

TAXI FAIR PREDICTION

STUDENT B 12/11/2018

Project Overview:

Why we are doing this project?

Many taxi companies are struggling to stay in the business due to effective pricing models of ride share companies like Uber & Lyft

Goal: Use City of Chicago Taxi trip dataset to build a more competitive pricing model (fare) based on trip_seconds & trip_miles.

ROI: Better profit to Taxi companies & Drivers

Dataset Overview:

Key Insights:

- 1-year data is distributed in 12 files ,one each for every month
- More than 1 million records in each file
- Have 20 key attributes to explain the data.

Challenges:

- Loading time due to the data volume
- Data Cleaning
- Data analysis due to the data volume.
- Training time.

taxi id trip_start_timestamp trip_end_timestamp trip_seconds trip_miles pickup_census_tract dropoff census tract pickup_community_area dropoff_community_area fare tips tolls extras trip_total payment_type company pickup_latitude pickup_longitude dropoff_latitude dropoff longitude

Infrastructure

- ITVersity Lab Big data systems
- Spark Job History Server UI
- Putty to submit spark jobs
- Winscp to deploy code
- Ambari File view to upload dataset files to hadoop cluster.
- Achieved faster run times to train models
- Run time Around 30 mins to train all models

O HDFS	Metrics Heatmaps	Config History			
O YARN	Metric Actions • La	st 1 hour ·			
MapReduce2			1 Commission		10
😐 Tez	HDFS Disk Usage	DataNodes Live	HDFS Links	Memory Usage	Network Usage
Hive		125755	NameNode	37.2 GB	195.3 KB
HBase	44%	5/5	Secondary NameNode		
🖽 Pig			5 Datawooes	18.6 GB	97.6 KB
Sqoop			More		لىلل
Oozie			-		
C ZooKeeper	CPU Usage	Cluster Load	NameNode Heap	NameNode RPC	NameNode CPU W
Storm	100%	10			
Flume	506		24%	0.04 ms	0.3%
Ambari Infra					
Amban Metrics	And Anthenet States				
 Kafka 	NameNode Uptime	HBase Master Heap	HBase Links	HBase Ave Load	HBase Master Upti
O Spark			HBase Master		
Spark2	47 1 d	10%	5 Region Servers	164.4	47.1 0
A Mahout	47.1 U	(ION)	Master Web UI	104.4	47.13
C. Clinker			More		

Sport 2.3.0.2.6.5.0-292 History	Server					
Event log directory: hdfs:///spark2-history/						
Last updated: 2018-12-11 04:59:08						
Client local time zone: America/Chicago						
Show 20 • entries						
App ID	App Name	Attempt ID	Started	Completed	Duration	Spark User
application_1540458187951_23124	ChicagoTaxiTrips		2018-12-10 14:45:19	2018-12-10 15:38:00	53 min	chodisettiramakrishna
application_1540458187951_23108	ChicagoTaxiTrips		2018-12-10 13:05:26	2018-12-10 13:58:21	53 min	chodisettiramakrishna
application_1540458187951_23088	ChicagoTaxiTrips		2018-12-10 12:09:32	2018-12-10 13:02:39	53 min	chodisettiramakrishna
application 1540458187951 23083	ChicagoTaviTrips		2018-12-10 11:42:14	2018-12-10 12:07:50	26 min	chodisettiramakrishna

spark-submit --master yarn --deploy-mode client --conf spark.ui.port=50865 --conf spark.dynamicAllocation.enabled=false --num-executors 8 --executor-memory 1536M --executor-cores 2 FinalProject2.py

application 1540458187951 23083

Data Analysis:

- Dropped null values and zero values to avoid bias in the results.
- Outlier analysis on trip_seconds, trip_miles using box plot and scatter plot.
- Identified trip_seconds are less than trip_miles.
- Considered trip miles range 0 to 2000
- Considered trip_seconds range 0 to 60000 seconds

Data analysis- trip_seconds



Data analysis- trip_miles



Model Analysis with different run conditions

- Run 1 Filtered Zeros , null values
- Run 2- Filtered Zeros , null values, trip_miles < 2000,trip_seconds < 60000</p>
- Run 3- Filtered Zeros , null values, trip_seconds >trip_miles, trip_miles < 2000,trip_seconds < 60000</p>

Model 1 :Elastic net

- Mix of Lasso & Ridge regression
- Objective is to find the best regularization parameters to reduce the SSE of Elastic Net

minimize $\left\{ SSE + \lambda_1 \sum_{j=1}^p \beta_j^2 + \lambda_2 \sum_{j=1}^p |\beta_j| \right\}$

- When lambda 1 and 2=0 then we get least square parameter estimate.
- When lambda 1 >0 and lambda =0 then it is Lasso Regression.
- When lambda =0 and lambda>0 then it is Ridge regression.

	RMSE	
RUN 1	RUN 2	RUN 3
19.3432	19.1767	19.1767

```
lm_model = lm.fit(train)
print("Coefficients: " + str(lm_model.coefficients))
print("Intercept: " + str(lm_model.intercept))
trainingSummary = lm_model.summary
print("RMSE: %f" % trainingSummary.rootMeanSquaredError)
print("r2: %f" % trainingSummary.r2)
```

lm_predictions = lm_model.transform(test)
lm_predictions.select("prediction","label","features").show(10)

print("Root Mean Squared Error (RMSE) on test data = %g" % test_result.rootMeanSquaredError)
lm_predictions = lm_model.transform(test)
lm_predictions.select("prediction","label","features").show(5)

MODEL 2: SIMPLE TREE MODEL

- In this technique, we split the population or sample into two or more homogeneous sets (or sub-populations) based on most significant splitter / differentiator in input variables.
- It works for both categorical and continuous input and output variables.

#*************************************		
#*************************************		
dt_taxiTrips_Tree = DecisionTreeRegressor(featuresCol = 'features', labelCol =	'label',	ma

dt_tax11r1ps_1ree = Dec1s1on1reekegressor(featuresto1 = f	reatures , lapellol =	lapel, maxuepth = 3)
decision_Model = dt_taxiTrips_Tree.fit(train)		
<pre>predictions = decision_Model.transform(test)</pre>		
predictions.show()		

	RMSE	
RUN 1	RUN 2	RUN 3
17.8708	17.9874	17.9874

MODEL 3: RANDOM FOREST

Random forests or **random decision forests** are an ensemble learning method

for classification, regression and other tasks that operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees to add text

rmse = dt_evaluator.evaluate(predictions)
print("Root Mean Squared Error (RMSE) on test data = %g" % rmse)
predictions.show(5)

	RMSE	
RUN 1	RUN 2	RUN 3
17.6018	17.7166	16.3496

OUTCOMES FROM THREE MODELS

- Out of three models ,Random forest model comes with low RMSE compared to others.
- RMSE was lower for Random forest even for the three runs

K-MEAN CLUSTERING

- Assemblers & transformation
- Feature Scaling
- Building the Model

52 # K-means 34 print("K- means started") 55 assembler = VectorAssembler(inputCols=['trip_seconds','trip_miles','fare'], outputCol='features') i6 assembled_data = assembler.transform(df_taxiTrips) 37 assembled_data.show(5) 58 scaler = StandardScaler(inputCol='features', outputCol='scaledFeatures') 59 scaler_model = scaler.fit(assembled_data) \$1 scaled_data = scaler_model.transform(assembled_data) Sealed_data.printSchema() scaled_data.show(5) \$4 55 kmeans = KMeans().setK(2).setSeed(1) 36 model = kmeans.fit(scaled data) >7 predictions = model.transform(scaled_data) 38 >9 predictions.groupBy('prediction').count().show()

CLUSTER OPTIMIZATION-SILHOUETTE METHOD

- Provides a measure of how close each point in one cluster is to points in the neighboring clusters.
- This metric (silhouette width) ranges from -1 to 1 for each observation

in your data

```
Close to 1Matched to assigned clusterClose to 0Borderline match between two clustersClose to -1May be assigned to wrong cluster
```

Silhouette with squared euclidean distance

RUN 1	RUN 2		RUN	13		
0.9992568	368	0.8225	i96	0.832	2390	661
				·	*	*
	1 2	3 4	5 6 k	7	8	9

```
evaluator = ClusteringEvaluator()
silhouette = evaluator.evaluate(predictions)
print("Silhouette with squared euclidean distance = " + str(silhouette))
centers = model.clusterCenters()
print("Cluster Centers: ")
for center in centers:
    print(center)
```

```
evaluator = ClusteringEvaluator()
cost = np.zeros(20)
silhouette = np.zeros(20)
for k in range(2,20):
```

```
kmeans = KMeans().setK(k).setSeed(1).setFeaturesCol("features")
model = kmeans.fit(scaled_data.sample(False,0.1, seed=42))
cost[k] = model.computeCost(scaled_data)
predictions = model.transform(scaled_data)
silhouette[k] = evaluator.evaluate(predictions)
```

```
import matplotlib.pyplot as plt
%matplotlib inline
fig, ax = plt.subplots(1,1, figsize =(8,6))
ax.plot(range(2,20),cost[2:20])
ax.set_xlabel('k')
ax.set_ylabel('cost')
```

```
import matplotlib.pyplot as plt
%matplotlib inline
fig, ax = plt.subplots(1,1, figsize =(8,6))
ax.plot(range(2,20),silhouette[2:20])
ax.set_xlabel('k')
ax.set_ylabel('silhouette')
```

