

Department of Computer Science

Missouri State University

CSC735 - Data Analytics

**Predicting NYC Yellow Taxi Fares Using Big Data
Analytics**

Project Report

Submitted by

Section: 01 Group number: 32

Shadman Sakib ID: M03543624

Khalid Hasan ID: M03543550

Submitted on December 23, 2023

1 Abstract

Due to the increase in ride-sharing app usage in New York City, the demand for yellow taxi cabs has been on the downside. This is due to the insightful information provided to customers by the app regarding the expected fare and time to reach their destination. To compete in this aspect we wish to provide a solution to this problem by building a model that can predict the fare of yellow taxis in New York City. In our experiment, we have used Apache Spark as a cluster computer to train our regression models. We have used Linear Regression, Decision Tree, and Random Forest models to predict the fare. From our evaluation of these three models, we have found Random Forest to be the best-performing.

2 Introduction

The heart of traffic in New York City revolves around taxi rides, which play a central role in the daily commute for many New Yorkers. These rides provide valuable insights into traffic patterns, potential roadblocks, and other relevant factors. Accurately predicting the duration of a taxi journey holds significant importance, as users consistently seek precise information on the time required to travel between locations. With the increasing prevalence of app-based taxi services offered by popular providers such as Ola and Uber, maintaining competitive pricing is essential to attract users and ensure their preference for these platforms. Anticipating the duration and cost of trips is beneficial for users in effectively planning their journeys. This information also assists drivers in selecting optimal routes, thereby reducing travel time. In this research study, real-time data provided by customers at the commencement or booking of a ride was utilized to predict fare. In our experiment, we have used an open-source distributed computing system that provides a fast and general-purpose cluster-computing framework for big data processing: Apache Spark, to predict the continuous value of taxi fare based on a handful of features from our dataset. We have trained and tested three different machine learning models to predict and evaluate them. The models used are Linear Regression, Decision Tree, and Random Forest. From our evaluation, we have found Random Forest outperforms the other 2 models by a small margin.

3 Objective

We wish to propose a fare prediction based on the location of pick-up and destination and other key factors such as the rate code and distance. We will be utilizing the powerful and popular open-source big data processing framework, Apache Spark. In this research, we have developed a fare prediction model for the New York City Yellow Taxi rides. We have used a few regression techniques to predict the continuous value of price based on a handful of features we have extracted from the raw data itself.

4 Problem Statement

Due to the increasing number of ride-sharing apps these days, taxi businesses are at a disadvantage. Ride-sharing apps often offer lower prices than traditional taxis due to lower overhead costs. This makes them more attractive to cost-conscious passengers, which can lead to a decline in taxi ridership. Ride-sharing apps use surge pricing during high-demand periods, which can incentivize more drivers to be on the road during those times. In contrast, traditional taxis typically charge fixed rates regardless of demand. Ride-sharing apps have leveraged technology to enhance the overall user experience, from seamless payment processing to predictive ride recommendations. If customers were transparent as to how much they would be charged beforehand maybe taxis would be in demand too. Ensuring price predictability is crucial, as it facilitates improved cost management, mitigates instances of unfair pricing, and empowers customers with the information needed to compare prices with those of competitors.

5 Literature Review

In a comparative study, the effectiveness of gradient boosting using XGBoost was evaluated against a deep learning technique known as multi-layer perceptron (MLP) for predicting trip duration [1]. The findings revealed that XGBoost, demonstrated superior performance over the MLP model when accounting for all variables. Nonetheless, the authors observed that the MLP model could potentially be enhanced through auto-tuning, albeit with the trade-off of requiring additional time.

6 Data Set

The trip data from the New York City Taxi and Limousine Commission, which includes observations on around 1 billion taxi journeys in New York City between 2009 and 2016, is the source of all the data used in this study. The data for yellow taxi rides in January 2020 were used for the primary analyses in this study. A random subset of 640,508 observations—of which 80% are used for training and 20% for testing—was used to construct the models. The dataset consists of 19 columns from which we have decided to use only those that show a high Correlation to taxi fare. We have decided to use the trip distance, pick-up point, destination, and rate code ID for different sections within New York City.

Table						
	summary ▲	trip_distance ▲	PULocationID ▲	DOLocationID ▲	RateCodeID ▲	fare_amount ▲
1	count	6405008	6405008	6405008	6339567	6405008
2	mean	2.929643933309735	164.73225778952968	162.6626908194338	1.0599077192495954	12.694108119770615
3	stddev	83.1591059732502	65.54373944111758	69.91260629496094	0.811843207190649	12.127295340046553
4	min	-30.62	1.0	1.0	1.0	-1238.0
5	max	210240.07	265.0	265.0	99.0	4265.0
5 rows						

Figure 1: Summary of features selected

7 Methodology

To carry out the project we have followed several data analytics steps including data cleaning, data pre-processing, model definition, model tuning, and data predictions. We have used Apache Spark's computing system which provides a fast and general-purpose cluster-computing framework for big data processing. Spark's ability to scale horizontally across a cluster of machines enables handling large datasets seamlessly, making it suitable for big data machine learning projects.

7.1 Data Preprocessing

Clean and well-processed data helps in building more accurate models. Removing inconsistencies, errors, or outliers ensures that the model is trained on reliable and representative data. In this step, we have removed invalid and redundant data generated due to error.

1. We have removed rows that have a negative trip distance and fare amount.
2. We have removed any rows that have null in the RateCode ID.
3. We have removed outliers using the IQR method where any values lower than or greater than 1.5 times the interquartile range below and above the first and third quartiles respectively are removed. You can see in figure 2 how we have reduced and normalized our data to a much more linear scale after removing the outliers.
4. We have split the data into train and test data sets in the ratio of 80:20.

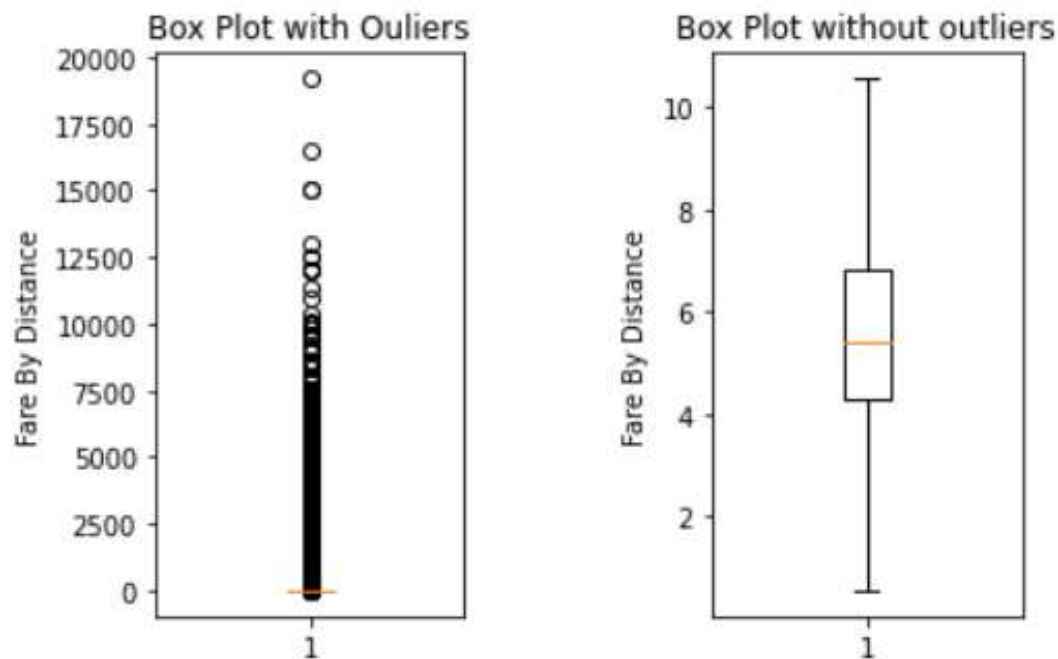


Figure 2: Outliers

7.2 Model Definition

We have applied the relevant regression models to the processed data to get the prediction. We have used three different models to compare and contrast among them. We have leveraged Linear, Decision Tree, and Random Forest Tree Regression models. Each model has been pushed through a pipeline to streamline and automate the process of building, training, and deploying models. Each pipeline includes three stages: an Assembler, a Standard scalar, and a Regressor as shown in figure 3.

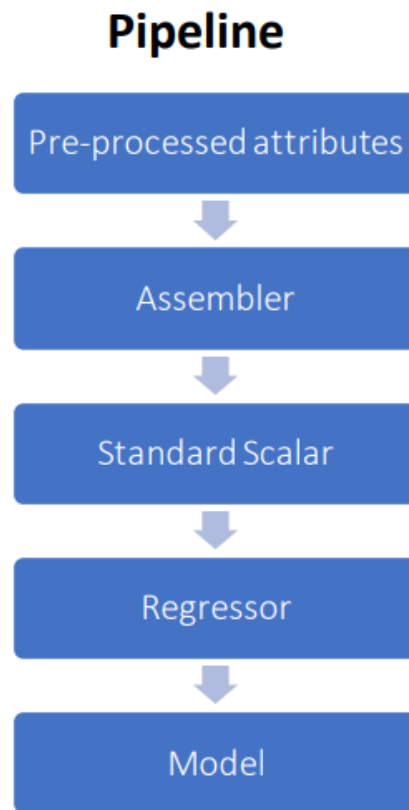


Figure 3: Pipeline

7.3 Model Tuning

To get the optimized model, we have tuned the models with different sets of hyperparameters. Spark's ParmGridBuilder is used to construct parameter grids for hyperparameter tuning. Figure 4 demonstrates the hyperparameters we have tuned to get the best model.

Linear Regressor	Decision Tree Regressor	Random Forest Regressor
<ul style="list-style-type: none"> ElasticNetParam: [0.0, 0.5, 1.0] RegParam: [0.1, 0.5, 0.9] MaxIter: [50, 100, 200] 	<ul style="list-style-type: none"> MaxDepth: [3, 5, 10] MinInfoGain: [0, 0.5] MinInstancesPerNode: [1, 2, 5, 10] 	<ul style="list-style-type: none"> NumTrees: [10, 20, 30, 40] MaxDepth: [5, 10]

Figure 4: hyperparamters for 3 models

7.4 Model Training

Used Spark's TrainValidationSplit model to split our training dataset into a training set and a validation set, with a train ratio of 0.75. We set up one TrainValidationSplit model

for each regressor with the pipeline, the parameter grid, and the evaluation process. Then we fitted the training data to the TrainValidationSplit model. After finishing the model training stage with the different combinations of hyperparameters, we find our best-trained models with optimized hyperparameters as shown in figure 5.

Model	Training Time	Parameters and hyper-parameters
Linear Regression	27.54 minutes	numIterations: 7 Co-efficients: [2.723277, 4.589777E-4, -2.8002004E-4, 0.825135] Intercept: 3.537987
Decision Tree	20.93 minutes	depth=10 numNodes=1663
Random Forest Tree	42.18 minutes	numTrees=30

Figure 5: results of training

7.5 Performance Evaluation

To evaluate the performance of our models we need proper evaluation metrics to find out the accuracy of the predictions and understand how close they are to the actual values. In our case, we have used 3 different metrics for a systematic comparison. We have used root mean squared error (RMSE), R-squared, and mean absolute error (MAE). The RMSE formula is given by:

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

The R-squared formula is given by:

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

The MAE formula is given by:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

8 Experiments and Results

We conducted an experimental study to evaluate our regression models, particularly the effects of corresponding features and regression techniques on the prediction of taxi fare. Our experiment consists of two analyses where each analysis visualizes the whole scenario of fare prediction based on a given test data set. These analyses are (1) Actual vs predicted fare analysis and (2) Residual analysis. We analyzed the prediction of a testing data set given by our defined three regression models: (1) Linear regression, (2)

Decision tree regression, and (3) Random forest tree regression.

8.1 Actual vs Predicted Fare Analysis

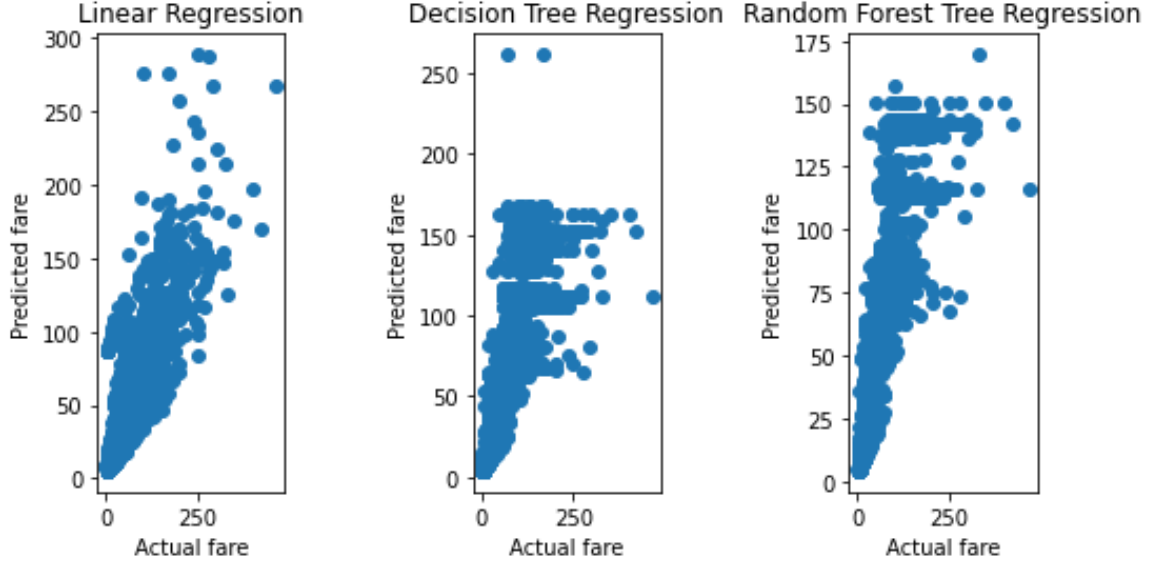


Figure 6: Actual vs Predicted Fare Analysis

In a regression analysis, an “actual vs predicted fare” analysis involves comparing the predicted values generated by a regression model to the actual observed values. This type of analysis helps evaluate how well the model performs in predicting the target variable (fare, in this case) based on input features. Each point on the scatter plot represents a data point from the test set, and the closer the points are to a diagonal line, the better the model’s predictions match the actual values.

Figure 6 shows three scatter plots for each regression model where the x-axis represents the actual fares, and the y-axis represents the predicted fares. The plot for linear regression shows that the points deviate more than other models. For example, the model predicts above 250 for a considerable amount of less than 200. In contrast, we see the decision tree regression model deviates less and we get a reasonable amount of predicted value above 150 for 200 to 250 actual values. However, the random forest tree shows the most excellence. We can easily find out that this model presents a visible diagonal line and most of the points center that line which means the points do not deviate from the diagonal line. This analysis shows the random forest tree regression works better in the case of the taxi fare data set.

8.2 Residual Analysis

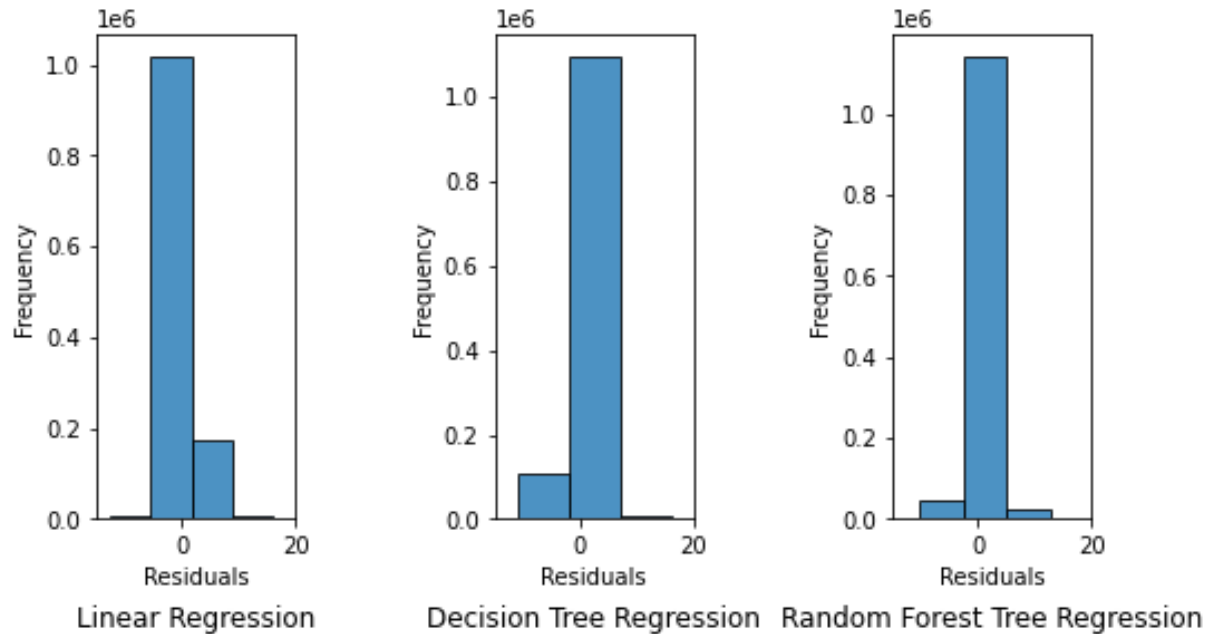


Figure 7: Residual analysis Analysis

Residual analysis is an essential step in evaluating the performance of a regression model. Residuals are the differences between the observed values (actual outcomes) and the predicted values generated by the regression model. Analyzing residuals helps to understand how well the model fits the data and identify areas for improvement. The histogram provides insights into the distribution of residuals. A normal distribution suggests that the model is performing well.

Figure 8 shows the histogram of residuals for the prediction values of each regression model. We can see four blocks in the case of the linear regression model which implies a greater standard deviation compared to other models. Furthermore, the frequency of the residuals at 0-mean is around 10^6 whereas the decision tree model has above 10^6 and the random forest tree regression model has a greater frequency at 0-mean with a smaller standard deviation. Moreover, the linear regression model is a bit right-skewed, the decision tree model is a bit left-skewed while the random forest model follows the normal distribution of models. These pieces of information suggest that the random forest tree regression model fits the given taxi fare data set better compared to other regression models.

9 Evaluation

Regression Model ▲	RMSE ▲	MAE ▲	R2 ▲
Linear	2.898132	1.565766	0.929683
Decision Tree	2.416508	1.287255	0.951112
Random Forest Tree	2.399815	1.303674	0.951785

Figure 8: Model Evaluation

We have used the common metrics Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R-squared (R2) which are used to evaluate the performance of regression models. RMSE is a metric that measures the average magnitude of the residuals (the differences between predicted and actual values), emphasizing larger errors. It is calculated as the square root of the mean squared error (MSE). Moreover, MAE is a metric that measures the average absolute difference between predicted and actual values, giving equal weight to all errors. It is less sensitive to outliers compared to RMSE. On top of that, R-squared is a metric that represents the proportion of the variance in the dependent variable (target) that is explained by the independent variables (features) in the model. It ranges from 0 to 1, where 0 indicates that the model does not explain any variability, and 1 indicates perfect prediction.

In the case of RMSE evaluation, the linear regression model gives more error compared to other models whereas the random forest tree model performs better. Furthermore, the decision tree regression model shows better performance in the scale of MAE followed by the random forest tree model. In addition to that, the random forest tree model has the maximum R2 score in comparison to others which has made it the best regression model. Thus, after comparing the evaluation metrics of our defined regression models, we can realize that each regression model shows excellence but the Random Forest Tree emerges as the best.

10 Conclusion

Here we have seen how to use Apache Spark and leverage its ability to handle big data with 3 different regression models. Overall we have found good results. In the future, we can compile all the recent data of these trips of the yellow taxi in New York, not just for a single month, and analyze trends and seasonality of the fare concerning the different features recorded. We can also extend the work by using neural networks or deep learning techniques.

References

- [1] Poongodi, M., Malviya, M., Kumar, C. et al. New York City taxi trip duration prediction using MLP and XGBoost. Int J Syst Assur Eng Manag 13 (Suppl 1), 16–27 (2022). <https://doi.org/10.1007/s13198-021-01130-x> [Accessed 27/12/2022]
-

Appendices

A Code for the Implementation

```
1 // Databricks notebook source
2 // MAGIC %md
3 // MAGIC
4 // MAGIC ## Overview
5 // MAGIC
6 // MAGIC This notebook is generated for Data Analytics project. The project is to
   predict fare of yellow taxi in NewYork. Various regression model are used to
   build a model that is trained on the dataset received from Kaggle. The dataset
   is specically for January 2020 named yellow_tripdata_2020_01.csv. It is
   approximately 600 MB in size.
7 // MAGIC
8 // MAGIC Due to the big data size, this project has been conducted using Spark.
   Scala is mainly used as a programming language for data access and manipulating
   . Python is used particularly for plot graphs.
9
10
11
12 import org.apache.spark.sql.functions.col
13
14 // # File location and type
15 val file_location = "/FileStore/tables/yellow_tripdata_2020_01.csv"
16 val file_type = "csv"
17
18 // # CSV options
19 val infer_schema = "false"
20 val first_row_is_header = "true"
21 val delimiter = ","
22
23 // # The applied options are for CSV files. For other file types, these will be
   ignored.
24 val df = spark.read.format(file_type)
25   .option("inferSchema", infer_schema)
26   .option("header", first_row_is_header)
27   .option("sep", delimiter)
28   .load(file_location)
29
30 df.cache()
31
32 var selectedData = df.select(
33   col("trip_distance").cast("double"),
```

```

34     col("PULocationID").cast("double"),
35     col("DOLocationID").cast("double"),
36     col("RateCodeID").cast("double"),
37     col("fare_amount").cast("double")
38 )
39
40 display(selectedData.describe())
41 println(selectedData.count())
42
43
44
45 // Count noisy data
46 println(selectedData.where(col("trip_distance")<=0 || col("fare_amount")<=0 || col
    ("RateCodeID").isNull).count)
47
48 // Removing noisy data observed from the statistics
49 selectedData = selectedData.selectExpr("*, "fare_amount / trip_distance as
    fareByDistance")
50                                     .where(col("trip_distance")>0 && col("fare_amount")>0 &&
    col("RateCodeID").isNotNull)
51
52 // display(selectedData.orderBy('trip_distance.asc))
53 println(selectedData.count)
54
55 // Create a view or table
56 selectedData.createOrReplaceTempView("selected_data")
57
58
59
60 // MAGIC %python
61 // MAGIC
62 // MAGIC pyData = spark.sql(
63 // MAGIC     """
64 // MAGIC     select * from selected_data
65 // MAGIC     """
66 // MAGIC )
67
68
69
70 // MAGIC %python
71 // MAGIC
72 // MAGIC # Use Python to create plots
73 // MAGIC import matplotlib.pyplot as plt
74 // MAGIC
75 // MAGIC # Data
76 // MAGIC y_values = pyData.select("fareByDistance").rdd.flatMap(lambda x: x).
    collect()
77 // MAGIC
78 // MAGIC # Create a box plot with all values
79 // MAGIC plt.subplot(1, 2, 1)
80 // MAGIC plt.subplots_adjust(wspace=0.8)
81 // MAGIC plt.boxplot(y_values)
82 // MAGIC plt.title("Box Plot with Outliers")
83 // MAGIC plt.ylabel("Fare By Distance")
84 // MAGIC
85 // MAGIC # Create a box plot without outliers
86 // MAGIC plt.subplot(1, 2, 2)
87 // MAGIC plt.boxplot(y_values, 0, '')
88 // MAGIC plt.title("Box Plot without outliers")

```

```

89 // MAGIC plt.ylabel("Fare By Distance")
90 // MAGIC
91 // MAGIC # Show the plot
92 // MAGIC plt.show()
93
94
95
96 import org.apache.spark.sql.DataFrame
97 import org.apache.spark.sql.functions._
98
99 def findOutliers(df: DataFrame, columns: Array[String]): DataFrame = {
100   // Identifying the numerical columns in a Spark DataFrame
101   // val numericColumns = df.dtypes.filter(_._2 == "Integer").map(_._1)
102   val numericColumns = columns
103
104   // Define the UDF to check if a value is an outlier
105   val isOutlier = udf((value: Int, Q1: Double, Q3: Double) => {
106     val IQR = Q3 - Q1
107     val lowerThreshold = Q1 - 1.5 * IQR
108     val upperThreshold = Q3 + 1.5 * IQR
109     if (value < lowerThreshold || value > upperThreshold) 1 else 0
110   })
111
112   var updatedDF = df
113
114   // Using the 'for' loop to create new columns by identifying the outliers for
115   // each feature
116   for (column <- numericColumns) {
117     val Q1 = updatedDF.stat.approxQuantile(column, Array(0.25), 0)(0)
118     val Q3 = updatedDF.stat.approxQuantile(column, Array(0.75), 0)(0)
119
120     val isOutlierCol = s"is_outlier_$column"
121
122     updatedDF = updatedDF
123       .withColumn(isOutlierCol, isOutlier(col(column), lit(Q1), lit(Q3)))
124   }
125
126   // Selecting the specific columns which we have added above
127   val selectedColumns = updatedDF.columns.filter(_.startsWith("is_outlier"))
128
129   // Adding all the outlier columns into a new column "total_outliers"
130   updatedDF = updatedDF.withColumn("total_outliers", selectedColumns.map(col).
131     reduce(_ + _))
132
133   // Dropping the extra columns created above
134   updatedDF = updatedDF.drop(selectedColumns: _*)
135
136   updatedDF
137 }
138
139 // Usage:
140 // val outliersDF = findOutliers(yourDataFrame)
141 // outliersDF.show()
142
143 // remove outliers
144 val outliersDF: DataFrame = findOutliers(selectedData, columns=Array("
  fareByDistance"))

```

```

145 println(outliersDF.where(col("total_outliers") > 0).count)
146
147 selectedData = outliersDF.where(col("total_outliers") === 0)
148 println(selectedData.count())
149
150
151
152 import org.apache.spark.ml.regression._
153 import org.apache.spark.ml.feature.{VectorAssembler, StandardScaler}
154 import org.apache.spark.sql.DataFrame
155 import org.apache.spark.ml.{Pipeline, PipelineModel}
156
157 // Step 4: Split the DataFrame into a training set and a test set (80% train, 20%
158 // test)
159 val Array(trainData, testData) = selectedData.randomSplit(Array(0.8, 0.2))
160
161 // val rFormula = new RFormula()
162
163 // Step 5: Prepare features and labels using a VectorAssembler
164 val assembler = new VectorAssembler()
165   .setInputCols(Array("trip_distance", "PULocationID", "DOLocationID", "RateCodeID"
166     " "))
167   .setOutputCol("features")
168   .setHandleInvalid("skip")
169
170 val sScalar = new StandardScaler().setInputCol("features")
171
172 // val regressionClasses = List("LinearRegression", "DecisionTreeRegressor", "
173 // RandomForestRegressor", "GBTRRegressor")
174
175 // Step 6: Create a Linear Regression model
176 val lR = new LinearRegression()
177   .setLabelCol("fare_amount")
178   .setFeaturesCol("features")
179   .setPredictionCol("prediction")
180
181 val dTR = new DecisionTreeRegressor()
182   .setLabelCol("fare_amount")
183   .setFeaturesCol("features")
184   .setPredictionCol("prediction")
185
186 val rFR = new RandomForestRegressor()
187   .setLabelCol("fare_amount")
188   .setFeaturesCol("features")
189   .setPredictionCol("prediction")
190
191 println(lR.explainParams())
192 println(dTR.explainParams())
193 println(rFR.explainParams())
194
195 val pipelineLR = new Pipeline().setStages(Array(assembler, sScalar, lR))
196 val pipelineDTR = new Pipeline().setStages(Array(assembler, sScalar, dTR))
197 val pipelineRFR = new Pipeline().setStages(Array(assembler, sScalar, rFR))
198
199
200 // specifying different combinations of hyperparameters to select the best model

```

```

    using an Evaluator, testing their predictions
201 import org.apache.spark.ml.tuning.ParamGridBuilder
202 import org.apache.spark.ml.evaluation.RegressionEvaluator
203
204 // Params for linear regressor
205 val paramsLR = new ParamGridBuilder()
206   .addGrid(lR.elasticNetParam, Array(0.0, 0.5, 1.0))
207   .addGrid(lR.regParam, Array(0.1, 0.5, 0.9))
208   .addGrid(lR.maxIter, Array(50, 100, 200))
209   .build()
210
211 // Params for Decision Tree regressors
212 val paramsDTR = new ParamGridBuilder()
213   .addGrid(dTR.maxDepth, Array(3, 5, 10))
214   .addGrid(dTR.minInfoGain, Array(0, 0.5))
215   .addGrid(dTR.minInstancesPerNode, Array(1, 2, 5, 10))
216   .build()
217
218 // Params for Random forests regressors
219 val paramsRFR = new ParamGridBuilder()
220   .addGrid(rFR.numTrees, Array(10, 20, 30, 40))
221   .addGrid(rFR.maxDepth, Array(5, 10))
222   .build()
223
224
225
226 // Evaluator using RegressionEvaluator using rmse
227 val rmse_evaluator = new RegressionEvaluator()
228   .setLabelCol("fare_amount")
229   .setPredictionCol("prediction")
230   .setMetricName("rmse")
231
232 // Evaluator using RegressionEvaluator using mae
233 val mae_evaluator = new RegressionEvaluator()
234   .setLabelCol("fare_amount")
235   .setPredictionCol("prediction")
236   .setMetricName("mae")
237
238 // Evaluator using RegressionEvaluator using r2
239 val r2_evaluator = new RegressionEvaluator()
240   .setLabelCol("fare_amount")
241   .setPredictionCol("prediction")
242   .setMetricName("r2")
243
244
245
246 // Define Train Validation Split
247 import org.apache.spark.ml.tuning.TrainValidationSplit
248
249 // Linear regressor
250 val tVS_LR = new TrainValidationSplit()
251   .setTrainRatio(0.75) // also the default.
252   .setEstimatorParamMaps(paramsLR).setEstimator(pipelineLR)
253   .setEvaluator(rmse_evaluator)
254
255 // Decision Tree regressor
256 val tVS_DTR = new TrainValidationSplit()
257   .setTrainRatio(0.75) // also the default.
258   .setEstimatorParamMaps(paramsDTR).setEstimator(pipelineDTR)

```

```

259     .setEvaluator(rmse_evaluator)
260
261 // Random Forests regressor
262 val tVS_RFR = new TrainValidationSplit()
263     .setTrainRatio(0.75) // also the default.
264     .setEstimatorParamMaps(paramsRFR).setEstimator(pipelineRFR)
265     .setEvaluator(rmse_evaluator)
266
267
268 // Get TrainValidationSplitModel for linear regressor
269 val tVS_LR_Model = tVS_LR.fit(trainData)
270
271
272
273 // Get TrainValidationSplitModel for decision tree regressor
274 val tVS_DTR_Model = tVS_DTR.fit(trainData)
275
276
277
278 // Get TrainValidationSplitModel for random forests regressor
279 val tVS_RFR_Model = tVS_RFR.fit(trainData)
280
281
282
283 // Get best model and statistics
284 import org.apache.spark.ml.regression.{LinearRegressionModel,
    DecisionTreeRegressionModel, RandomForestRegressionModel, GBTRegressionModel}
285
286 val bestPipelineLRModel = tVS_LR_Model.bestModel.asInstanceOf[PipelineModel]
287 val bestLRModel = bestPipelineLRModel.stages.last.asInstanceOf[
    LinearRegressionModel]
288
289 val summary = bestLRModel.summary
290 println(summary)
291 println(s"numIterations: ${summary.totalIterations}")
292 println(s"Co-efficients: ${bestLRModel.coefficients} Intercept: ${bestLRModel.
    intercept}")
293 summary.residuals.show()
294 println(summary.objectiveHistory.toSeq.toDF.show())
295 println(summary.objectiveHistory)
296 println(summary.rootMeanSquaredError)
297 println(summary.r2)
298
299 val bestPipelineDTRModel = tVS_DTR_Model.bestModel.asInstanceOf[PipelineModel]
300 val bestDTRModel = bestPipelineDTRModel.stages.last.asInstanceOf[
    DecisionTreeRegressionModel]
301
302 val summaryLR = bestLRModel.summary
303 println(summaryLR)
304
305 val bestPipelineRFRModel = tVS_RFR_Model.bestModel.asInstanceOf[PipelineModel]
306 val bestRFRModel = bestPipelineRFRModel.stages.last.asInstanceOf[
    RandomForestRegressionModel]
307
308
309
310 // Step 7: Make predictions on the test data
311 // val testPreprocessed = assembler.transform(testData)
312 val predictionsLR = bestPipelineLRModel.transform(testData)

```



```

313 predictionsLR.show(false)
314
315 val predictionsDTR = bestPipelineDTRModel.transform(testData)
316 predictionsDTR.show()
317
318 val predictionsRFR = bestPipelineRFRModel.transform(testData)
319 predictionsRFR.show()
320
321
322
323 // RMSE
324 val rmseLR = rmse_evaluator.evaluate(predictionsLR)
325 println(s"Root Mean Squared Error (RMSE) for LR: $rmseLR")
326
327 val rmseDTR = rmse_evaluator.evaluate(predictionsDTR)
328 println(s"Root Mean Squared Error (RMSE) for DTR: $rmseDTR")
329
330 val rmseRFR = rmse_evaluator.evaluate(predictionsRFR)
331 println(s"Root Mean Squared Error (RMSE) for RFR: $rmseRFR")
332
333
334 // MAE
335 val maeLR = mae_evaluator.evaluate(predictionsLR)
336 println(s"Mean Absolute Error (MAE) for LR: $maeLR")
337
338 val maeDTR = mae_evaluator.evaluate(predictionsDTR)
339 println(s"Mean Absolute Error (MAE) for DTR: $maeDTR")
340
341 val maeRFR = mae_evaluator.evaluate(predictionsRFR)
342 println(s"Mean Absolute Error (MAE) for RFR: $maeRFR")
343
344
345 // R2
346 val r2LR = r2_evaluator.evaluate(predictionsLR)
347 println(s"R-squared (r2) Error for LR: $r2LR")
348
349 val r2DTR = r2_evaluator.evaluate(predictionsDTR)
350 println(s"R-squared (r2) Error for DTR: $r2DTR")
351
352 val r2RFR = r2_evaluator.evaluate(predictionsRFR)
353 println(s"R-squared (r2) Error for RFR: $r2RFR")
354
355
356
357 import org.apache.spark.sql.{SparkSession, Row}
358 import org.apache.spark.sql.types._
359
360 val evaluatedDataMap = Seq(
361   Map("Regression Model" -> "Linear", "RMSE" -> rmseLR, "MAE" -> maeLR, "R2" ->
     r2LR),
362   Map("Regression Model" -> "Decision Tree", "RMSE" -> rmseDTR, "MAE" -> maeDTR, "
     R2" -> r2DTR),
363   Map("Regression Model" -> "Random Forest Tree", "RMSE" -> rmseRFR, "MAE" ->
     maeRFR, "R2" -> r2RFR)
364 )
365
366 // Define the schema based on the keys and types of the first map
367 val schema = new StructType()
368   .add("Regression Model", StringType)

```

```

369     .add("RMSE", DoubleType)
370     .add("MAE", DoubleType)
371     .add("R2", DoubleType)
372
373 // Convert the sequence of maps to a sequence of Rows
374 val rows = evaluatedDataMap.map { rowMap =>
375     Row.fromSeq(schema.map(field => rowMap.getOrElse(field.name, null)))
376 }
377
378 // Create a DataFrame
379 val evaluationMatrix = spark.createDataFrame(spark.sparkContext.parallelize(rows),
380     schema)
381
382 display(evaluationMatrix.select(col("Regression Model"), round('RMSE, 6).as("RMSE
383     "), round('MAE, 6).alias("MAE"), round('R2, 6).alias("R2")))
384
385 // create sql table view to access data while using python
386 predictionsLR.createOrReplaceTempView("predictions_lr")
387 predictionsDTR.createOrReplaceTempView("predictions_dtr")
388 predictionsRFR.createOrReplaceTempView("predictions_rfr")
389
390
391
392 // MAGIC %python
393 // MAGIC
394 // MAGIC import numpy as np
395 // MAGIC
396 // MAGIC # Create python dataframe
397 // MAGIC
398 // MAGIC predictionsLR = spark.sql(
399 // MAGIC     """
400 // MAGIC     select * from predictions_lr
401 // MAGIC     """
402 // MAGIC )
403 // MAGIC
404 // MAGIC predictionsDTR = spark.sql(
405 // MAGIC     """
406 // MAGIC     select * from predictions_dtr
407 // MAGIC     """
408 // MAGIC )
409 // MAGIC
410 // MAGIC predictionsRFR = spark.sql(
411 // MAGIC     """
412 // MAGIC     select * from predictions_rfr
413 // MAGIC     """
414 // MAGIC )
415 // MAGIC
416 // MAGIC
417 // MAGIC # Data Preparation of actual and predicted fares
418 // MAGIC models = {
419 // MAGIC     "lr": "Linear Regression",
420 // MAGIC     "dtr": "Decision Tree Regression",
421 // MAGIC     "rfr": "Random Forest Tree Regression"
422 // MAGIC }
423 // MAGIC
424 // MAGIC actual_fares = {
425 // MAGIC     "lr": np.array(predictionsLR.select("fare_amount").rdd.flatMap(lambda

```

```

        x: x).collect()),
426 // MAGIC      "dtr": np.array(predictionsDTR.select("fare_amount").rdd.flatMap(
        lambda x: x).collect()),
427 // MAGIC      "rfr": np.array(predictionsRFR.select("fare_amount").rdd.flatMap(
        lambda x: x).collect()),
428 // MAGIC }
429 // MAGIC predicted_fares = {
430 // MAGIC      "lr": np.array(predictionsLR.select("prediction").rdd.flatMap(lambda
        x: x).collect()),
431 // MAGIC      "dtr": np.array(predictionsDTR.select("prediction").rdd.flatMap(
        lambda x: x).collect()),
432 // MAGIC      "rfr": np.array(predictionsRFR.select("prediction").rdd.flatMap(
        lambda x: x).collect())
433 // MAGIC }
434 // MAGIC residuals = {
435 // MAGIC      "lr": actual_fares["lr"] - predicted_fares["lr"],
436 // MAGIC      "dtr": actual_fares["dtr"] - predicted_fares["dtr"],
437 // MAGIC      "rfr": actual_fares["rfr"] - predicted_fares["rfr"]
438 // MAGIC }
439
440
441
442 // MAGIC %python
443 // MAGIC # Use Python to create plots
444 // MAGIC import matplotlib.pyplot as plt
445 // MAGIC
446 // MAGIC # Create scatter plots
447 // MAGIC # Increase total size by setting figsize
448 // MAGIC plt.figure(figsize=(8, 4))
449 // MAGIC for i, k in enumerate(actual_fares.keys()):
450 // MAGIC     plt.subplot(1, 3, i+1)
451 // MAGIC     plt.subplots_adjust(wspace=1)
452 // MAGIC     plt.scatter(actual_fares[k], predicted_fares[k])
453 // MAGIC     plt.xlabel("Actual fare")
454 // MAGIC     plt.ylabel("Predicted fare")
455 // MAGIC     plt.title(models[k])
456 // MAGIC plt.show()
457 // MAGIC
458 // MAGIC # Create histogram of residuals
459 // MAGIC # Increase total size by setting figsize
460 // MAGIC plt.figure(figsize=(8, 4))
461 // MAGIC for i, k in enumerate(residuals.keys()):
462 // MAGIC     plt.subplot(1, 3, i+1)
463 // MAGIC     plt.subplots_adjust(wspace=1)
464 // MAGIC     plt.hist(residuals[k], bins=60, edgecolor='black', alpha=0.8)
465 // MAGIC     # Set x-axis range from -20 to 20
466 // MAGIC     plt.xlim(-15, 20)
467 // MAGIC     plt.xlabel('Residuals')
468 // MAGIC     plt.ylabel('Frequency')
469 // MAGIC     plt.title(models[k], y=-.25)
470 // MAGIC
471 // MAGIC plt.show()

```