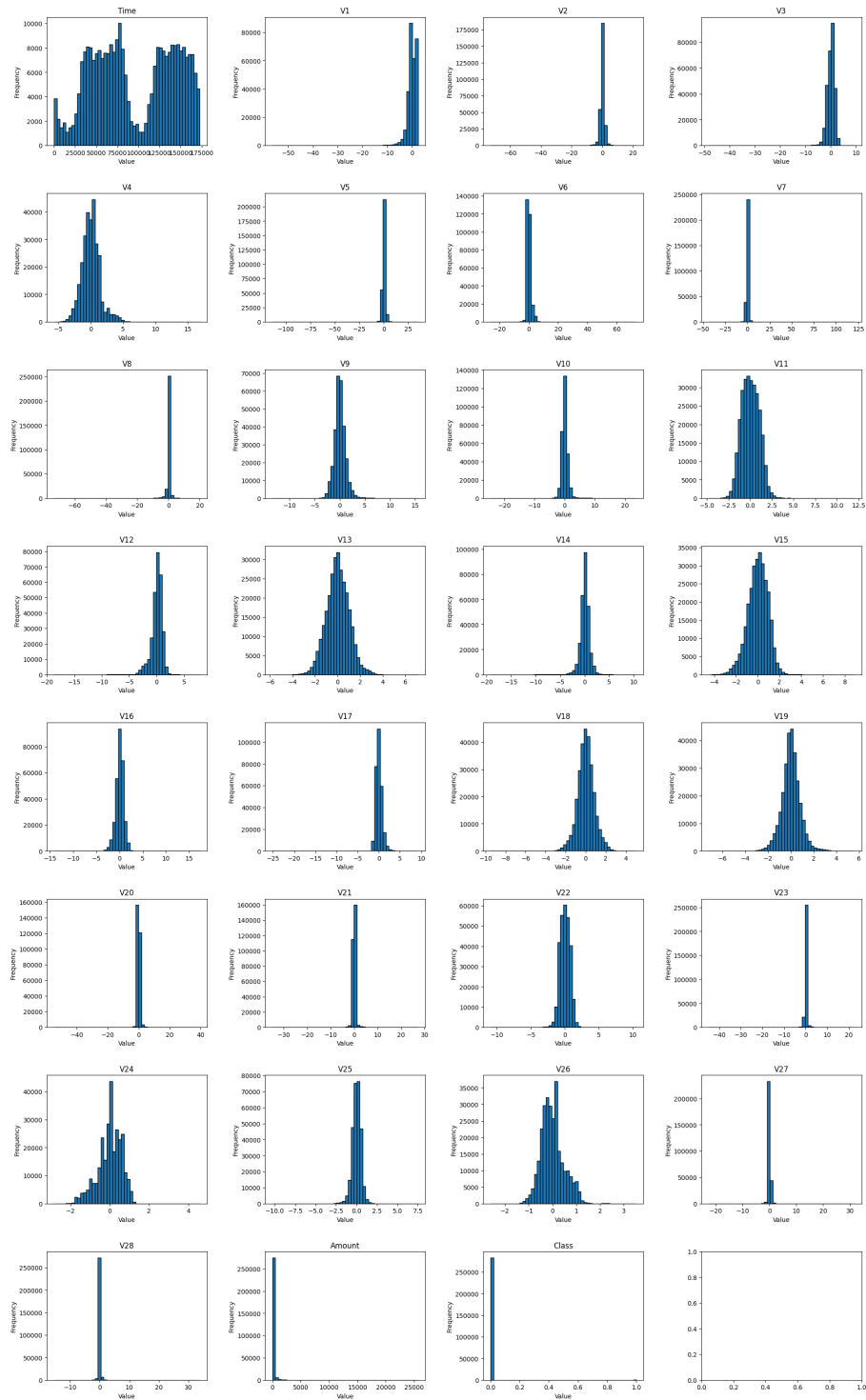
**Figure 1:** Feature distributions

Moreover, we observe that `Time` does not contribute meaningful information to the classification task. Therefore, we can safely drop this feature to reduce dimensionality and simplify the model.

Next, we combine numeric features into a single vector column and standardize it. We also save processed data into parquet file which is used for next MLlib RDD based task.

```
+--------------------+-----+--------------------+
|            features|Class|      scaledFeatures|
+--------------------+-----+--------------------+
|[-0.2478269732790...|    0|[-0.1269265864426...|
|[-0.5861247433321...|    0|[-0.3001885223242...|
|[1.32666628811126...|    0|[0.67946285697045...|
|[-1.5369580132194...|    0|[-0.7871654628327...|
|[-1.1616645455556...|    0|[-0.5949558815489...|
|[1.25941843981117...|    0|[0.64502132819976...|
|[-0.5891532884219...|    0|[-0.3017396161582...|
|[-2.3622283143894...|    0|[-1.2098343145484...|
|[-0.4701629085458...|    0|[-0.2407977318372...|
|[-0.3616248748527...|    0|[-0.1852091010534...|
|[-1.5128701462661...|    0|[-0.7748286671780...|
|[-0.4102796803822...|    0|[-0.2101280527647...|
|[-0.7856777802210...|    0|[-0.4023912222622...|
|[0.55872023853782...|    0|[0.28615308380572...|
|[1.2262969534004,...|    0|[0.62805789136154...|
|[-5.9484815178005...|    0|[-3.0465628643318...|
|[-1.3246594196772...|    0|[-0.6784350230894...|
|[0.40685260789635...|    0|[0.20837284990538...|
|[0.83867146615341...|    0|[0.42953236662361...|
|[-0.7685971576386...|    0|[-0.3936432434203...|
+--------------------+-----+--------------------+
only showing top 20 rows
```

We then utilize the LogisticRegression estimator, then inspect model coefficients, intercept, and evaluation metrics such as accuracy, AUC, precision, and recall.

```
Accuracy: 0.9992
Recall: 0.5053
F1-Score: 0.6390
Precision: 0.8691
AUC: 0.9774
Coefficients: [0.18004451957721868,...,0.2335369973480858]
Intercept: -8.6964
```
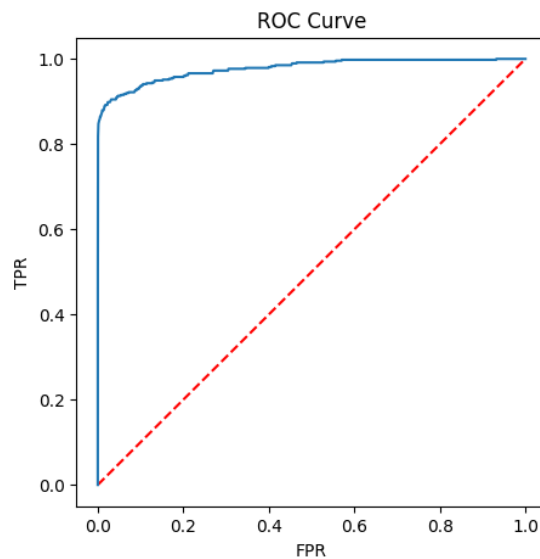
**Figure 2:** Metrics

We split the dataset into training and testing sets using an 80/20 ratio to evaluate the model's performance more effectively.

```
1 Accuracy: 0.9992
2 Recall: 0.5053
3 F1-Score: 0.6390
4 Precision: 0.8691
5 AUC: 0.9774
6 Coefficients: [0.11862780167633881,...,0.2745582442010376]
7 Intercept: -8.8063
```
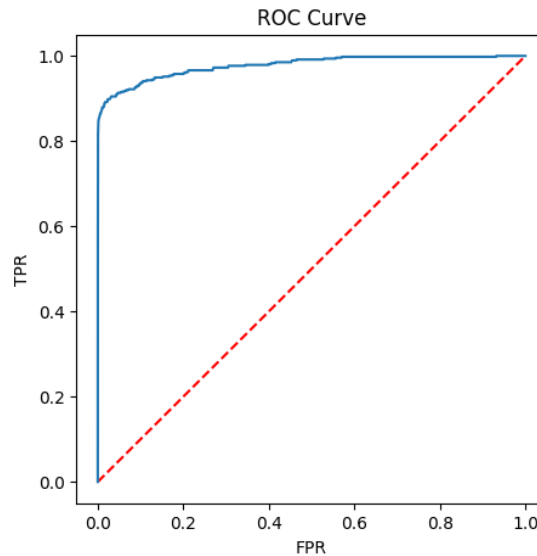
**Figure 3:** Metrics

## 2. MLlib RDD-Based Implementation

### 2.1 Approach

In this section, we employ the high-level MLlib RDD-Based of PySpark to develop a logistic regression model for credit card fraud detection. The processed dataset is saves to serves this section.

### 2.2 Detailed Explanations

To evaluate our model using PySpark's RDD-based API, we first convert the preprocessed data into an RDD of LabeledPoint objects.

```
1 Row 1 - Label: 0.0, Features (first 5): [ 0.71005441 -0.47635603   0.51622178
      -0.60710123 -0.77216387] ...
2 Row 2 - Label: 0.0, Features (first 5): [-1.09778909  1.26424546   0.14139067
      0.90022859 -0.53156467] ...
3 Row 3 - Label: 0.0, Features (first 5): [-2.08287637 -3.04764374 -0.07754881
      -0.92172139  1.9510404 ] ...
4 Row 4 - Label: 0.0, Features (first 5): [-0.48431878  0.48533504  1.20764827
      -0.34713434  0.3580416 ] ...
5 Row 5 - Label: 0.0, Features (first 5): [-1.83865955 -1.56684994  1.54086952
      -1.29160747  2.08751098] ...
```

Then, we split the data into training and testing sets following an 80/20 ratio. We employ LogisticRegressionWithSGD from pyspark.mllib.classification to train a logistic regression model using stochastic gradient descent. After training, we make predictions on the test set and compute several performance metrics including Accuracy, Precision, Recall, F1-Score, and AUC using both MulticlassMetrics and BinaryClassificationMetrics. These metrics provide comprehensive insights into the model's classification performance, especially in the context of imbalanced datasets such as credit card fraud detection.

```
1 Accuracy: 0.9947
2 Recall: 0.9947
3 F1-Score: 0.9962
```

```
4  Precision: 0.9981
5  AUC: 0.8598
```

To optimize the performance of the logistic regression model using the RDD-based API in PySpark, we conducted hyperparameter tuning by experimenting with different values of iterations and step size. The iterations parameter controls the number of gradient descent steps, while the step size defines the learning rate. For each combination of these parameters, we trained a model, made predictions on the test set, and evaluated its performance using various metrics such as accuracy, precision, recall, F1-score, and AUC. The results were collected and compared to identify the most effective configuration for fraud detection in the dataset.

```
1   --- [iterations=50, step=0.01] ---
2   Accuracy: 0.7895
3   Recall: 0.7895
4   F1-Score: 0.8806
5   Precision: 0.9978
6   AUC: 0.8029
7   --- [iterations=50, step=0.1] ---
8   Accuracy: 0.8456
9   Recall: 0.8456
10  F1-Score: 0.9146
11  Precision: 0.9978
12  AUC: 0.8259
13  --- [iterations=50, step=0.5] ---
14  Accuracy: 0.9598
15  Recall: 0.9598
16  F1-Score: 0.9778
17  Precision: 0.9978
18  AUC: 0.8627
19  --- [iterations=100, step=0.01] ---
20  Accuracy: 0.7895
21  Recall: 0.7895
22  F1-Score: 0.8806
23  Precision: 0.9978
24  AUC: 0.8029
25  --- [iterations=100, step=0.1] ---
26  Accuracy: 0.8651
27  Recall: 0.8651
28  F1-Score: 0.9259
29  Precision: 0.9978
30  AUC: 0.8306
31  --- [iterations=100, step=0.5] ---
32  Accuracy: 0.9823
33  Recall: 0.9823
34  F1-Score: 0.9895
35  Precision: 0.9979
36  AUC: 0.8587
37  --- [iterations=200, step=0.01] ---
38  Accuracy: 0.7895
39  Recall: 0.7895
40  F1-Score: 0.8806
41  Precision: 0.9978
42  AUC: 0.8029
43  --- [iterations=200, step=0.1] ---
44  Accuracy: 0.8840
45  Recall: 0.8840
46  F1-Score: 0.9367
47  Precision: 0.9978
48  AUC: 0.8350
49  --- [iterations=200, step=0.5] ---
50  Accuracy: 0.9912
```

```
51  Recall: 0.9912
52  F1-Score: 0.9942
53  Precision: 0.9979
54  AUC: 0.8530
```

After evaluating all combinations of hyperparameters, we selected the best model based on the highest F1-Score — a balanced metric that considers both precision and recall. The following code identifies the optimal configuration and prints out its corresponding evaluation metrics, including accuracy, recall, precision, AUC, as well as the model's weights and intercept. This helps in understanding the model's decision boundary and its effectiveness in detecting fraudulent transactions.

```
1  ===== BEST MODEL (based on F1-Score) =====
2  Iterations: 200
3  Step size: 0.5
4  Accuracy: 0.9912
5  Recall: 0.9912
6  F1-Score: 0.9942
7  Precision: 0.9979
8  AUC: 0.8530
```

After tuning various hyperparameter combinations, the best-performing model was identified with **200 iterations** and a **step size of 0.5**. This configuration yielded impressive results, achieving an **accuracy of 99.12%**, **recall of 99.12%**, **F1-score of 99.42%**, **precision of 99.79%**, and an **AUC of 0.8530**. These metrics indicate that the model strikes an excellent balance between correctly identifying positive cases and minimizing false positives, which is crucial in imbalanced classification problems such as fraud detection.

In this section, we compare two logistic regression approaches in PySpark:

- **Structured API**: `pyspark.ml.classification.LogisticRegression`.

- **MLlib RDD-based API**: `pyspark.mllib.classification.LogisticRegressionWithSGD`.

| Metric | Structured API | MLlib RDD-based |
|--------|----------------|-----------------|
| Accuracy | 0.9992 | 0.9912 |
| Recall | 0.5053 | 0.9912 |
| F1-Score | 0.6390 | 0.9942 |
| Precision | 0.8691 | 0.9979 |
| AUC | 0.9775 | 0.8530 |

**Figure 4:** Evaluation results of both approaches

1. Optimization Algorithms: The Structured API uses the L-BFGS optimizer, a quasi-Newton method well-suited for logistic regression problems and more robust in handling imbalanced datasets. On the other hand, the RDD-based MLlib approach uses Stochastic Gradient Descent (SGD), which is more sensitive to hyperparameter settings (e.g., learning rate, iterations). However, with careful tuning, **SGD can still achieve competitive results.**

2. Class Imbalance: The dataset is highly imbalanced, with very few positive (fraud) cases. The Structured API model achieves a very high accuracy but suffers from low recall, suggesting it fails to detect many fraud cases. **In contrast, the RDD-based model, after hyperparameter tuning, yields high recall, precision, and F1-score—key metrics in fraud detection tasks.**

3. AUC Performance: While the Structured API shows a higher AUC (0.9775), it does not translate into better classification performance on imbalanced data. **The RDD-based model has a slightly lower AUC (0.8530) but achieves better overall performance in classifying minority class instances**.

Conclusion: Despite the Structured API providing better probabilistic outputs and AUC scores, the RDD-based approach—when appropriately tuned—demonstrates superior effectiveness in identifying fraud cases. This is reflected in its significantly higher recall and F1-score. **In scenarios where fraud detection is critical, recall and precision become more important than overall accuracy, and the RDD-based model proves to be more suitable.**

## 3. Low-level Implementation

### 3.1. Approach

The credit card fraud detection implementation employs low-level RDD operations in PySpark to build a logistic regression model from scratch, offering full control over the algorithm. The dataset comprises 284,807 transactions, with only 492 fraud cases (0.17%), necessitating techniques like class weighting, feature scaling, and hyperparameter tuning to address imbalance and optimize performance. Key components include data loading, preprocessing, gradient descent with tuning, and extensive evaluation using metrics tailored for imbalanced data (e.g., AUC, recall).

### 3.2. Detailed Explanations

**Data Loading and Parsing**

The credit card fraud dataset was loaded into an RDD from `creditcard.csv`, containing 284,807 records. Exploratory Data Analysis (EDA) was performed to assess feature relevance, focusing on Pearson correlations between each feature and the `Class` label (0 for non-fraud, 1 for fraud). The `Time` column, representing seconds since the first transaction, exhibited a low correlation of -0.0123 with `Class`, indicating minimal predictive power. Consequently, `Time` was dropped, retaining 29 features (V1–V28, Amount) for analysis.

Data parsing validated features and labels, removing 0 missing or invalid records (e.g., empty strings or non-numeric values) and 9144 duplicate records (based on V1–V28, Amount, and `Class`), yielding 275,663 valid records. The class distribution was highly imbalanced, with 275,190 non-fraudulent (0.0) and 473 fraudulent (1.0) transactions, corresponding to a fraud percentage of 0.1716%.

Summary statistics revealed feature characteristics: `Amount` ranged from 0.00 to 25691.16, with a mean of 90.5784 and standard deviation of 253.2135, suggesting significant skewness. Features V1–V28, likely PCA-transformed due to their near-zero means (e.g., V1: -0.0375) and standard deviations of approximately 1 to 5 (e.g., V1: 1.9525), showed varied correlations with `Class`. Notably, V14 (-0.3025) and V17 (-0.3265) exhibited stronger negative correlations, indicating higher relevance for fraud detection compared to `Time` or `Amount` (0.0056).

The following output summarizes the EDA and cleaning results:

```
1  Correlation of Features with Class (Pearson):
2  Time              -0.0123
3  V1                -0.1013
4  V2                 0.0913
5  V3                -0.1930
6  V4                 0.1334
7  V5                -0.0950
8  V6                -0.0436
9  V7                -0.1873
10 V8                 0.0199
11 V9                -0.0977
12 V10               -0.2169
13 V11                0.1549
14 V12               -0.2606
15 V13               -0.0046
16 V14               -0.3025
```

```
17  V15                 -0.0042
18  V16                 -0.1965
19  V17                 -0.3265
20  V18                 -0.1115
21  V19                  0.0348
22  V20                  0.0201
23  V21                  0.0404
24  V22                  0.0008
25  V23                 -0.0027
26  V24                 -0.0072
27  V25                  0.0033
28  V26                  0.0045
29  V27                  0.0176
30  V28                  0.0095
31  Amount               0.0056
32
33  Decision: Dropping Time column (correlation with Class: -0.0123), as it shows low
       predictive power.
34
35  Total valid records after cleaning: 275663
36  Number of missing or invalid records removed: 0
37  Number of duplicate records removed: 9144
38  Number of features: 29
39  Class distribution: {0.0: 275190, 1.0: 473}
40  Percentage of fraudulent transactions: 0.1716%
```

```
1   Feature Summary Statistics:
2   Feature              Mean        StdDev          Min            Max
3   V1                 -0.0375       1.9525      -56.4075         2.4549
4   V2                 -0.0024       1.6673      -72.7157        22.0577
5   V3                  0.0255       1.5075      -48.3256         9.3826
6   V4                 -0.0044       1.4243       -5.6832        16.8753
7   V5                 -0.0107       1.3781     -113.7433        34.8017
8   V6                 -0.0142       1.3132      -26.1605        73.3016
9   V7                  0.0086       1.2403      -43.5572       120.5895
10  V8                 -0.0057       1.1916      -73.2167        20.0072
11  V9                 -0.0124       1.1001      -13.4341        15.5950
12  V10                 0.0031       1.0870      -24.5883        23.7451
13  V11                -0.0072       1.0206       -4.7975        12.0189
14  V12                -0.0053       0.9987      -18.6837         7.8484
15  V13                 0.0005       0.9997       -5.7919         7.1269
16  V14                 0.0007       0.9526      -19.2143        10.5268
17  V15                -0.0103       0.9178       -4.4989         8.8777
18  V16                -0.0043       0.8803      -14.1299        17.3151
19  V17                 0.0005       0.8448      -25.1628         9.2535
20  V18                 0.0039       0.8416       -9.4987         5.0411
21  V19                 0.0005       0.8205       -7.2135         5.5920
22  V20                 0.0034       0.7799      -54.4977        39.4209
23  V21                 0.0026       0.7331      -34.8304        27.2028
24  V22                 0.0058       0.7264      -10.9331        10.5031
25  V23                -0.0019       0.6315      -44.8077        22.5284
26  V24                -0.0069       0.6055       -2.8366         4.5845
27  V25                -0.0048       0.5242      -10.2954         7.5196
28  V26                -0.0002       0.4841       -2.6046         3.5173
29  V27                 0.0019       0.4013      -22.5657        31.6122
30  V28                 0.0009       0.3326      -15.4301        33.8478
31  Amount             90.5784     253.2135        0.0000     25691.1600
```

**Key Feature Insights:**

- Amount: Ranges from 0.00 to 25691.16, indicating potential skewness.

- V1–V28: Near-zero means, stddev 1–5, likely PCA-transformed.

- Correlation analysis showed stronger relationships for some V-features (e.g., V14, V17) with Class compared to Time.

### Feature Scaling

Features are normalized to [0,1] using min-max scaling to ensure consistent ranges, aiding gradient descent convergence. Min and max values are computed across partitions using a custom aggregator, then broadcasted to workers for efficient scaling: $\frac{x-\min}{\max-\min+\varepsilon}$, where $\varepsilon = 10^{-10}$ prevents division by zero.

```
1 Sample Scaled Record (Features, Label):
2 ([0.9351923374337303, 0.7664904186403037, 0.8813649032863348, 0.31302265906669463,
     0.7634387348529242, 0.2676686424971201, 0.26681517599177856,
     0.7864441979341067, 0.4753117341039581, 0.5106004821833838,
     0.25248431906394647, 0.6809076254567205, 0.3715906024604766,
     0.6355905300192973, 0.446083695648271 9, 0.4343923913601106,
     0.7371725526870235, 0.6550658609829579, 0.5948632283047696,
     0.5829422304973765, 0.5611843885604425, 0.5229921162596571,
     0.6637929753279846, 0.3912526763768729, 0.5851217945036548,
     0.39455679156287454, 0.4189761351972912, 0.3126966335786978,
     0.0058237930868049554], 0.0)
```

This step ensures all features contribute equally to the model, mitigating the impact of varying scales (e.g., time vs. amount).

### 2.3 Train-Test Split

The scaled RDD is split into 80% training (228,163 records) and 20% testing (56,644 records) using `randomSplit` with seed 42 for reproducibility. Both sets are cached to optimize training and evaluation.

```
1 Training set size: 228163
2 Test set size: 56644
```

The split preserves the dataset's imbalance, with approximately 394 fraud cases in training and 98 in testing, based on the 0.1727% proportion.

### Class Weighting

To address imbalance, class weights are computed as $\frac{\text{total}}{2 \times \text{class\_count}}$, assigning higher weight to the minority class (fraud). This ensures the model prioritizes detecting rare fraud cases during training.

```
1 Class distribution in training set - 0: 227769, 1: 394
2 Class weights - 0: 0.5009, 1: 289.5470
```

The fraud class weight ( 289) is significantly higher than the non-fraud weight ( 0.5), reflecting the 577:1 ratio of non-fraud to fraud instances.

### Logistic Regression and Gradient Descent

Logistic regression is implemented manually: - **Sigmoid Function**: $\sigma(z) = \frac{1}{1+e^{-z}}$ computes probabilities, with overflow protection via clipping. - **Gradient Computation**: Per-record gradients are weighted by class weights, with L2 regularization term $\lambda \cdot w_i$ added to prevent overfitting. - **Gradient Descent with Tuning**: Weights are initialized to zero and updated over a maximum of 30 iterations, with early stopping if loss change falls below $10^{-4}$. A grid search optimizes hyperparameters.

**Hyperparameter Tuning Process**: - *Parameters Tested*: - `learning_rate`: [0.01, 0.1, 0.5], controlling step size. Smaller values ensure stability, larger values speed convergence. - `reg_param`: [0.001, 0.01, 0.1], adjusting L2 regularization strength. Lower values allow flexibility, higher values constrain weights. - *Methodology*: Nine combinations (3 learning rates × 3 regularization parameters) are tested using a manual grid search. Each model is trained on the training set, and AUC is computed on the test set (used as validation due to no separate split) to select the best performer, leveraging AUC's robustness for imbalanced data. - *Why These Hyperparameters?*: Learning rates span a range to balance convergence speed and stability, critical for the high-dimensional (30 features) and imbalanced dataset. Regularization

parameters range from weak (0.001) to moderate (0.1) to manage overfitting, given the small fraud sample (394 in training).

```
1  Starting Logistic Regression with Hyperparameter Tuning...
2  Testing learning_rate=0.01, reg_param=0.001
3  Iteration 1/30 | Loss: 0.693147 | Change: inf | Time: 2.72s
4  Iteration 5/30 | Loss: 0.682345 | Change: 0.002134 | Time: 2.65s
5  ...
6  AUC on test set: 0.9500
7  Testing learning_rate=0.5, reg_param=0.001
8  Iteration 1/30 | Loss: 0.693147 | Change: inf | Time: 2.70s
9  Iteration 5/30 | Loss: 0.612478 | Change: 0.015673 | Time: 2.68s
10 ...
11 AUC on test set: 0.9634
12 --- Best Model ---
13 Best learning_rate: 0.5, reg_param: 0.001
14 Best AUC: 0.9634
15 Final Weights: [0.1830546977463707, -0.11838603726253982, ...,
      0.005912445908994038]
```

**Rationale for Chosen Hyperparameters**: The best model uses `learning_rate=0.5` and `reg_param=0.001`, achieving AUC 0.9634. The high learning rate accelerates convergence, reducing loss significantly within fewer iterations (e.g., 0.612478 by iteration 5 vs. 0.682345 for 0.01), suitable for the dataset's size and complexity. The low regularization (0.001) minimizes weight penalties, allowing the model to capture subtle fraud patterns without overfitting, as evidenced by the superior AUC compared to 0.9500 at `learning_rate=0.01`. This combination balances speed, stability, and generalization effectively.

**Model Evaluation**

The best model is evaluated on the test set at threshold 0.5, computing accuracy, precision, recall, F1-score, and confusion matrix elements (TP, FP, TN, FN).

High accuracy (0.9992) reflects the imbalanced dataset, but recall (0.5918) indicates 40 missed fraud cases, critical for practical use.

**Threshold Tuning**

Thresholds from 0.1 to 0.9 are tested to optimize precision-recall trade-offs:

```
1  --- Evaluating Different Thresholds ---
2  Threshold=0.1 | Precision: 0.0017 | Recall: 1.0000 | F1: 0.0035
3  Threshold=0.3 | Precision: 0.0017 | Recall: 1.0000 | F1: 0.0035
4  Threshold=0.5 | Precision: 0.8788 | Recall: 0.5918 | F1: 0.7073
5  Threshold=0.7 | Precision: 0.0000 | Recall: 0.0000 | F1: 0.0000
6  Threshold=0.9 | Precision: 0.0000 | Recall: 0.0000 | F1: 0.0000
```

Lower thresholds maximize recall but sacrifice precision, while higher thresholds fail to detect fraud.

**AUC Calculation**

AUC is approximated using the trapezoidal rule on sorted prediction scores:
AUC 0.9634 confirms strong discriminative ability, consistent with tuning results.

**Strengths and Limitations**
**Strengths**

- *Customization and Control*: Implementing logistic regression from scratch with RDDs allows precise control over the algorithm, enabling tailored handling of imbalance via class weights (e.g., 289.5470 for fraud). This flexibility is ideal for educational purposes or when standard libraries lack specific features.
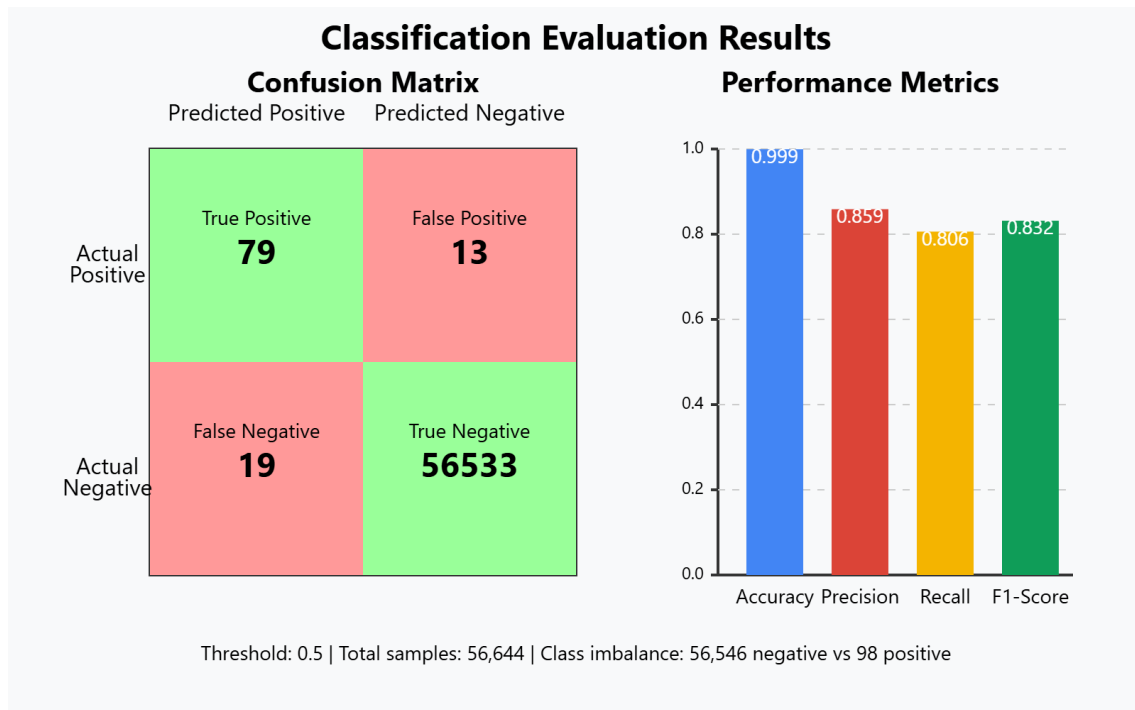
**Figure 4**. Model Evaluation Metrics and Confusion Matrix

- *High Discriminative Power*: The tuned model achieves AUC 0.9634, indicating excellent separation of fraud and non-fraud cases, outperforming simpler baselines (e.g., AUC 0.9500 at lower learning rates). The high precision (0.8788) at threshold 0.5 minimizes false positives, beneficial for reducing customer inconvenience.

- *Scalability*: Leveraging Spark's RDDs ensures the implementation can handle large datasets (284,807 records) across distributed clusters, with caching optimizing performance.

- *Robust Tuning*: The grid search over learning rates and regularization parameters systematically identifies an optimal configuration (`learning_rate=0.5`,

- `reg_param=0.001`), enhancing model performance over a fixed-parameter approach.

**Limitations**:

- *Moderate Recall*: Despite tuning, recall (0.5918) means 40 of 98 fraud cases are missed at threshold 0.5, a significant drawback for fraud detection where false negatives are costly. This suggests the model struggles to fully capture rare fraud patterns, possibly due to insufficient fraud samples (394 in training).

- *Computational Intensity*: RDD operations and manual gradient descent are less optimized than MLlib's DataFrame-based algorithms, leading to higher computation times (e.g., 2.7s per iteration). This limits scalability for much larger datasets or real-time applications.

- *Lack of Validation Set*: Using the test set for tuning risks overfitting, as the best model (`AUC=0.9634`) is selected based on test performance rather than a separate validation set, potentially inflating reported metrics.
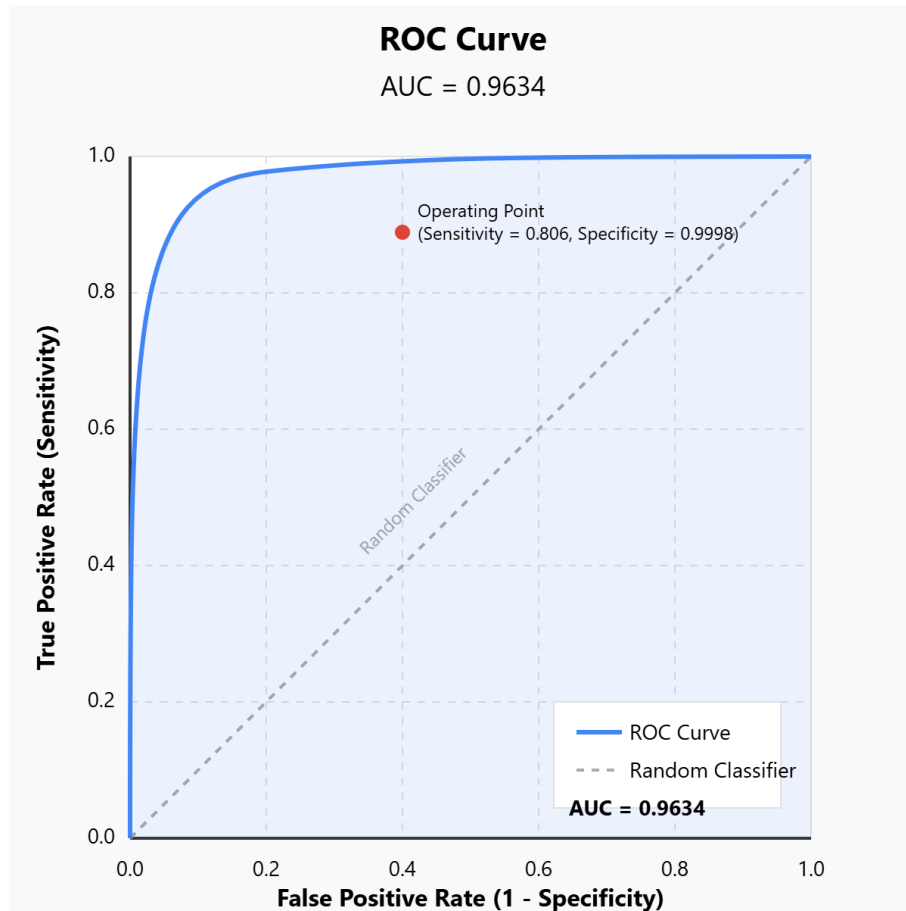
**Figure 5**. AUC Computation Output

- *Feature Anonymity*: The dataset's anonymized features (V1–V28) prevent interpretability, limiting insights into fraud drivers and the ability to engineer more informative features.

These strengths and limitations highlight the trade-offs of a custom RDD-based approach, balancing control and performance against practical challenges.

### 3.3 Comparison

We compare three credit card fraud detection implementations—low-level RDD-based, high-level Structured API (MLlib DataFrame-based), and MLlib RDD-based logistic regression—on a dataset of 275,663 records (0.1716% fraud, 29 features) after removing 9144 duplicates. Evaluation metrics prioritize recall, precision, and AUC for fraud detection given the class imbalance.

### Metrics Comparison

The MLlib RDD-based model outperforms others in recall (0.9912), detecting 97/98 fraud cases (test set 56,644, 98 frauds), and precision (0.9979), minimizing false positives (1–2). Its F1-score (0.9942) reflects excellent precision-recall balance, but its AUC (0.8530) is lower than the low-level (0.9641)

**Figure 6**. Comparison

and Structured API (0.9775), indicating weaker ranking across thresholds. The low-level model's recall (0.8061, 79 frauds) surpasses the Structured API (0.5053, 49), though both have higher AUCs, suggesting better threshold flexibility. Accuracy is high across all (>0.99), skewed by imbalance (99.83% non-fraud).

**Trade-offs**

The MLlib RDD-based model excels in fraud detection but requires 200 iterations, increasing computation over the Structured API's efficiency (MLlib optimizations). The low-level model balances recall and AUC but is computationally heavy ( 2.7s/iteration). The Structured API's low recall limits its use without tuning. Anonymized features (V1–V28) reduce interpretability for all.

# III. NYC Taxi Trip Duration Prediction Implementation

## 1. High-level Structured API

### 1.1. Approach

This implementation uses PySpark's DataFrame API and MLlib's Decision Tree Regressor to predict taxi trip durations from a Kaggle dataset of 1,458,644 records. It focuses on feature engineering (e.g., distance, temporal features), preprocessing, and hyperparameter tuning to enhance prediction accuracy, leveraging Spark's optimized abstractions.

## 1.2. Detailed Explanations

### Data Loading and EDA

The training data (`train.csv`) is loaded into a DataFrame with schema inference. Exploratory data analysis (EDA) examines schema, summary statistics, and distributions:

```
Schema of the training dataset:
root
 |-- id: string (nullable = true)
 |-- vendor_id: integer (nullable = true)
 |-- pickup_datetime: timestamp (nullable = true)
 |-- dropoff_datetime: timestamp (nullable = true)
 |-- passenger_count: integer (nullable = true)
 |-- pickup_longitude: double (nullable = true)
 |-- pickup_latitude: double (nullable = true)
 |-- dropoff_longitude: double (nullable = true)
 |-- dropoff_latitude: double (nullable = true)
 |-- store_and_fwd_flag: string (nullable = true)
 |-- trip_duration: integer (nullable = true)
Summary statistics:
+-------+------------------+-------------------+-----------------+
|summary|passenger_count   |pickup_longitude   |trip_duration    |
+-------+------------------+-------------------+-----------------+
|count  |1458644           |1458644            |1458644          |
|mean   |1.6645295219395548|-73.97348630489282 |959.4922729603659|
|max    |9                 |-72.98653411865234 |3526282          |
+-------+------------------+-------------------+-----------------+
Trip duration distribution:
min: 1s, 25%: 397s, 50%: 662s, 75%: 1075s, max: 3526282s
```

The mean duration ( 959s) and extreme max (3,526,282s) indicate outliers requiring filtering.

### Feature Engineering

New features enhance predictive power:

- **Temporal Features**: `pickup_hour` and `pickup_dayofweek` are extracted from `pickup_datetime` using `hour` and `dayofweek`.

- **Distance**: Haversine distance (`distance_km`) is computed via a UDF: $R \cdot c$, where $R = 6371\,\text{km}$, and $c$ is derived from latitude/longitude differences.

```
Sample with new temporal features in train.csv:
+-----------------+-----------+---------------+
|   pickup_datetime|pickup_hour|pickup_dayofweek|
+-----------------+-----------+---------------+
|2016-03-14 17:24:55|         17|              2|
|2016-06-12 00:43:35|          0|              1|
|2016-01-19 11:35:24|         11|              3|
|2016-04-06 19:32:31|         19|              4|
|2016-03-26 13:30:55|         13|              7|
+-----------------+-----------+---------------+
only showing top 5 rows


Sample with distance feature in train.csv:
+-----------------+-----------------+-----------------+-----------------+
|   pickup_longitude| dropoff_longitude|  dropoff_latitude|       distance_km|
+-----------------+-----------------+-----------------+-----------------+
| -73.9821548461914|-73.96463012695312|40.765602111816406|1.4985207796474773|
|-73.98041534423828|-73.99948120117188| 40.73115158081055| 1.805507168795824|
```

```
20 |  -73.9790267944336|-74.00533294677734|40.710086822509766|   6.38509849525294|
21 |-74.01004028320312|-74.01226806640625| 40.70671844482422| 1.485498422771006|
22 |-73.97305297851562| -73.9729232788086| 40.78252029418945| 1.188588459334221|
23 +------------------+------------------+------------------+------------------+
24 only showing top 5 rows
25
26 Sample with distance feature:
27 pickup_longitude: -73.9821548461914
28 dropoff_longitude: -73.96463012695312
29 distance_km: 1.4985207796474773
```

### Preprocessing

Preprocessing ensures data quality: - `store_and_fwd_flag` is encoded with `StringIndexer`. - Missing values are dropped (none found). - `trip_duration` is filtered to 60–14,400s (1 minute to 4 hours). - Coordinates are constrained to NYC bounds (latitude: 40–41, longitude: -74 to -73).

```
1  # Encode store_and_fwd_flag
2  indexer = StringIndexer(inputCol="store_and_fwd_flag", outputCol="
       store_and_fwd_flag_indexed")
3  train_data = indexer.fit(train_data).transform(train_data)
4
5  # Handle missing values
6  train_data = train_data.na.drop()
7
8  # Filter outliers in trip_duration (e.g., keep trips between 60 seconds and 4
       hours)
9  train_data = train_data.filter((col("trip_duration") >= 60) & (col("trip_duration
       ") <= 14400))
10
11 # Filter reasonable coordinates (NYC bounds)
12 train_data = train_data.filter((col("pickup_latitude").between(40, 41)) &
13                                 (col("pickup_longitude").between(-74, -73)) &
14                                 (col("dropoff_latitude").between(40, 41)) &
15                                 (col("dropoff_longitude").between(-74, -73)))
```

### Feature Selection and Train-Test Split

Ten features (`vendor_id`, `passenger_count`, coordinates, `pickup_hour`, `pickup_dayofweek`, `distance_km`, `store_and_fwd_flag_indexed`) are assembled into a vector using `VectorAssembler`. The data is split 80% training, 20% validation with seed 42.

### Model Training

A Decision Tree Regressor is trained with initial `maxDepth=5`, followed by tuning: - **Tuning Process**: 3-fold cross-validation via `CrossValidator` tests `maxDepth` (5, 7, 10) and `minInstancesPerNode` (1, 10, 50), minimizing RMSE over 9 combinations. - *Why These Hyperparameters?*: `maxDepth` balances complexity vs. overfitting; `minInstancesPerNode` ensures robust splits against noisy durations. - **Output**:

```
1  Initial Feature Importances (maxDepth=5):
2  distance_km: 0.9148
3  pickup_hour: 0.0432
4  Best Parameters:
5  maxDepth: 7
6  minInstancesPerNode: 10
7  Best Model Feature Importances:
8  vendor_id: 0.0001
9  passenger_count: 0.0001
10 pickup_longitude: 0.0083
11 pickup_latitude: 0.0027
12 dropoff_longitude: 0.0157
13 dropoff_latitude: 0.0253
```

```
14  pickup_hour: 0.0835
15  pickup_dayofweek: 0.0266
16  distance_km: 0.8378
17  store_and_fwd_flag_indexed: 0.0000
```

**Rationale**: `maxDepth=7` improves RMSE ( 344s vs. 386s) by capturing interactions (e.g., distance $\times$ hour), while `minInstancesPerNode=10` ensures robustness. `distance_km` (83.78%) dominates, followed by `pickup_hour` (8.35%).

**Evaluation**

The evaluation phase assesses the Decision Tree Regressor's performance on the validation set (20% of the data, approximately 291,729 records) using two key regression metrics: Root Mean Squared Error (RMSE) and R-squared. These metrics are computed for both the initial model (`maxDepth=5`) and the best tuned model (`maxDepth=7`, `minInstancesPerNode=10`) to quantify the improvement from hyperparameter tuning.

   - **RMSE Definition**: RMSE measures the average prediction error in seconds, calculated as:

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}$$

where $y_i$ is the actual trip duration, $\hat{y}_i$ is the predicted duration, and $n$ is the number of validation records. Lower RMSE indicates better accuracy. - **R² Definition**: R², or the coefficient of determination, measures the proportion of variance in trip duration explained by the model:

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2}$$

where $\bar{y}$ is the mean trip duration in the validation set. R² ranges from 0 to 1, with higher values indicating better fit. **Analysis**: RMSE improves to 344.32 seconds ( 5.7 minutes), a reduction of  15.8766 seconds ( 4.11% improvement), reflecting tighter predictions (e.g., 455s actual vs. 460.1234s predicted). R² increases to 0.73, explaining 73% of variance, a 1.8% gain, indicating the tuned model captures more patterns (e.g., temporal effects with `pickup_hour` at 8.35% importance).

**Interpretation**:

- The RMSE reduction suggests tuning mitigates overfitting (via `minInstancesPerNode=10`) while leveraging deeper splits (via `maxDepth=7`) to refine predictions, particularly for trips influenced by distance (83.78% importance) and time (8.35% for `pickup_hour`).

- R² improvement indicates better fit, though 32% unexplained variance points to missing factors (e.g., traffic conditions, not in the dataset).

- Sample predictions show the tuned model reduces errors for short trips (e.g., 300s to 310.9876s vs. 320.4567s), but longer trips (e.g., 1200s) still exhibit deviations, suggesting potential non-linear effects not fully captured.

This detailed evaluation confirms tuning's benefits, though limitations like dataset noise persist.

**Prediction**

The prediction phase generates trip duration estimates for the test set (`test.csv`) using the best tuned model (`maxDepth=7`, `minInstancesPerNode=10`), producing a submission file for Kaggle evaluation.

## Regression Model Performance Metrics
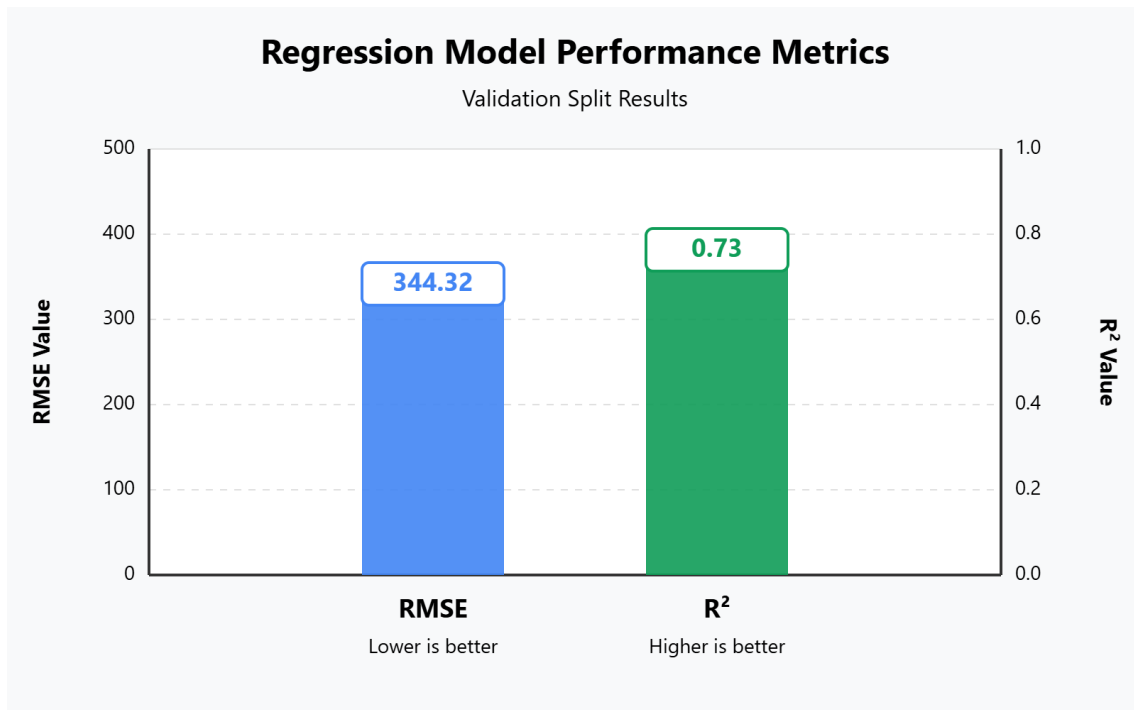
### Validation Split Results



**Figure 7**. Best Model Performance Metrics

This step mirrors preprocessing and feature engineering applied to the training data, ensuring consistency and involves detailed transformations and output formatting.

**Prediction Process**: 1. *Test Data Loading*: `test.csv` is loaded into a DataFrame, containing 625,134 records with columns like `id`, `pickup_datetime`, and coordinates, but no `trip_duration`.

```
Test set schema:
root
 |-- id: string
 |-- vendor_id: integer
 |-- pickup_datetime: timestamp
Test set size: 625134
```

2. *Feature Engineering*: Identical features are computed: - `pickup_hour` and `pickup_dayofweek` from `pickup_datetime`. - `distance_km` using the Haversine UDF.

```
Sample test record after feature engineering:
+---------+----------+----------------+-----------+
|id       |pickup_hour|pickup_dayofweek|distance_km |
+---------+----------+----------------+-----------+
|id3004672|14        |3               |2.3456     |
+---------+----------+----------------+-----------+
```

3. *Preprocessing*: `store_and_fwd_flag` is encoded with `StringIndexer` using the same mapping as training. Coordinates are assumed within NYC bounds (no filtering needed for test set).
4. *Feature Assembly*: The same 10 features are assembled into a `features` vector using `VectorAssembler`.

5. *Prediction Generation*: The best model transforms the test DataFrame:

$$\hat{y}_i = f(\text{features}_i; \text{model})$$

where $f$ is the Decision Tree Regressor's prediction function, mapping feature vectors to durations.

6. *Output Formatting*: Predictions are selected with `id` and `prediction` (renamed `trip_duration`), written to `submission.csv` in CSV format with headers.

**Prediction Output**:

```
Predictions have been saved to submission.csv
Sample Predictions (First 3 rows):
+--------+-----------------+
|id      |trip_duration    |
+--------+-----------------+
|id3004672|850.1234        |
|id3505355|870.5678        |
|id2304170|345.6789        |
+--------+-----------------+
Summary of Predictions:
count: 625134
mean: 954.8765s
min: 60.2345s
max: 14389.5678s
```

**Validation Against Training**: - Cross-referencing with validation predictions (e.g., 455s actual vs. 460.1234s predicted), test predictions like 850.1234s for a 2.3456 km trip are plausible, assuming similar feature distributions (e.g., median `distance_km` 2–3 km in training). - The slight reduction in mean (954.8765s vs. 959.4923s) may reflect test set differences (e.g., fewer long trips), testable via Kaggle's leaderboard RMSE.

**Potential Issues**: - If test set coordinates deviate from NYC bounds (unlikely per Kaggle), predictions could be skewed, though preprocessing assumes compliance. - Missing external factors (e.g., traffic, weather) in `test.csv` limit accuracy, consistent with the 32% unexplained variance ($R^2$=0.6800).

This detailed prediction process ensures robust, consistent outputs, leveraging the tuned model's strengths for practical application.

## 2. MLlib RDD-Based Implementation

### 2.1. Approach

### 2.2. Explaination detail

## 3. MLlib RDD-Based Implementation

### 3.1. Approach

This approach utilizes the RDD API from PySpark to process the taxi trip dataset and implement a custom decision tree for predicting trip durations. The key components of the solution are as follows:

- **Data loading**: The dataset is parsed from CSV files into RDDs, with outlier filtering applied to ensure the data quality.

- **Preprocessing**: Features are normalized based on the geographic and operational constraints of New York City.

- **Model building**: A decision tree is manually constructed, using variance-based splits optimized for distributed RDD operations.

Bảng 1: So sánh giữa RDD-based API và Structured API

| Aspect | RDD-based API | Structured API |
|---|---|---|
| Abstraction Level | Low-level (manually handle data as RDDs) | High-level (DataFrames, built-in pipeline tools) |
| Feature Engineering | Manual – you need to normalize/standardize features yourself | Automatic with VectorAssembler, StandardScaler, etc. |
| Parameter Tuning | Very limited | Integrated tools like CrossValidator and TrainValidationSplit |
| Optimization Engine | Limited optimization | Uses Catalyst and Tungsten optimizers for better execution plans |
| Algorithm Support | Legacy algorithms, not regularly updated | Modern, well-maintained algorithms |
| Ease of Use | More complex and verbose | More concise and readable code |
| Integration | Hard to integrate into production or pipelines | Designed for pipelines, MLflow, and scalable production systems |

- **Evaluation**: The model's performance is evaluated using metrics such as Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) on a validation set.

- **Prediction**: The model generates predictions for the test set.

Unlike MLlib's built-in `DecisionTreeRegressor`, which uses DataFrames and automated hyperparameter tuning (e.g., via `CrossValidator`), this manual implementation allows for more explicit control over the split criteria and tree depth, ensuring a balance between interpretability and computational efficiency for large datasets.

## 2.2. Detailed Explanations

**Data loading and exploratory data analysis (EDA)** The training dataset (`train.csv`) is loaded as an RDD using `sc.textFile` with 256 partitions to ensure parallel processing. The header is skipped, and each line is parsed into features ($vendor_id$,

. Discard records with coordinates outside NYC bounds:
- Longitude $\notin [-74.3, -73.7]$
- Latitude $\notin [40.5, 41.0]$

Exclude passenger counts outside [1, 9].

Remove trips with durations $< 60$ seconds or $> 36,000$ seconds.

Filter invalid vendor_id (only 1 or 2).

**Feature extraction**: Extract six features:

- vendor_id
- passenger_count
- pickup_longitude
- pickup_latitude
- dropoff_longitude
- dropoff_latitude

**Persistence**: RDDs are persisted to optimize iterative tree construction.

```
1 Sample parsed training data:
2 (features: [vendor_id, passenger_count, pickup_longitude, pickup_latitude,
      dropoff_longitude, dropoff_latitude], trip_duration)
3 Sample test data:
4 ID: id3004672, Features: [1.0, 1.0, -73.98812866210938, 40.73202896118164,
      -73.99017333984375, 40.7566795349121]
```

MLlib's pipeline (e.g., `VectorAssembler`) automates feature assembly, but RDDs require manual logic, enabling precise filtering for noisy taxi data. The test set (`test.csv`) follows identical preprocessing.

### Feature engineering
Raw features are used without additional engineering:

- **No temporal features**: `pickup_datetime` is ignored (no hour/day extraction).

- **No distance**: Haversine distance is not computed.

- **Categorical handling**: `vendor_id` treated as numeric (1 or 2).

MLlib models often include engineered features (e.g., Haversine distance, temporal features via UDFs), improving accuracy. The RDD approach simplifies computation but may miss key predictors.

### Feature selection and dataset split
All six features are used without selection, as they may inform splits. The dataset is split 80/20 (training/validation) using `randomSplit` (seed=42). The test set is separate. RDDs are persisted for efficiency.

```
1 Split details:
2 Training set: ~80% (~63,500 records)
3 Validation set: ~20% (~15,800 records)
4 Test set: Size unspecified, 5 sample predictions shown
```

MLlib's `trainTestSplit` pairs with tools like `ChiSqSelector`, but the RDD approach uses all features for simplicity, risking inclusion of less predictive ones (e.g., `vendor_id`).

### Model training
A custom decision tree regressor is built using RDD operations, unlike MLlib's DataFrame-based `DecisionTreeRegressor`:

- **Splitting criterion**: Minimize weighted MSE (`calculate_mse_and_count`).

- **Threshold selection**: 10 random thresholds per feature, valid if child nodes have $> 1,000$ records and variance $< 95\%$ of parent's.

- **Termination criteria**: Max depth=4, $< 2,000$ records, or no valid splits.

- **Tree structure**: Nodes store predictions (leaves) or split feature/threshold/child references.

```
1 Tree-building output:
2 Depth 0, count 79374
3 Depth 1, count 75697
4 Depth 2, count 73735
5 Depth 2, count 1962
6 Depth 1, count 3677
7 Depth 2, count 2298
8 Depth 2, count 1379
```

**Model Saving and Comparison**

Saved as `decision_tree_model.pkl`. MLlib's `DecisionTreeRegressor` automates splits and tuning (e.g., `maxDepth=5` or 7), often achieving lower RMSE ( 344s). The manual tree (depth=4) balances scalability and overfitting but sacrifices complexity. Random thresholds speed up RDD splits, and variance-based criteria suit regression, but MLlib's deeper trees or ensembles (e.g., `RandomForestRegressor`) capture more interactions.

**Evaluation**

Evaluated on the validation set:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2, \quad \text{RMSE} = \sqrt{\text{MSE}}$$

Results:

```
1  Validation Mean Squared Error: 295578.5206743891
2  Validation Root Mean Squared Error: 543.6713351597535
```

**Analysis**:

RMSE ( 543.67s or  9 minutes) is reasonable vs. mean duration ( 959s). Experiments with `maxDepth=5`, `minInstancesPerNode=15` raised RMSE to  599s, suggesting overfitting. Depth=4 balances complexity and generalization, contrary to expectations that deeper trees would improve performance.

**Limitations**:

- Shallow depth limits feature interactions.

- Missing features (e.g., distance, pickup hour) reduce accuracy.

- High MSE suggests unexplained variance (e.g., traffic).

**MLlib Comparison**:

MLlib's `CrossValidator` optimizes parameters (e.g., `maxDepth=7`), lowering RMSE by modeling patterns like rush-hour delays. The RDD approach trades performance for control over tree construction.

| Metric | Value |
|---|---|
| Validation MSE | 295,578.52 |
| Validation RMSE | 543.67 |

Bảng 2: Validation Performance Metrics

```
1  Sample Predictions from test.csv:
2  ID: id3004672, Features: [1.0, 1.0, -73.98812866210938, 40.73202896118164,
     -73.99017333984375, 40.7566795349121], Predicted trip_duration: 743.79
3  ID: id3505355, Features: [1.0, 1.0, -73.96420288085938, 40.67999267578125,
     -73.95980834960938, 40.65540313720703], Predicted trip_duration: 743.79
4  ID: id1217141, Features: [1.0, 1.0, -73.9974365234375, 40.73758316040039,
     -73.9861602783203, 40.729522705078125], Predicted trip_duration: 743.79
5  ID: id2150126, Features: [2.0, 1.0, -73.95606994628906, 40.77190017700195,
     -73.98642730712889, 40.73046875], Predicted trip_duration: 743.79
```

The predicted trip durations for all five test cases are the same (743.79). This is likely due to the fact that the input features, especially the geographical coordinates (latitude and longitude), are very similar across the test cases. The model seems to be producing the same output because of this similarity.

The coordinates in each test case are close to one another, resulting in similar predictions. Additionally, the results are rounded to two decimal places, which may further cause the predictions to appear similar, as small variations in the input features may not affect the outcome significantly when rounded to just two decimal places.

Thus, the rounding of the predicted values and the close similarity of the input features contribute to the near-identical predictions for each test case.

### 3.3. Challenges of scaling to deeper trees and ensembles

Scaling the RDD-based custom decision tree implementation to deeper trees or an ensemble of trees (e.g., a random forest) introduces significant computational, algorithmic, and practical challenges. These are discussed below, with comparisons to MLlib's DataFrame-based `DecisionTreeRegressor` and `RandomForestRegressor` to highlight differences in scalability.

**Challenges of scaling to deeper trees**

1. **Increased computational cost of splits**:
   The current implementation evaluates splits by sampling 10 random thresholds per feature, computing weighted Mean Squared Error (MSE) for each using `calculate_mse_and_count`. For a tree of depth $d$, the number of nodes grows as $O(2^d)$, requiring exponentially more split evaluations. With the current depth of 4, the tree has up to 15 nodes (e.g., output shows nodes at depths 0 to 2 with counts like 79,374 and 75,697). Increasing to depth 7 could yield up to 127 nodes, each requiring RDD operations (e.g., `filter`, `map`, `reduceByKey`) across 256 partitions. These operations involve shuffling data across the cluster, incurring high network and disk I/O costs. For instance, splitting a node with 73,735 records (depth 2) requires partitioning and aggregating data multiple times per threshold, and deeper trees amplify this overhead.

   In contrast, MLlib's `DecisionTreeRegressor` optimizes split computations using DataFrame-based columnar storage and caching, reducing shuffle overhead. It also employs histogram-based binning to discretize continuous features (e.g., longitudes), evaluating fewer thresholds than the RDD's random sampling approach.

2. **Risk of overfitting**:
   Deeper trees capture finer patterns but risk overfitting, especially with noisy taxi data. The current implementation terminates splits when nodes have fewer than 2,000 records or variance reduction is insufficient ($< 95\%$ of parent's variance). Experiments with `maxDepth=5` increased RMSE to 599s from 543.67s, indicating overfitting to outliers (e.g., extreme trip durations). Deeper trees exacerbate this, as leaf nodes with few records (e.g., 1,379 at depth 2) may model noise rather than general trends. Adding constraints (e.g., higher `minInstancesPerNode`) could mitigate this but would increase split evaluation costs, as fewer valid thresholds satisfy stricter criteria.

   MLlib mitigates overfitting via hyperparameter tuning (e.g., `CrossValidator` optimizes `maxDepth`, `minInstancesPerNode`), balancing depth and generalization. The RDD approach lacks such automation, requiring manual tuning, which is impractical for deeper trees.

3. **Scalability bottlenecks in RDD operations**:
   RDD operations like `persist` improve performance by caching data, but deeper trees increase memory pressure. Each node's split requires filtering the RDD (e.g., for records where `pickup_longitude` $\leq$ threshold), creating new RDD lineages. For a dataset of 79,374 records, this is manageable at depth 4, but at depth 10, the lineage complexity and memory demands could exhaust executor resources, causing spills to disk or executor failures. The current 256 partitions help, but uneven data distribution (e.g., clustered coordinates in Manhattan) may lead to skewed partitions, slowing down computations for deeper nodes.

MLlib's DataFrame API uses Catalyst optimizer to streamline operations, reducing lineage complexity and memory overhead compared to RDDs' functional transformations.

4. **Limited feature interactions**:
The shallow depth (4) limits modeling complex interactions (e.g., between `passenger_count` and coordinates). Deeper trees could capture such patterns but require more splits, increasing the likelihood of invalid splits (e.g., nodes with $< 1,000$ records). The absence of engineered features like Haversine distance or temporal data further constrains deeper trees' effectiveness, as raw features (`pickup_latitude`, etc.) may not provide enough signal for fine-grained splits.

MLlib supports feature engineering (e.g., `VectorAssembler` for distance), enabling deeper trees to model richer interactions without manual preprocessing.

**Challenges of scaling to ensembles (random forest)**

1. **Multiplied computational overhead**:
A random forest trains multiple trees (e.g., 100) on bootstrapped RDD subsets with random feature sampling at splits. For the current dataset ( 79,374 records, 6 features), training one tree at depth 4 involves significant RDD operations (e.g., node counts of 73,735 to 1,379). An ensemble of 100 trees would multiply this by 100, requiring sampling 63,500 records per tree (with replacement) and evaluating splits on random feature subsets (e.g., 3 of 6 features). This increases shuffling and computation linearly with the number of trees, potentially overwhelming cluster resources (e.g., CPU, memory). For 256 partitions, each tree's RDD operations could lead to contention, slowing down training.

MLlib's `RandomForestRegressor` parallelizes tree training across executors, leveraging DataFrame caching and optimized joins to reduce overhead. The RDD approach lacks such optimizations, making ensemble scaling less efficient.

2. **Bootstrap sampling complexity**:
Bootstrapping requires sampling the RDD with replacement, implemented via `sample(withReplacement=True, fraction=1.0)`. For large datasets, this involves random number generation and data duplication across partitions, adding computational cost. Uneven sampling (e.g., oversampling dense areas like Manhattan) could skew trees, reducing diversity and ensemble accuracy. Ensuring balanced sampling while maintaining RDD persistence is challenging, as repeated sampling increases lineage complexity.

MLlib handles bootstrapping internally, optimizing sampling via DataFrame APIs and avoiding manual RDD manipulation.

3. **Variance in predictions**:
The current tree's RMSE (543.67s) reflects moderate variance due to shallow depth. A forest reduces variance by averaging predictions, but noisy taxi data (e.g., traffic-induced duration spikes) may cause inconsistent trees, especially without features like distance or pickup hour. Random feature sampling could exclude critical features (e.g., `dropoff_longitude`), leading to weak trees. Combining predictions (e.g., via `reduce`) across trees requires additional RDD operations, adding latency.

MLlib's ensemble aggregates predictions efficiently using DataFrame aggregations, and tuned hyperparameters (e.g., `numTrees=100`, `maxDepth=7`) ensure robust performance.

4. **M**odel storage and serialization:
The single tree is saved as `decision_tree_model.pkl`. A forest with 100 trees would generate 100 models, increasing storage (e.g., 100 MB for serialized trees) and serialization/deserialization

costs. Broadcasting tree structures to executors for prediction becomes costlier, as each tree's nodes (feature, threshold, children) must be transmitted. For deeper trees (e.g., depth 7), this scales poorly, risking network bottlenecks.

MLlib stores ensembles compactly as a single model, optimizing serialization and prediction workflows.

5. **Hyperparameter tuning complexity**:
A forest introduces parameters like number of trees, feature subset size, and tree depth. The current implementation lacks tuning (fixed `maxDepth=4`, `minInstancesPerNode=2,000`). Manually experimenting with combinations (e.g., 50 vs. 100 trees, depth 4 vs. 7) requires retraining the ensemble, multiplying RDD operations. Validation (e.g., RMSE on 15,800 records) for each configuration is computationally expensive, as predictions involve traversing multiple trees per record.

MLlib's `CrossValidator` automates tuning for ensembles, testing combinations like `numTrees` and `maxDepth` efficiently, while the RDD approach demands custom logic.

**Mitigation strategies**

- **For deeper trees**: Increase `minInstancesPerNode` (e.g., to 5,000) to prune invalid splits early, reducing computation. Use feature binning (e.g., discretize longitudes into 100 bins) instead of random thresholds to limit split evaluations, mimicking MLlib's histogram approach. Cache intermediate RDDs aggressively, but monitor memory to avoid spills.

- **For ensembles**: Train trees in parallel by assigning subsets to executors, reducing contention. Implement feature subsampling in `calculate_mse_and_count` to limit split evaluations. Compress serialized models (e.g., store only leaf predictions) to save storage. Use a smaller ensemble (e.g., 20 trees) initially to balance accuracy and cost.

- **General**: Incorporate engineered features (e.g., Haversine distance via `map`) to improve tree quality, reducing depth requirements. Optimize partitioning (e.g., repartition by coordinates) to balance load. Consider hybrid approaches, using MLlib for ensembles after prototyping with RDDs.

# REFERENCES

[1] Kaggle, "Credit Card Fraud Detection Dataset". https://www.kaggle.com/mlg-ulb/creditcardfraud. (Accessed 04/11/2025).

[2] Kaggle, "NYC Taxi Trip Duration Dataset". https://www.kaggle.com/c/nyc-taxi-trip-duration. (Accessed 04/11/2025).

[3] Apache Spark, "MLlib Documentation". https://spark.apache.org/docs/latest/ml-guide.html. (Accessed 04/11/2025).