

# New York Taxi Trip Duration

## BigData Analysis

Adarsh Dudhat  
University of Calgary  
Calgary, Canada  
[adarsh.dudhat@ucalgary.ca](mailto:adarsh.dudhat@ucalgary.ca)

Ayman Shahriar  
University of Calgary  
Calgary, Canada  
[ayman.shahriar@ucalgary.ca](mailto/ayman.shahriar@ucalgary.ca)

Anh Tuan Pham  
University of Calgary  
Calgary, Canada  
[anhuan.pham@ucalgary.ca](mailto:anhuan.pham@ucalgary.ca)

**Abstract**—New York city is one of the major cities of the world, and taxi cabs are a popular form of transport for its residents. That is why in this project, we have explored ways of using machine learning models to predict the duration of taxi trips when given details such as start and end coordinates, and the starting time. The data was cleaned, and visualized, and a positive linear correlation between the distance and trip duration could be seen in certain visualizations. A K-means clustering model was created to divide the segments into groups based on important trip data. A multiple linear regression model was then created to predict the trip duration, but this model did not have satisfactory performance. Finally, a gradient boosted trees model was created to predict the trip duration, and the performance of this model is much better than the linear regression model. This project aims to be particularly useful towards those, who want to check their trip time in NYC. For future work, we would consider experimenting with different machine learning models.

## I. INTRODUCTION

New York City (NYC) is one of the biggest cities in the world, with a population of around 18,000,000 people. Due to its large population, traffic congestion is one of the most critical problems in New York. In addition, taxi cabs are a popular choice of transport when traveling within this city. That is why in this project, we chose to predict the trip duration for a taxi in New York using the “NYC taxis trip duration dataset” which can be found [here](#) and so people can get an idea of how long their taxi trip will take and plan their trips beforehand.

The current dataset is one released by the NYC Taxi and Limousine Commission, which includes pickup time, geographic coordinates, number of passengers, and several other variables. The dataset contains more than 510,000 rows.

Since the dataset is mostly from a Kaggle challenge (where people compete to use the provided data to create the best visualizations, machine learning models and reports), there are a lot of machine learning models that were built using this dataset. However, we found out that the accuracy of these models is not very high, which means that there is a possibility of developing models that can achieve higher accuracy and reliability. Moreover, the current machine learning models were built using mostly Tensorflow and other CNN techniques, while we manipulate the data and create a model using SkLearn and PySpark.

In our solution, we have created several visualizations that show a relationship between trip duration and other variables, also we demonstrated how the mean passenger count, total trips, distance and duration vary on different weekdays. Also, we have created 3 machine learning models to either gain insights into the structure of the data or try to predict the duration of a trip. That includes regression

models like linear regression, gradient boosted trees, and K-means clustering.

## II. BACKGROUND AND RELATED WORK

#### *A. Technical Background*

To understand the report, it requires a good understanding of the dataset, supervised and unsupervised machine learning techniques such as linear regressions, gradient boosted trees and KMeans clustering.

### ***B. Review of existing work pertinent to Project***

We used only one dataset, as opposed to combining several datasets, so there was not much work that needed to be done prior to starting with the project. But it was important to understand the dataset and values of each label column.

### III. METHODOLOGY

## A. Experiment Setup

To run the notebook, users should have the following python packages installed in their environment: warning, PySpark, Numpy, pandas, matplotlib, and so on. To complete the setup, please check README.md

### ***B. Data Cleaning and Preprocessing***

To clean the data, the team has decided to check for values that are not valid such as `NaN`, `trip_duration < 0`, and many others.

```
In [57]: 1 # Change Datatype for each column.
2
3 def change_type(df, column_name, column_type, is_time):
4     if column_type == "int" or column_type == "float":
5         new_df = df.withColumn(column_name, df[column_name].cast(column_type))
6     elif column_type == "bool":
7         new_df = df.withColumn(column_name, functions.when(df[column_name] == 'N', 0).otherwise(1))
8     elif is_time == True:
9         new_df = df.withColumn(column_name, to_timestamp(df[column_name], "yyyy-MM-dd HH:mm:ss"))
10    else:
11        new_df = df
12
13    return new_df

In [58]: 1 # Calling change_type function for each column.
2
3 column_types = ["int", "time", "time", "int", "float", "float", "float", "float", "bool", "int"]
4 count = 0
5
6 for col in df.dtypes:
7     if column_types[count] == "time":
8         df = change_type(df, col[0], column_types[count], True)
9     else:
10        df = change_type(df, col[0], column_types[count], False)
11    count += 1

In [59]: 1 # Printing Datatype for all columns.
2
3 for col in df.dtypes:
4     print(col[0] + " : " + str(col[1]))
```

Fig. 1. Changing column type from String.

While preparing the data for the machine learning models, there were several operations performed on the

dataset. We have changed the type of each column from string to either integer, double, string or timestamp according to its values as shown in the above figure (Fig. 1).

```
In [60]: 1 # Calculate distance between two places.
2
3 def distance(pickUpLat1, dropOffLat2, pickUpLon1, dropOffLon2):
4     longitude = radians(dropOffLon2) - radians(pickUpLon1)
5     x = sin(longitude / 2)**2 + cos(radians(pickUpLat1)) * cos(radians(dropOffLat2)) * sin(longitude / 2)**2
6
7     temp = 2 * asin(sqrt(x))
8     earth_radius = 6371
9     return(temp * earth_radius)

In [61]: 1 # Calculating distance between pickup and dropoff latitude and longitude. In (K.M.)
2
3 distance_func = udf(distance, FloatType())
4 df = df.withColumn("distance", distance_func(df.pickup_latitude, df.dropoff_latitude, df.pickup_longitude, df.dropoff_longitude))


```

Fig. 2. Calculating distance for each trip based on its geographies.

We have calculated the distance between pickup latitude, longitude and dropoff latitude and longitude for all the trips as shown in the above figure (Fig. 2).

```
In [62]: 1 # filter the data where distance is greater than 0.
2 df = df.filter(df.distance > 0)

In [63]: 1 # Remove all trips that have trip_duration more than 7200 seconds which is equivalent to more than 2 hours and less
2 # Remove all trips which total distance is less than 1.5 and more than 60 KM because no taxi accepts rides as the c
3
4 df = df.filter((df.trip_duration > 180) & (df.trip_duration < 7200) & (df.distance > 1.5) & (df.distance < 60))

In [64]: 1 df.show()

+-----+
| vendor_id | pickup_datetime | dropoff_datetime | passenger_count | pickup_longitude | pickup_latitude | dropoff_longitude |
+-----+
| 2 | 2016-03-14 17:24:55 | 2016-03-14 17:32:30 | 1 | -73.982155 | 40.767937 | -73.984663 | | | |
| 40 | 765664 | 0 | 455 | 12.4444735 | 0 | -73.99948 |
| 1 | 2016-05-12 00:43:35 | 2016-05-12 00:54:06 | 1 | -73.980415 | 40.738564 | -73.99948 |
| 40 | 73115 | 0 | 663 | 12.6599078 | 0 | -73.99948 |
| 2 | 2016-01-19 11:35:24 | 2016-01-19 12:10:48 | 1 | -73.97903 | 40.76394 | -74.00533 |
| 40 | 71080 | 0 | 2124 | 3.669969 | 0 | -73.96902 | 40.75784 | -73.957405 |
| 1 | 2016-05-17 22:34:59 | 2016-05-17 22:40:46 | 1 | -73.96928 | 40.79778 | -73.92247 |
| 40 | 75896 | 0 | 341 | 1.6197599 | 0 | -73.99948 | 40.7384 | -73.98579 |
| 2 | 2016-05-21 07:54:58 | 2016-05-21 08:20:49 | 1 | -73.96928 | 40.79778 | -73.92247 |
| 40 | 73215 | 0 | 231 | 1.15014 | 6.528439 | 0 | -73.99948 | 40.7384 | -73.98579 |
| 1 | 2016-05-27 23:12:23 | 2016-05-27 23:13:38 | 1 | -73.98265 | 40.76384 | -74.00223 |
| 40 | 7310188 | 0 | 255 | 1.9105701 | 0 | -73.98265 | 40.76384 | -74.00223 |
| 2 | 2016-05-15 11:16:11 | 2016-05-15 11:14:59 | 4 | -73.99153 | 40.74944 | -73.95654 |
| 40 | 77063 | 0 | 1128 | 4.8806133 | 0 | -73.99153 | 40.74944 | -73.95654 |


```

Fig. 3. Filter the trips based on distance and Example DF after filtering.

We keep only trips whose trip duration is between 3 minutes to 2 hours as taxi trips in New York City might not take more than 2 hours at any time. In addition to that we believe that taxi rides in New York city may always be less than 60 kilometers as no taxi accepts rides as the destination is more than 60 kilometers far away. (Fig. 3)

```
Add timeline for each ride

In [65]: 1 df = df.withColumn("week_day", date_format(df["pickup_datetime"], "E"))
2 df = df.withColumn("year", year(df["pickup_datetime"]))
3 df = df.withColumn("month", month(df["pickup_datetime"]))
4 df = df.withColumn("quarter_of_year", quarter(df["pickup_datetime"]))
5 df = df.withColumn("hour", date_format(df["pickup_datetime"], "H"))
6
7 df = df.drop("pickup_datetime")
8 df = df.drop("dropoff_datetime")

In [66]: 1 df.show()

+-----+
| vendor_id | passenger_count | pickup_longitude | pickup_latitude | dropoff_longitude | dropoff_latitude | store_and_fwd_flag | trip_duration | distance | week_day | year | month | quarter_of_year | hour |
+-----+
| 2 | 1 | -73.982155 | 40.767937 | -73.96463 | 40.765602 | 0 | 17 | 1 | Sun | 2016 | 3 | 1 | 1 |
| 455 | 2.4444735 | Mon | 1 | -73.980415 | 40.738564 | -73.99948 | 40.73115 | 0 | 1 | 1 | 1 | 1 | 1 |
| 663 | 2.6599078 | Sun | 2016 | 0 | 663 | 12.6599078 | 0 | 1 | 1 | 1 | 1 | 1 |
| 2124 | 3.669969 | Tue | 2016 | 1 | -73.97903 | 40.76394 | -74.00533 | 40.710087 | 0 | 1 | 1 | 1 | 1 |
| 341 | 1.6197599 | Fri | 2016 | 6 | -73.96902 | 40.75784 | -73.95705 | 40.765896 | 0 | 1 | 1 | 1 | 1 |
| 1 | 6.528439 | Sat | 2016 | 5 | -73.99948 | 40.79778 | -73.92247 | 40.76056 | 0 | 1 | 1 | 1 | 1 |
| 255 | 1.9105701 | Fri | 2016 | 5 | -73.99948 | 40.7384 | -73.98579 | 40.732815 | 0 | 1 | 1 | 1 | 1 |
| 1274 | 2.7310188 | Tue | 2016 | 5 | -73.98265 | 40.76384 | -74.00223 | 40.73299 | 0 | 1 | 1 | 1 | 1 |


```

Fig. 4. Adding timeline for each ride and Example DF.

Last but not the least, we added quarter, year and day of the week to the dataset to get more accurate predictions from the machine learning model. (Fig. 4)

## C. Data Visualization

From the figure below (Fig. 5) we can see that the total number of trips is higher on Friday and Saturday with the Average Passenger count higher on Sunday as well because people are more likely to go out with the family at the start of the weekend.

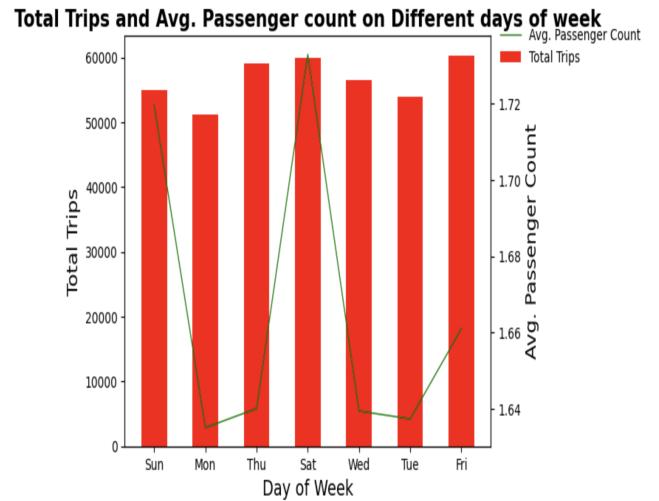


Fig. 5. Total trips and average passengers count on different week days.

This figure (Fig. 5) also shows us that as the weekend arrives people are more likely to take taxi trips to their destination with their families and friends.

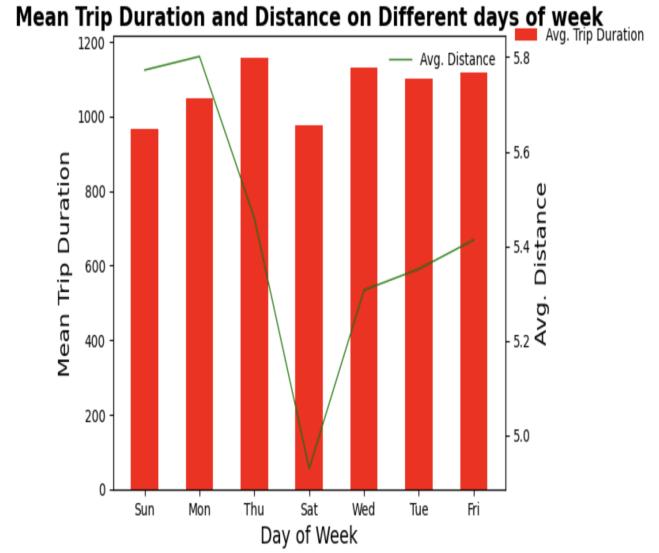


Fig. 6. Average trip duration and distance on different week days.

We can conclude from the above figure (Fig. 6) that trip duration is higher on the weekdays compared to weekends.

This figure (Fig. 6) helps us to match our initial expectation that trip duration is lower on weekends as there would be less traffic because most of people are not working.

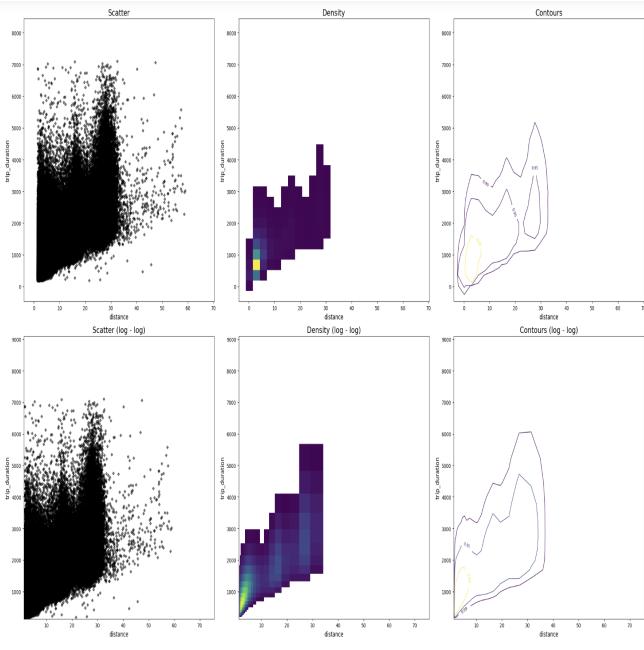


Fig. 7. Scatter, Density and Contours(log) between distance and trip duration.

While predicting and getting detailed insights of the data regarding any duration of the trip the most important variable to consider is a distance and also most people rely on the distance of two places to predict their duration of the trip.

From the figure above (Fig. 7), we can say that there is a linear positive correlation between the distance and the duration of the trip. Also, from the density and Contours plots, we can say that most of the data has distance and trip duration at small scale.

Around 70% of the trip data has a distance less than 10 kilometers and its trip duration is less than 2000 seconds which is nearly an hour. Hence, we can conclude that data contains some outliers as most of the data is compacted and there are some trips whose trip duration and distance is very huge.

#### **D. Data Preparation for Regression**

As mentioned above, we have done some data processing to remove potential outliers. As both regression models require all features to be compressed into a vector, so we have converted a week\_day column into an index column (Fig. 8).

```
In [11]: # Convert week day string column into index column.
label_index = StringIndexer(inputCol = 'week_day', outputCol = 'week_day_index')
df = label_index.fit(df).transform(df)
df.show()

+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|vendor_id|passenger_count|pickup_longitude|dropoff_longitude|dropoff_latitude|store_and_fld_flag|trip_duration|distance|
+-----+-----+-----+-----+-----+-----+-----+-----+
|        2|           1|     -73.982155|      40.769793|      -73.96463|       40.765602|          0|    455.2|4.444735
|       3|           1|     -73.980415|      40.738564|      -73.99948|       40.73115|          0|   663.2|4.559078
|       1|           2|           4.0|           4.0|           4.0|           4.0|          0|    2124.3|3.66994
|       2|           1|     -73.97903|      40.76394|      -74.00533|       40.710087|          0|    110.0|4.444735
|Tue| [2016] |       1|           1|     -73.96902|      40.75784|      -73.957405|       40.765896|          0|    341.1|4.619759
|Fri| [2016] |       6|           2|     -73.96928|      40.79978|      -73.92247|       40.76056|          0|    1551.6|5.628439
|Sat| [2016] |       2|           1|           1.0|           1.0|           1.0|           1.0|          0|          0|        0.0
+-----+-----+-----+-----+-----+-----+-----+-----+
```

Fig. 8. Converting week day string values column into integers.

Also, we have vectorised all the features and converted them into a PySpark dataframe as shown in the figure below (Fig. 9.)

```
In [12]: # Creating a list of input columns headers
input_col = [i[0] for i in df.dtypes if i[1] == "int" or i[1] == "double"]
input_col.remove("trip_duration")

# Creating a dataframe of input columns.
input_features = VectorAssembler(inputCols = input_col, outputCol = "input_features")
input_column = input_features.transform(df)
input_column.select("input_features").show()

23/01/15 23:58:49 WARN package: truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting spark.sql.debug.maxToStringFields'.
+-----+
| input_features |
+-----+
|[2,0,1,0,-73,9821,...|
|[1,0,1,0,-73,9804,...|
|[2,0,1,0,-73,9790,...|
|[1,0,1,0,-73,9690,...|
|[2,0,1,0,-73,9692,...|
|[1,0,1,0,-73,9994,...|
|[2,0,1,0,-73,9926,...|
|[2,0,1,0,-73,9915,...|
|[2,0,1,0,-73,9629,...|
|[2,0,1,0,-73,9321,...|
|[1,0,1,0,-74,0039,...|
|[1,0,1,0,-73,9803,...|
|[2,0,1,0,-73,9795,...|
|[1,0,1,0,-73,9935,...|
|[2,0,1,0,-73,9932,...|
|[2,0,1,0,-73,9554,...|
|[1,0,1,0,-73,9807,...|
|[2,0,1,0,-73,9942,...|
|[1,0,1,0,-73,9821,...|
|[1,0,1,0,-73,9709,...|
+-----+
only showing top 20 rows

In [13]: data_frame = input_column.select("input_features", "trip_duration")
data_frame.show(5)

+-----+-----+
| input_features|trip_duration |
+-----+-----+
|[2,0,1,0,-73,9821,...|      455|
|[1,0,1,0,-73,9804,...|      663|
+-----+
```

Fig. 9. DataFrame for the regression machine learning model.

In both regression models, we have used the trip\_duration column as an output column and other columns as an input column. And after all processing, 70% of the DataFrame split into training and 30% DataFrame was for testing. So, both regression models will have enough data to identify trends in the trip\_duration of NYC taxi trips.

## *E. Data Preparation for Clustering*

KMeans clustering requires all features to be an indexer, we have converted a `week_day` column into an index column (Fig. 10).

```
# Convert a week_day string column to an index.
df = df.drop(df['store_and_fwd_flag'])
week_day_col = StringIndexer(inputCol = 'week_day', outputCol = 'week_day_index')
df = week_day_col.fit(df).transform(df)
df.show(5)

+-----+-----+-----+-----+-----+-----+-----+-----+
|vendor_id|passenger_count|pickup_longitude|pickup_latitude|dropoff_longitude|dropoff_latitude|trip_duration|distance|week_day|year|month|
|quarter_of_year|hour|week_day_index|
+-----+-----+-----+-----+-----+-----+-----+-----+
| 2 | 1 | 1 | -73.982155 | 40.767937 | -73.96463 | 40.765602 | 455 | 2.4444735 | Mon | [2016] | 3 | |
| 1 | 17 | 1 | 6.0 | -73.980415 | 40.738564 | -73.99948 | 40.731115 | 663 | 2.6599078 | Sun | [2016] | 6 |
| 2 | 0 | 1 | 4.0 | -73.97903 | 40.76394 | -74.00553 | 40.710087 | 2124 | 3.669969 | Tue | [2016] | 1 |
| 2 | 11 | 1 | 5.0 | -73.980415 | 40.738564 | -73.99948 | 40.731115 | 663 | 2.6599078 | Sun | [2016] | 6 |
| 1 | 22 | 0 | 0.0 | -73.980415 | 40.738564 | -73.99948 | 40.731115 | 663 | 2.6599078 | Sun | [2016] | 6 |
```

Fig. 10. Converting week day string values column into integers.

All other input data requirements are mostly similar to the regression models, most of the data processing has been done in the similar way for Clustering algorithms.

#### IV. ALGORITHMS

### *A. Linear Regression*

To identify the best parameters and overcome overfitting problem for the linear regression machine learning model, we used CrossValidator for `regParam`, `elasticNetParam`, `maxIter` and `fitIntercept` input parameters for the model.

```
In [21]: # print best parameter of the model.
print("Best maxBins Param: ", model.bestModel._java_obj.getBestParam())
print("Best maxDepth Param: ", model.bestModel._java_obj.getBestParam())
print("Best maxIter Param: ", model.bestModel._java_obj.getBestParam())
print("Best fitIntercept Param: ", model.bestModel._java_obj.getFitIntercept())
print("Intercept: ", model.bestModel.intercept)

# Accuracy and coefficients of the linear equation.
print("Coefficients of Determination(Accuracy on Training data): ", str(round(model.bestModel.summary.r2*100, 2)), "%")
print("Coefficients: ", model.bestModel.coefficients)

Best Regression Param: 0.01
Best Elastic Net Param: 0.25
Best maxBins Param: 20
Best fitIntercept: True
Intercept: -19235.005973344687
Coefficient of Determination(Accuracy on Training data): 50.12 %
Coefficients: [0.749325910707641, 2.2572463385882204, -325.336075706018, -217.6755819562678, -536.4434332987057, -862.7021553118498, 75.327
96525925, 82.03011447964624, 0.0, 18.879690528403284, 22.08299544444144, 5.383600297732522, -12.10626879857068]
```

Fig. 11. Best Parameters for linear regression model.

Above figure (Fig. 11) shows the best input parameters for the linear regression model. With these parameters, the accuracy of the model on training DataFrame is roughly 50%.

```
In [22]: # Predict.
prediction = model.transform(test)

prediction.show(5)

|Stage 7523>          (0 + 1) / 1|
+-----+-----+
| input_features|trip_duration| prediction|
+-----+-----+
|[1,0,1,0,-74.23101...| 1880 | 4318.493897415774|
|[1,0,1,0,-74.2015...| 340 | 885.8448166915368|
|[1,0,1,0,-74.0184...| 667 | 1244.345431175642|
|[1,0,1,0,-74.0180...| 3551 | 1487.224294263888|
|[1,0,1,0,-74.0179...| 1054 | 1244.150761453415|
+-----+-----+
only showing top 5 rows
```

```
In [23]: # Evaluator for testing data.
test_pred_evaluator = RegressionEvaluator(predictionCol = "prediction", labelCol = "trip_duration", metricName = "r2")

# Print RMSE and accuracy of the model on testing data.
print("RMSE: ", evaluator.evaluate(prediction))
print("Coefficient of Determination(Accuracy on Test Data): ", str(round(test_pred_evaluator.evaluate(prediction)*100, 4)), "%")

RMSE: 500.76489579072614
|Stage 7525>          (0 + 1) / 1|
Coefficient of Determination(Accuracy on Test Data): 50.3202 %
```

Fig. 12. Prediction, RMSE and accuracy on test data.

Figure 12. shows the root mean square error (RMSE) on training DataFrame which is nearly 500 and accuracy of the model on testing data is slightly less than 50.5%.

## B. Gradient Boosted Trees Regression

Similar to the linear regression model, to identify the best parameters and overcome overfitting, we used CrossValidator on maxDepth, maxIter and maxBins input parameters of the gradient boosted trees model.

```
In [27]: # Print best depth, bins, maxIter, impurity and lossType.
print("Best maxDepth Param: ", model_gbt.bestModel._java_obj.getBestDepth())
print("Best maxBins Param: ", model_gbt.bestModel._java_obj.getBestNumBins())
print("Best maxIter Param: ", model_gbt.bestModel._java_obj.getBestMaxIter())
print("Best Impurity Param: ", model_gbt.bestModel._java_obj.getBestImpurity())
print("Best LossType Param: ", model_gbt.bestModel._java_obj.getBestLossType())

Best maxDepth: 10
Best maxBins: 64
Best maxIter: 20
Best Impurity: variance
Best LossType: squared
```

Fig. 13. Best Parameters for Gradient boosted trees model.

Figure 13 shows the best input parameters for the gradient boosted trees regression model and with these input parameters we achieved the accuracy of 75% on testing data with root mean square error (RMSE) nearly 360.

```
In [29]: # Evaluator for testing data for gradient boosted trees regression model.
test_pred_evaluator_gbt = RegressionEvaluator(predictionCol = "prediction", labelCol = "trip_duration", metricName = "r2")

# Print RMSE and accuracy of the model on testing data for gradient boosted trees regression.
print("RMSE: ", gbt_evaluator.evaluate(prediction_gbt))
print("Coefficient of Determination(Accuracy on Test Data): ", str(round(test_pred_evaluator_gbt.evaluate(prediction_gbt)*100, 4)), "%")

RMSE: 362.0912865178779
|Stage 18732>          (0 + 1) / 1|
Coefficient of Determination(Accuracy on Test Data): 74.0255 %
```

Fig. 14. RMSE and coefficient of determination of the gradient boosted tress regression model.

Hence, we can conclude from the evidence that Gradient boosted trees regression model is better machine learning model compared to linear regression because of lower RMSE and higher Coefficient of Determination.

## C. KMeans Clustering

### 1. Finding optimal number of clusters

The K-means clustering is an unsupervised machine learning technique that groups similar data points together based on underlying similarities and patterns. Aim of the KMeans clustering was to find out if data points within the same clusters have similar trip durations or not. If any correlation was found, it becomes easier to find out which properties of the dataset are strongly influencing the trip duration.

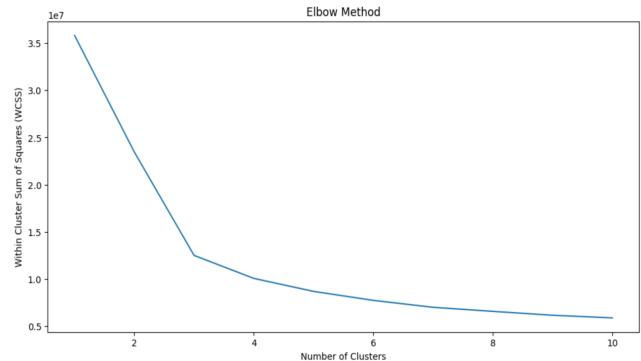


Fig. 15. Elbow Method (Determination of optimal clusters).

While working with data that has multiple dimensions, it is not possible to visualize the data and decide on the optimal number of clusters to use for the K-means clustering model. However, there is a mathematical way to determine the optimal number of clusters, called the “elbow method”. To use this method, we have to create clustering models using a different range of clusters, and for each cluster, we need to compute the Within Cluster Sum of Squares (WCSS). WCSS is the within-cluster-sum-of-squares which means the sum of squared distance between each point and the centroid of the cluster. With this method in (Fig. 15), we found the optimal number of clusters is k = 4.

```
In [87]: #Using StandardScaler to set all the columns in the dataframe to the same scale.
std_scaler = StandardScaler()
X = std_scaler.fit_transform(input_col_pandas_df)

X = pd.DataFrame(X, columns = input_col_pandas_df.columns)
X
```

```
Out[87]: vendor_id  passenger_count  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude  distance  year  month  quarter_of_year  hour  v
0  0.932331  -0.508690  -0.295293  0.608663  0.085584  0.949561  -0.462010  0.0  -0.35667  -1.013854  0.527376
1  -1.075580  -0.508690  -0.257719  -0.345022  -0.723479  -0.442761  -0.447515  0.0  1.469916  0.988333  -2.057387
2  0.932331  -0.508690  -0.227810  0.478802  -0.859290  -1.16052  -0.283727  0.0  -1.506056  -1.013856  -0.384886
3  -1.075580  1.778467  -0.011651  0.280767  0.253316  0.029553  -0.615457  0.0  1.469916  0.988333  1.287956
4  0.932331  -0.508690  -0.017265  1.577407  1.064351  0.257718  0.178782  0.0  0.874721  0.988333  -0.993061
```

Fig. 16. Set all column data at the same scale for KMeans.

To fit the data into KMeans cluster and set all the columns data at the same scale, we used StandardScaler data preprocessing function and used output to fit in our KMeans model.

```
In [89]: # Graph out the size of each cluster.
input_col_pandas_df['cluster'] = clusters_prediction
cluster_counts = input_col_pandas_df.groupby(['cluster']).count()

plt.figure(figsize=(12, 6), dpi=100)
plt.bar(['0', '1', '2', '3'], cluster_counts.values, width = 0.3)

plt.xlabel('Cluster')
plt.ylabel('Total Data points Count in a Cluster')
plt.show()
```

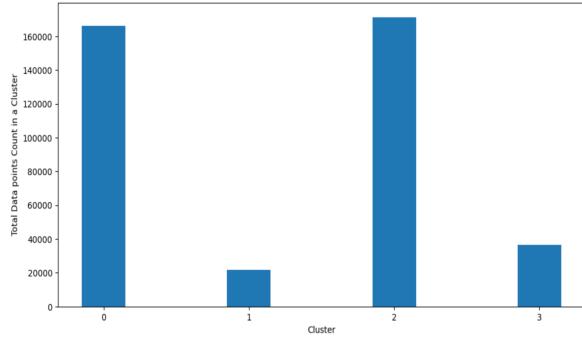


Fig. 17. Total data points in each cluster.

As we have seen before that almost 70% of the data points are nearly clustered together, that leads to the conclusion that most of the data points might reside in the same cluster. As we can see in the above figure (Fig. 17) that most of the data points are in cluster 0 and 2 as we expected.

## 2. Cluster diagnoses

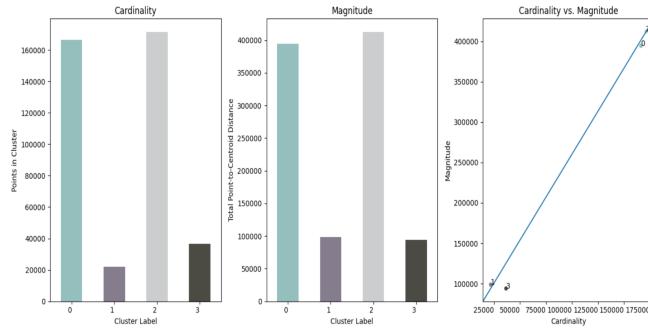


Fig. 18. Cluster cardinality, magnitude and Cardinality vs. Magnitude.

Now, it is important to get good cluster diagnoses. There are three main diagnoses that we have prepared which includes cluster cardinality, cluster magnitude and cluster cardinality vs. cluster magnitude. Cluster cardinality means the total number of data points within each cluster. While cluster magnitude shows the total points to centroid distance per cluster. That is used to check how widely data points are spread in each cluster. As we can see in the above figure (Fig. 18) clusters with higher cardinality tend to have higher magnitude.

## 3. Cluster characteristics

To make it easier for a business person and a data scientist to identify the key characteristics for each cluster and as it is important for them to know what it stands for, we have generated a box plot (Fig. 19) that shows each feature distribution per cluster.

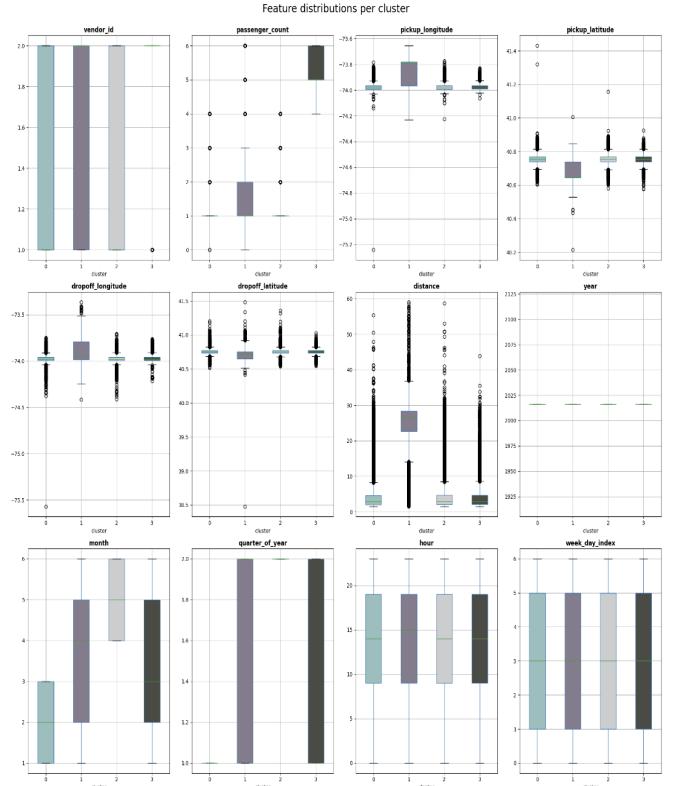


Fig. 19. Box plot for the feature distribution for each cluster.

From this box plot (Fig. 19), we can conclude that passenger\_count, pickup\_longitude, dropoff\_longitude, distance, and month are the main features for a data point to reside in a particular cluster.

```
In [94]: # Calculates mean for each feature in each cluster.
X_mean = pd.concat([pd.DataFrame(X.mean()).drop('cluster'), columns=['mean']], X.groupby('cluster').mean().T, axis = 1)
X_mean
```

	0	1	2	3
vendor_id	-0.098066	0.055865	-0.109343	0.925092
passenger_count	-0.300208	0.021665	-0.290389	2.711194
pickup_longitude	-0.158842	2.642086	-0.154917	-0.130937
pickup_latitude	0.131116	-2.222794	0.127117	0.139509
dropoff_longitude	-0.101386	1.869185	-0.095520	0.082480
dropoff_latitude	0.063276	-1.059515	0.060757	0.060466
distance	-0.187411	3.043984	-0.172868	-0.163478
year	0.000000	0.000000	0.000000	0.000000
month	-0.884719	0.023120	0.861630	-0.027033
quarter_of_year	-0.103856	0.021634	0.986333	-0.023090
hour	0.004716	0.037538	-0.006411	-0.013838
week_day_index	-0.007553	0.083785	-0.000668	-0.013494

Fig. 20. Mean for each feature per cluster.

Last but not the least radar plot, which is another best way to summarize all relevant information in one plot. In order to avoid crashes from the visualization and maintain the scale of the plot, we have used the StandardScale DataFrame and calculated the mean for each feature per cluster (Fig. 20) and generated the radar plot as shown in the figure below (Fig. 21).

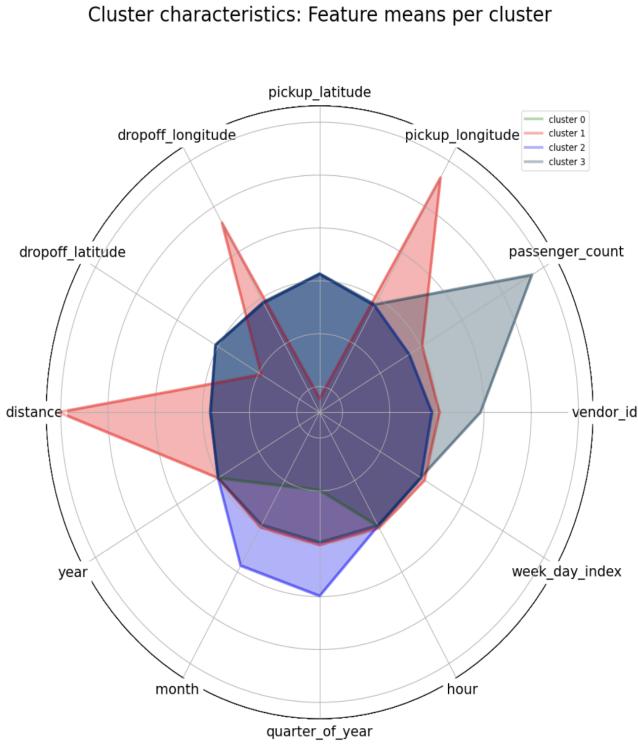


Fig. 21. Radar plot.

From the above plot (Fig. 21), the decision for a datapoint to be in cluster 1 and 3 highly depends on the distance, dropoff\_longitude, pickup\_logitude and passenger\_count respectively. While data points in cluster 2 highly depended on the month and quarter of the year.

## V. RESULT, CONCLUSION AND FUTURE WORK

This entire project involves 6 major steps: Exploring the dataset, Data Cleaning, Data Preprocessing, Data Visualization, Clustering and Regression Machine Learning models.

By visualizing a distance and trip duration on both linear and log-log plots, we saw that these two properties display a certain amount of positive linear correlation. Alongside, roughly 70% of the NYC trip has a distance less than 10 kilometers and their trip duration is less than 2000 seconds which is less than half an hour. That led us to the conclusion that data has a lot of outliers in terms of distance and trip duration. In addition to this, we reached the conclusion that the number of trips and average passenger counts increase as weekends arrive. On the other hand, trip duration seems to be higher during weekdays (Fig. 6).

As mentioned before, we used CrossValidator to find the best possible parameters and generalize a pattern in the data for our machine learning models. The root mean square error (RMSE) for the linear regression model was 500.76 and coefficient of determination was 50.32%. Hence, the linear regression model is not the best machine learning model for predicting the NYC taxi trip duration.

On the other hand, with the gradient boosted trees regression, the root mean square error is 362 and coefficient of determination (accuracy) is roughly 75% because of the better hyper parameter tuning and optimization on different

loss functions. Hence, we can conclude that gradient boosted trees regression provides us better prediction compared to linear regression model.

For the clustering, it was difficult to see the differences in clusters. Although, we have performed several iterations while changing the parameters and by plotting several graphs to identify the patterns. Still we have seen some of the features like distance, passenger\_counts, pickup\_longitude, and dropoff\_longitude really help us to clusterized the data as shown in the above figures (Fig. 18, 19, 21).

To use this huge data set and build machine learning models, it requires to have really good hardwares with good unified memory, better SSD storage and most importantly a great GPU. We have limited hardware resources, hence we were not able to experiment our machine learning model with several different hyper tuning parameters. So, for the future iteration, it is worthwhile to experiment with different hyper tuning parameters for the regression model.

While we were only able to use the Elbow method to determine the best number of clusters for our data because of the limited resources. For the future, we would like to use silhouette analysis, Agglomerative hierarchical Clustering and Gaussian Mixtures model to identify the best number of clusters and do better cluster diagnoses.

## VI. REFERENCES

- <https://www.macrotrends.net/cities/23083/new-yorkcity/population#:~:text=The%20current%20metro%20area%>
- <https://www.kaggle.com/competitions/nyc-taxi-trip-duration>
- <https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1>
- <https://towardsdatascience.com/best-practices-for-visualizing-your-cluster-results-20a3baac7426>
- <https://towardsdatascience.com/common-mistakes-in-cluster-analysis-and-how-to-avoid-them-eb960116d773>