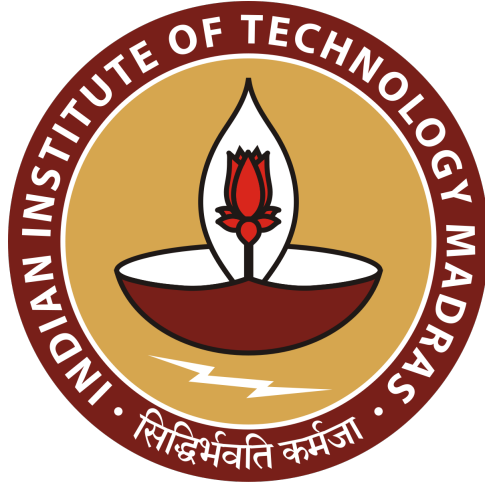


INDIAN INSTITUTE OF TECHNOLOGY MADRAS



Big Data Laboratory - CS4830 Final Project- Team Event Horizon

G Prashant (BS17B011)
Jeswant Krishna K (AE17B030)
Avinash G Bagali (AE17B110)

May 16, 2021

Contents

1	Project Overview and Objective	3
2	Task 1: Batch Computation	3
2.1	Data Exploration and Preprocessing	3
2.1.1	Feature Selection	4
2.1.2	Handling Categorical Variables	4
2.1.3	Imputation	5
2.1.4	Standard Scaling	5
2.1.5	Assembly	6
2.2	Model Training	6
2.2.1	Hyperparameter Tuning and Model Selection	6
2.2.2	Screenshots of Google Cloud Console Dataproc Job Log	7
3	Task 2: Real-time Computation	9
3.1	Latency	9
3.2	Real-time Prediction Results	9
4	Inferences and Conclusion	9
5	PySpark Codes	11
5.1	Pyspark Code for RandomForestClassifier Model	11
5.2	Pyspark Code for LogisticRegression Model	13
5.3	Pyspark Code for NaiveBayes Model	15
5.4	Code for Kafka Producer	17
5.5	Code for Kafka Subscriber	18

1 Project Overview and Objective

The task of this project is to perform big data processing and analysis using tools such as Apache Spark and Kafka. The dataset assigned to our team is “**NYC Parking Tickets**” dataset, which contains a large number of records pertaining to parking violations in the City of New York. Each record consists of several details including Registration State, Vehicle Make, Issuing Agency, Violation Location, etc (<https://www.kaggle.com/new-york-city/nyc-parking-tickets>). The first task is to perform preprocessing and model training using batch computation using PySpark run on a Dataproc Cluster. Precisely, the goal is to train a supervised classification model for predicting the label “Violation County” using the given data, i.e., for predicting the County where a violation has occurred. The next task is to use the saved model stored in Google Cloud Storage (GCS) bucket to perform real-time analysis of streaming data using Apache Kafka. This report provides all the necessary implementational details for each task of the project, right from preprocessing, model training, hyperparameter tuning, model selection to real-time prediction with appropriate screenshots.

2 Task 1: Batch Computation

2.1 Data Exploration and Preprocessing

The dataset contains 51 features in total, most of them being categorical features. These features provide all the information about the parking ticket issued in the City of New York. Figure 1 provides a glimpse of the first batch of the dataset, displayed as a Pandas Dataframe.

Summons Number	Plate ID	Registration State	Plate Type	Issue Date	Violation Code	Vehicle Body Type	Vehicle Make	Issuing Agency	Street Code1	Street Code2	Street Code3	Vehicle Expiration Date	Issuer Code	Issuer Command	Issuer Squad
7922553869	33684MD	NY	COM	06/17/2016	42	VAN	FORD	T	34550	10410	10510	20170430	345557	T501	B
7922553870	60996MC	NY	COM	06/17/2016	69	VAN	CHEVR	T	34550	10410	10510	88888888	345557	T501	B
7922553882	87642ME	NY	COM	06/17/2016	69	VAN	CHEVR	T	34550	10410	10510	20170930	345557	T501	B
7922553894	71985KA	NY	COM	06/17/2016	69	VAN	DODGE	T	34530	0	0	20160701	345557	T501	B
7922553900	66460AN	NY	COM	06/17/2016	69	VAN	GMC	T	34550	0	0	20180531	345557	T501	B
...
8359500927	FFM1969	NY	PAS	06/08/2016	38	SUBN	TOYOT	T	25590	56890	78820	20161128	362197	T201	P
8359500940	2482082	IN	PAS	06/08/2016	51	VAN	FRUEH	T	73980	40404	40404	88880088	362197	T201	P
8359500952	65919JW	NY	COM	06/08/2016	52	VAN	FRUEH	T	0	0	0	88888888	362197	T201	P
8359500964	68718MG	NY	COM	06/08/2016	82	PICK	DODGE	T	0	0	0	20170531	362197	T201	P

Figure 1: Screenshot showing a few records and features of the first batch of the NYC dataset

The target label column is *Violation_County*, which specifies the County in which the violation occurred. It takes one of the following four values.

- BX
- K
- NY
- Q

Figure 2 shows a bar plot of the number of occurrences of each County in the first batch of the NYC Dataset.

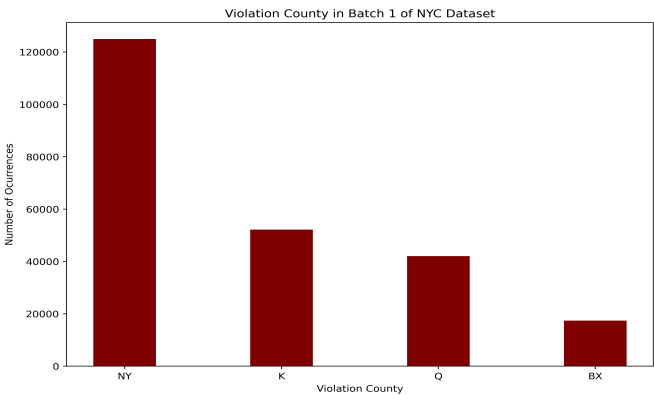


Figure 2: Violation County in Batch 1 of NYC Dataset

From the above figure, it can noticed that most of the violations have occurred in the NY County.

Based on a visual inspection of the first batch of the NYC dataset, it was evident that multiple features had missing values across all the samples of the batch. Further, there were some features which were irrelevant for the purpose of the prediction. Hence, preprocessing the dataset to extract useful features was undoubtedly necessary for this task. The following sections provide all details on the preprocessing steps performed.

2.1.1 Feature Selection

The following features had missing values across all of the data samples in a batch, and were hence removed.

- | | | |
|-------------------------------------|---------------------|-------------------------|
| • No Standing or Stopping Violation | • Longitude | • BBL |
| • Hydrant Violation | • Community Board | • NTA |
| • Double Parking Violation | • Community Council | • Plate ID |
| • Latitude | • Census Tract | • Unregistered Vehicle? |
| | • BIN | • Violation Legal Code |

Secondly, it was important to remove unnecessary features that might not be of potential importance for the purpose of predicting the Violation County. Given below are features that were removed either because they carried redundant information or because they were not required.

- | | | |
|---------------------------|------------------------|--------------------------|
| • Summons Number | • Date First Observed | • Days Parking In Effect |
| • Vehicle Expiration Date | • From Hours In Effect | • House Number |
| • Time First Observed | • To Hours In Effect | • Issuer Code |

After removing all of the above features, the dataset contained 24 features, which are mentioned below.

- | | | |
|---------------------------|-------------------------------------|-------------------------|
| • Registration State | • Issuer Command | • Sub Division |
| • Plate Type | • Issuer Squad | • Vehicle Color |
| • Issue Date | • Violation Time | • Vehicle Year |
| • Violation Code | • Violation In Front Of Or Opposite | • Meter Number |
| • Vehicle Body Type | • Street Name | • Feet From Curb |
| • Vehicle Make | • Intersecting Street | • Violation Post Code |
| • Issuing Agency | • Law Section | • Violation Description |
| • Street Codes 1, 2 and 3 | | |

2.1.2 Handling Categorical Variables

Now that the important features are selected, the next step is to handle categorical features and other non-numeric features. We considered all the features except *Violation Code*, *Street Codes (1, 2 and 3)*, *Vehicle Year* and *Feet From Curb* as categorical and non-numeric features. The features that involve representing Date and Time - *Issue Date* and *Violation Time* were transformed in a special manner for numerical representation. The details of the same given below.

Issue Date:

This feature specifies the date on which the parking ticket is issued, in MM/DD/YYYY format. We developed a User Defined Function (UDF) `dateParse` that will transform the date in this format to YYYYMMDD, of integer type. This way, it becomes a numerical feature, with the order of the dates being preserved, i.e., higher value indicates a later date. Figure 3 shows the numerical transformation of the feature *Issue Date*.

Issue Date	Issue Date
06/17/2016	20160617
06/17/2016	20160617
06/17/2016	20160617
06/17/2016	20160617
06/17/2016	20160617
06/17/2016	20160617
06/17/2016	20160617
06/27/2016	20160627
06/27/2016	20160627
06/27/2016	20160627
06/27/2016	20160627
06/27/2016	20160627
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615
06/15/2016	20160615

Figure 3: Numerical Transformation of Issue Date Feature

Violation Time

This feature specifies the time at which the parking ticket is issued, in HHMM or HHMP 12-hour format, where A and P represents AM and PM respectively. To preserve the order of this feature, a UDF was defined that can transform the time in the given format to HHMM 24-hour format, which is of integer type. This way, a higher value denotes a later time. Figure 4 shows the numerical transformation of the feature *Violation Time*.

Violation Time	Violation Time
1023A	1023
1037A	1037
1046A	1046
1115A	1115
1218P	1218
1225P	1225
1237P	1237
0636A	636
0639A	639
0646A	646
0654A	654
0705A	705
0805A	805
0822A	822
0827A	827
0831A	831
0907A	907
0918A	918
0945A	945
0950A	950

Figure 4: Numerical Transformation of Violation Time Feature

For the other remaining categorical features, it was clearly noticable that “One-Hot Encoding (OHE)” would be the best way to represent most of them, because they did not have an inherent order (for example, *Registration State* - NY, PA). Nevertheless, it is important to note that many of these features had a large number of unique values, which makes OHE an infeasible option. Therefore, Label Encoding was performed for these features using `StringIndexer` in PySpark. Moreover, even the target label *Violation County* was numerically encoded using `StringIndexer`.

2.1.3 Imputation

Imputation was performed to fill the missing values in the dataset with `Imputer` in PySpark. The missing values were imputed with the mode of the corresponding feature column, i.e., with the frequently occurring value. Imputing values with mode is more suitable for features that are non-numeric in nature.

2.1.4 Standard Scaling

The data was scaled by subtracting the feature-wise average and dividing by their standard deviation, as given for any feature x below,

$$\hat{x} = \frac{x - \bar{x}}{s_x}$$

where \bar{x} is the average and s_x is the standard deviation for feature x . This transformation was performed using `StandardScaler` in PySpark. This transformation was performed to scale the data to lie between 0 and 1 for better stability during model training.

2.1.5 Assembly

`VectorAssembler` was used to assemble all the 24 selected and transformed features, prior to model training. They were collectively named as `features` and the target column was named as `label`.

NOTE: Analysis of features and feature selection were performed with the first batch of the NYC Dataset as reference. Our assumption here is that our feature engineering steps will generalize well with other batches of the dataset.

2.2 Model Training

Three popular supervised machine learning models were used to train on the dataset - **Random Forest Classifier, Logistic Regression, Gaussian Naive Bayes**. Hyperparameter tuning was performed and the validation accuracies were considered for choosing the best model. Training multiple models and hyperparameter tuning were made easier using `TrainValidationSplit` in PySpark, where the training size was specified to be 90% of the total dataset. The remaining 10% is the validation set.

2.2.1 Hyperparameter Tuning and Model Selection

Model	Hyperparameters		Validation Accuracy
RandomForestClassifier	maxDepth	numTrees	
	20	50	99.979 %
	20	100	99.958 %
	30	50	99.979 %
	30	100	99.958 %
LogisticRegression	RegularizationParameter	ElasticNetParameter	
	0	0	69.32335 %
	0	0.5	69.32335 %
	0.5	0	55.00664 %
	0.5	0.5	47.21852 %
	1	0	50.22741 %
	1	0.5	47.21852 %
NaiveBayes	Smoothing		
	0.1		44.37758 %
	0.5		44.37758 %
	1.0		44.37758 %

Table 1: Hyperparameter Tuning of Models

As seen from the Table 1 the best model was Random Forest Classifier with hyperparameters of **maxDepth=20** and **numTrees=50** and has **ValidationAccuracy=99.979 %**. Even though one more model has same accuracy this was chosen as the best because of the lower complexity, in accordance with Occam’s Razor principle which states that simple models are better. The next part shows the screenshots of the logs obtained after training of each of the classification models. From this table, it can also be noticed that Logistic Regression model performs better when compared to with regularization (although much worse than Random Forests). Gaussian Naive Bayes model is insensitive to tuning the **smoothing** hyperparameter, with the performances being the worst compared to Random Forest Classifier and Logistic Regression.

2.2.2 Screenshots of Google Cloud Console Dataproc Job Log

RandomForestClassifier:

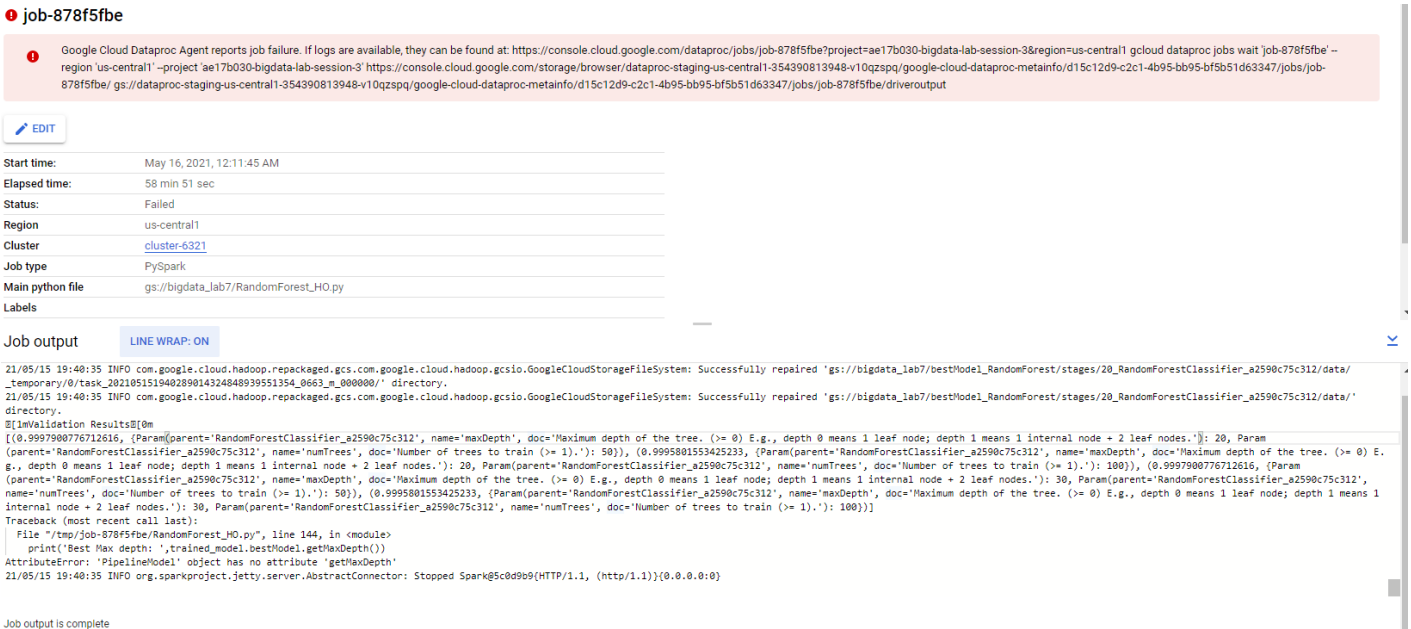


Figure 5: Screenshot showing completed training of RandomForestClassifier along with Validation Accuracy printed

Important Note :- The screenshot is showing job failure but actually this was due to an error in the print statement after the model training was completed.

LogisticRegression:



Figure 6: Screenshot showing completed training of LogisticRegression along with Validation Accuracy printed

NaiveBayes (Gaussian):

✔ job-3a495531

EDIT

Start time:	May 16, 2021, 2:54:44 AM
Elapsed time:	28 min 14 sec
Status:	Succeeded
Region	us-central1
Cluster	cluster-6321
Job type	PySpark
Main python file	gs://bigdata_lab7/NaiveBayes_HO.py
Labels	

EQUIVALENT REST

Job output

LINE WRAP: ON

```
21/05/15 21:52:46 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/18_StandardScaler_0f64231301f1/metadata/' directory.
21/05/15 21:52:48 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/18_StandardScaler_0f64231301f1/data/_temporary/0/task_202105152152463940856591456111345_0346_m_000000/' directory.
21/05/15 21:52:48 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/18_StandardScaler_0f64231301f1/data/' directory.
21/05/15 21:52:50 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/19_VectorAssembler_bfad8f6cd50/metadata/_temporary/0/task_202105152152483514461208583058921_0969_m_000000/' directory.
21/05/15 21:52:50 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/19_VectorAssembler_bfad8f6cd50/metadata/' directory.
21/05/15 21:52:52 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/20_NaiveBayes_6d5699067ff2/metadata/_temporary/0/task_202105152152503374743444237353590_0971_m_000000/' directory.
21/05/15 21:52:52 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/20_NaiveBayes_6d5699067ff2/metadata/' directory.
21/05/15 21:52:54 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/20_NaiveBayes_6d5699067ff2/data/_temporary/0/task_202105152152522059095823137562733_0351_m_000000/' directory.
21/05/15 21:52:54 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageFileSystem: Successfully repaired 'gs://bigdata_lab7/bestModel_NaiveBayes/stages/20_NaiveBayes_6d5699067ff2/data/' directory.
21/05/15 21:52:54 INFO org.apache.spark.executor.TaskExecutor: [Invalidation Results]
[[{"Param(parent='NaiveBayes_6d5699067ff2', name='smoothing', doc='The smoothing parameter, should be >= 0, default is 1.0'): 0.1}), {"Param(parent='NaiveBayes_6d5699067ff2', name='smoothing', doc='The smoothing parameter, should be >= 0, default is 1.0'): 0.5}), {"Param(parent='NaiveBayes_6d5699067ff2', name='smoothing', doc='The smoothing parameter, should be >= 0, default is 1.0'): 1.0})]]
21/05/15 21:52:54 INFO org.sparkproject.jetty.server.AbstractConnector: Stopped Spark@400d13064(HTTP/1.1, (http/1.1))[0.0.0.0:0]
```

Job output is complete

Figure 7: Screenshot showing completed training of NaiveBayes along with Validation Accuracy printed

3 Task 2: Real-time Computation

We read the data from the given Google Cloud Storage bucket and published it to the Kafka topic using publisher.py. Subscriber.py reads the data published to the Kafka topic and stores the data as a streaming Spark dataframe. We loaded the saved model stored in the GCS bucket and performed real-time predictions on it and measured the accuracy and F1-score batch-wise using the `foreachBatch` function.

3.1 Latency

In a period of 9.5 seconds, 100 rows were published to the Kafka topic ‘PROJECT’, out of which the results of 46 were printed to the console.

3.2 Real-time Prediction Results

Violation_County	label	prediction
NY	0.0	0.0
Q	2.0	2.0
K	1.0	1.0
NY	0.0	0.0
Q	2.0	2.0
K	1.0	1.0
K	1.0	1.0
NY	0.0	0.0
K	1.0	1.0
BX	3.0	3.0
K	1.0	1.0
NY	0.0	0.0
Q	2.0	2.0
Q	2.0	2.0
Q	2.0	2.0
Q	2.0	2.0
Q	2.0	2.0
NY	0.0	0.0
Q	2.0	2.0
K	1.0	1.0

Figure 8: Real-time predictions for Batch 1

Number of rows in batch 1 = 23
Accuracy on batch 1 is 1.0
F1 score on batch 1 is 1.0

Figure 9: Accuracy and F1-score for Batch 1

Violation_County	label	prediction
BX	3.0	3.0
K	1.0	1.0
NY	0.0	0.0
K	1.0	1.0
NY	0.0	0.0
NY	0.0	0.0
Q	2.0	2.0
NY	0.0	0.0
Q	2.0	2.0
K	1.0	1.0
NY	0.0	0.0
NY	0.0	0.0
NY	0.0	0.0
NY	0.0	0.0
Q	2.0	2.0
BX	3.0	3.0
NY	0.0	0.0
BX	3.0	3.0
NY	0.0	0.0
BX	3.0	3.0

Figure 10: Real-time predictions for Batch 2

Number of rows in batch 3 = 21
Accuracy on batch 3 is 1.0
F1 score on batch 3 is 1.0

Figure 11: Accuracy and F1-score for Batch 2

4 Inferences and Conclusion

- In this project, we have processed and analysed the massive-sized NYC Parking Ticket dataset using Apache PySpark by implementing supervised classification algorithms to predict the Violation County.
- Preprocessing the dataset is an extremely essential step prior to model training, due to the noisy and incomplete nature of the dataset.
- By implementing multiple classification models such as Random Forests, Logistic Regression and Gaussian Naive Bayes, and tuning the hyperparameters we observed that the best performing model was Random Forest Classifier with maximum depth 20 and number of trees 50 (with all other hyperparameters taking default values). The best model validation accuracy was 99.979%.
- Logistic Regression and Naive Bayes models yielded low validation accuracies.

- Given that the NYC dataset contains predominantly categorical and non-numeric features, it is evident that tree-based models like Decision Trees and Random Forests are well-suited for the task of predicting the violation County. On the other hand, Logistic Regression and probabilistic models like Naive Bayes cannot perform well unless the complexity is increased - for example, polynomial feature transformation.
- Training each class of algorithms (with hyperparameter tuning) took approximately 1 hour for the massively large training size of close to 40 million samples. This highlights the power of training models in the Cloud using Big Data pipelines.
- After saving the best model (Random Forest Classifier) in GCS bucket, we performed real-time prediction using Apache Kafka. Kafka producer is responsible for sending/publishing the rows of the dataset in batches as messages to a topic. Kafka consumer is responsible for retrieving the messages from the topic, thereby enabling real-time analysis.
- There is a high latency in the number of rows published to Kafka and the number of rows being processed. In our case, the latency seems to be caused primarily due to the Dataproc cluster and we can reduce this latency by scaling this cluster.

5 PySpark Codes

5.1 Pyspark Code for RandomForestClassifier Model

```
1 #Importing Requirements
2 from __future__ import print_function
3 from pyspark.context import SparkContext
4 from pyspark.ml.linalg import Vectors
5 from pyspark.ml.classification import RandomForestClassifier
6 from pyspark.sql.session import SparkSession
7 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
8 from pyspark.sql.types import DoubleType
9 from pyspark.ml.feature import StringIndexer, Imputer
10 from pyspark.ml import Pipeline
11 from pyspark.ml.feature import VectorAssembler
12 from pyspark.ml.feature import StandardScaler
13 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder, TrainValidationSplit
14 from pyspark.sql import functions as F
15
16 from pyspark.sql.types import IntegerType
17 from functools import reduce
18
19 sc = SparkContext()
20 spark = SparkSession(sc)
21
22 data_dir = "gs://bdl2021_final_project/nyc_tickets_train.csv"
23
24 df = spark.read.format("csv").option("recursiveFileLookup", "true").option("pathGlobFilter",
25     "*.csv").load(data_dir,inferSchema=True, header = True)
26
27 drop_cols = ['No Standing or Stopping Violation',
28     'Hydrant Violation', 'Double Parking Violation', 'Latitude',
29     'Longitude', 'Community Board', 'Community Council', 'Census Tract',
30     'BIN', 'BBL', 'NTA', 'Plate ID', 'Summons Number', 'Vehicle Expiration Date', 'Time
31     First Observed', 'Date First Observed', 'From Hours In Effect', 'To Hours In Effect', '
32     Unregistered Vehicle?', 'Days Parking In Effect', 'House Number', 'Violation Legal Code', '
33     Issuer Code']
34
35 req_cols = ["Registration State", "Plate Type", "Issue Date", "Violation Code", "Vehicle
36     Body Type", "Vehicle Make", "Issuing Agency", "Street Code1", "Street Code2", "Street
37     Code3",
38     "Issuer Command", "Issuer Squad", "Violation Time", "Violation In Front Of Or
39     Opposite", "Street Name", "Intersecting Street", "Law Section", "Sub Division", "
40     Vehicle Color",
41     "Vehicle Year", "Meter Number", "Feet From Curb", "Violation Post Code", "
42     Violation Description"]
43
44 drop_cols = list(map(lambda x: "_".join(x.split()), drop_cols))
45 req_cols = list(map(lambda x: "_".join(x.split()), req_cols))
46
47 oldColumns = df.schema.names
48 df = reduce(lambda df, idx: df.withColumnRenamed(oldColumns[idx], "_".join(oldColumns[idx].
49     split()))), range(len(oldColumns)), df)
50
51 df.drop(*drop_cols)
52
53 string_cols = req_cols.copy()
54 string_cols.remove("Violation_Code")
55 string_cols.remove("Street_Code1")
56 string_cols.remove("Street_Code2")
57 string_cols.remove("Street_Code3")
58 string_cols.remove("Vehicle_Year")
59 string_cols.remove("Feet_From_Curb")
60
61 #Remove Extra two columns
62 string_cols.remove("Issue_Date")
63 string_cols.remove("Violation_Time")
64
65 #Defining two udfs
66 @F.udf(returnType=IntegerType())
67 def dateParse(dstr):
68     BrStr = dstr.split("/")
69     NewBrStr = [BrStr[2], BrStr[0], BrStr[1]]
70     DateAsNo = int("".join(NewBrStr))
71     return DateAsNo
72
73 @F.udf(returnType=IntegerType())
74 def timeParse(tstr):
75     try:
76         NoPart = tstr[:-1]
77         dayLight = tstr[-1]
78         if(dayLight == 'A'):
79             return(int(NoPart))
80         else:
81             MM = NoPart[-2:]
82             HH = str(int(NoPart[:-2])%12+12)
83             TimeAsNo = int("".join([HH,MM]))
84             return TimeAsNo
85     except Exception:
86         return None
```

```

77
78 #Pass into both udf's
79 df = df.withColumn("Issue_Date", dateParse("Issue_Date"))
80 df = df.withColumn("Violation_Time", timeParse("Violation_Time"))
81
82 #Assign the Class Label column as label
83 indexer = StringIndexer(inputCol="Violation_County", outputCol="Violation_County_ind")
84 df = indexer.fit(df).transform(df)
85
86 df = df.withColumn("label", df["Violation_County_ind"])
87 df.select(F.col("Violation_Code").cast('int').alias("Violation_Code"))
88 df.select(F.col("Street_Code1").cast('int').alias("Street_Code1"))
89 df.select(F.col("Street_Code2").cast('int').alias("Street_Code2"))
90 df.select(F.col("Street_Code3").cast('int').alias("Street_Code3"))
91 df.select(F.col("Vehicle_Year").cast('int').alias("Vehicle_Year"))
92 df.select(F.col("Feet_From_Curb").cast('float').alias("Feet_From_Curb"))
93
94 # Label Encoding of Classes and removing the original column
95 indexers = [StringIndexer(inputCol=column, outputCol=column+"_ind").setHandleInvalid(value="
    skip") for column in string_cols]
96 req_cols = [col+"_ind" for col in string_cols] + [col for col in req_cols if col not in
    string_cols]
97
98 imputer = Imputer()
99 imputer.setStrategy("mode")
100 imputer.setInputCols(req_cols)
101 imputer.setOutputCols([col+"_imp" for col in req_cols])
102 req_cols = [col+"_imp" for col in req_cols]
103
104 assembler1 = VectorAssembler(inputCols = req_cols, outputCol = "features_old")
105 scaler = StandardScaler(inputCol="features_old", outputCol="scaledFeatures",withStd=True,
    withMean=True)
106 assembler2 = VectorAssembler(inputCols = ["scaledFeatures"], outputCol = "features")
107
108 # Construct a new Logistic Regression object and fit the training data.
109 rf = RandomForestClassifier()
110 #pipeline
111 pipe = Pipeline(stages = indexers + [imputer,assembler1,scaler,assembler2, rf])
112
113 # Create a grid of multiple values of the hyper-parameter regParam
114 paramGrid = ParamGridBuilder().addGrid(rf.maxDepth,[20,30]).addGrid(rf.numTrees,[50,100]).
    build()
115
116 #Create a CrossValidator Object
117 obj = TrainValidationSplit(estimator=pipe,
118     estimatorParamMaps=paramGrid,
119     evaluator=MulticlassClassificationEvaluator(metricName = '
    accuracy'),
120     trainRatio=0.9,
121     seed = 2021)
122
123 #Train the model with the CrossValidator Object
124 trained_model = obj.fit(df)
125
126 #SAVING...
127 trained_model.bestModel.save("gs://bigdata_lab7/bestModel_RandomForest")
128
129 # Acquire and print the best model details from the CrossValidator object
130 print("\033[1mValidation Results\033[0m")
131 val = list(zip(trained_model.validationMetrics, trained_model.getEstimatorParamMaps()))
132 print(val)

```

Listing 1: RandomForestClassifier using Pyspark

5.2 Pyspark Code for LogisticRegression Model

```
1 #Importing Requirements
2 from __future__ import print_function
3 from pyspark.context import SparkContext
4 from pyspark.ml.linalg import Vectors
5 from pyspark.ml.classification import LogisticRegression
6 from pyspark.sql.session import SparkSession
7 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
8 from pyspark.sql.types import DoubleType
9 from pyspark.ml.feature import StringIndexer, Imputer
10 from pyspark.ml import Pipeline
11 from pyspark.ml.feature import VectorAssembler
12 from pyspark.ml.feature import StandardScaler
13 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder, TrainValidationSplit
14 from pyspark.sql import functions as F
15
16 from pyspark.sql.types import IntegerType
17 from functools import reduce
18
19 sc = SparkContext()
20 spark = SparkSession(sc)
21
22 data_dir = "gs://bdl2021_final_project/nyc_tickets_train.csv"
23
24 df = spark.read.format("csv").option("recursiveFileLookup", "true").option("pathGlobFilter",
25     "*.csv").load(data_dir,inferSchema=True, header = True)
26
27 drop_cols = ['No Standing or Stopping Violation',
28     'Hydrant Violation', 'Double Parking Violation', 'Latitude',
29     'Longitude', 'Community Board', 'Community Council', 'Census Tract',
30     'BIN', 'BBL', 'NTA', 'Plate ID', 'Summons Number', 'Vehicle Expiration Date', 'Time
31     First Observed', 'Date First Observed', 'From Hours In Effect', 'To Hours In Effect', '
32     Unregistered Vehicle?', 'Days Parking In Effect', 'House Number', 'Violation Legal Code', '
33     Issuer Code']
34
35 req_cols = ["Registration State", "Plate Type", "Issue Date", "Violation Code", "Vehicle
36     Body Type", "Vehicle Make", "Issuing Agency", "Street Code1", "Street Code2", "Street
37     Code3",
38     "Issuer Command", "Issuer Squad", "Violation Time", "Violation In Front Of Or
39     Opposite", "Street Name", "Intersecting Street", "Law Section", "Sub Division", "
40     Vehicle Color",
41     "Vehicle Year", "Meter Number", "Feet From Curb", "Violation Post Code", "
42     Violation Description"]
43
44 drop_cols = list(map(lambda x: "_".join(x.split()), drop_cols))
45 req_cols = list(map(lambda x: "_".join(x.split()), req_cols))
46
47 oldColumns = df.schema.names
48 df = reduce(lambda df, idx: df.withColumnRenamed(oldColumns[idx], "_".join(oldColumns[idx].
49     split()))), range(len(oldColumns)), df)
50
51 df.drop(*drop_cols)
52
53 string_cols = req_cols.copy()
54 string_cols.remove("Violation_Code")
55 string_cols.remove("Street_Code1")
56 string_cols.remove("Street_Code2")
57 string_cols.remove("Street_Code3")
58 string_cols.remove("Vehicle_Year")
59 string_cols.remove("Feet_From_Curb")
60
61 #Remove Extra two columns
62 string_cols.remove("Issue_Date")
63 string_cols.remove("Violation_Time")
64
65 #Defining two udfs
66 @F.udf(returnType=IntegerType())
67 def dateParse(dstr):
68     BrStr = dstr.split("/")
69     NewBrStr = [BrStr[2], BrStr[0], BrStr[1]]
70     DateAsNo = int("".join(NewBrStr))
71     return DateAsNo
72
73 @F.udf(returnType=IntegerType())
74 def timeParse(tstr):
75     try:
76         NoPart = tstr[:-1]
77         dayLight = tstr[-1]
78         if(dayLight == 'A'):
79             return(int(NoPart))
80         else:
81             MM = NoPart[-2:]
82             HH = str(int(NoPart[:-2])%12+12)
83             TimeAsNo = int("".join([HH,MM]))
84             return TimeAsNo
85     except Exception:
86         return None
87
88 #Pass into both udf's
89 df = df.withColumn("Issue_Date", dateParse("Issue_Date"))
```

```

80 df = df.withColumn("Violation_Time", timeParse("Violation_Time"))
81
82 #Assign the Class Label column as label
83 indexer = StringIndexer(inputCol="Violation_County", outputCol="Violation_County_ind")
84 df = indexer.fit(df).transform(df)
85
86 df = df.withColumn("label", df["Violation_County_ind"])
87 df.select(F.col("Violation_Code").cast('int').alias("Violation_Code"))
88 df.select(F.col("Street_Code1").cast('int').alias("Street_Code1"))
89 df.select(F.col("Street_Code2").cast('int').alias("Street_Code2"))
90 df.select(F.col("Street_Code3").cast('int').alias("Street_Code3"))
91 df.select(F.col("Vehicle_Year").cast('int').alias("Vehicle_Year"))
92 df.select(F.col("Feet_From_Curb").cast('float').alias("Feet_From_Curb"))
93
94 # Label Encoding of Classes and removing the original column
95 indexers = [StringIndexer(inputCol=column, outputCol=column+"_ind").setHandleInvalid(value="
    skip") for column in string_cols]
96 req_cols = [col+"_ind" for col in string_cols] + [col for col in req_cols if col not in
    string_cols]
97
98 imputer = Imputer()
99 imputer.setStrategy("mode")
100 imputer.setInputCols(req_cols)
101 imputer.setOutputCols([col+"_imp" for col in req_cols])
102 req_cols = [col+"_imp" for col in req_cols]
103
104 assembler1 = VectorAssembler(inputCols = req_cols, outputCol = "features_old")
105 scaler = StandardScaler(inputCol="features_old", outputCol="scaledFeatures",withStd=True,
    withMean=True)
106 assembler2 = VectorAssembler(inputCols = ["scaledFeatures"], outputCol = "features")
107
108 # Construct a new Logistic Regression object and fit the training data.
109 lr = LogisticRegression(maxIter=25)
110 #pipeline
111 pipe = Pipeline(stages = indexers + [imputer,assembler1,scaler,assembler2, lr])
112
113 # Create a grid of multiple values of the hyper-parameter regParam
114 paramGrid = ParamGridBuilder().addGrid(lr.regParam,[0,0.5,1]).addGrid(lr.elasticNetParam
    ,[0,0.5]).build()
115
116 #Create a CrossValidator Object
117 obj = TrainValidationSplit(estimator=pipe,
118                             estimatorParamMaps=paramGrid,
119                             evaluator=MulticlassClassificationEvaluator(metricName = '
    accuracy'),
120                             trainRatio=0.9,
121                             seed = 2021)
122
123 #Train the model with the CrossValidator Object
124 trained_model = obj.fit(df)
125
126 #SAVING...
127 trained_model.bestModel.save("gs://bigdata_lab7/bestModel_LogisticRegression_FinTRY")
128
129 # Acquire and print the best model details from the CrossValidator object
130 print("\033[1mValidation Results\033[0m")
131 val = list(zip(trained_model.validationMetrics, trained_model.getEstimatorParamMaps()))
132 print(val)

```

Listing 2: LogisticRegression using Pyspark

5.3 Pyspark Code for NaiveBayes Model

```
1 #Importing Requirements
2 from __future__ import print_function
3 from pyspark.context import SparkContext
4 from pyspark.ml.linalg import Vectors
5 from pyspark.ml.classification import NaiveBayes
6 from pyspark.sql.session import SparkSession
7 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
8 from pyspark.sql.types import DoubleType
9 from pyspark.ml.feature import StringIndexer, Imputer
10 from pyspark.ml import Pipeline
11 from pyspark.ml.feature import VectorAssembler
12 from pyspark.ml.feature import StandardScaler
13 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder, TrainValidationSplit
14 from pyspark.sql import functions as F
15
16 from pyspark.sql.types import IntegerType
17 from functools import reduce
18
19 sc = SparkContext()
20 spark = SparkSession(sc)
21
22 data_dir = "gs://bdl2021_final_project/nyc_tickets_train.csv"
23
24 df = spark.read.format("csv").option("recursiveFileLookup", "true").option("pathGlobFilter",
25     "*.csv").load(data_dir,inferSchema=True, header = True)
26
27 drop_cols = ['No Standing or Stopping Violation',
28     'Hydrant Violation', 'Double Parking Violation', 'Latitude',
29     'Longitude', 'Community Board', 'Community Council', 'Census Tract',
30     'BIN', 'BBL', 'NTA', 'Plate ID', 'Summons Number', 'Vehicle Expiration Date', 'Time
31     First Observed', 'Date First Observed', 'From Hours In Effect', 'To Hours In Effect', '
32     Unregistered Vehicle?', 'Days Parking In Effect', 'House Number', 'Violation Legal Code', '
33     Issuer Code']
34
35 req_cols = ["Registration State", "Plate Type", "Issue Date", "Violation Code", "Vehicle
36     Body Type", "Vehicle Make", "Issuing Agency", "Street Code1", "Street Code2", "Street
37     Code3",
38     "Issuer Command", "Issuer Squad", "Violation Time", "Violation In Front Of Or
39     Opposite", "Street Name", "Intersecting Street", "Law Section", "Sub Division", "
40     Vehicle Color",
41     "Vehicle Year", "Meter Number", "Feet From Curb", "Violation Post Code", "
42     Violation Description"]
43
44 drop_cols = list(map(lambda x: "_".join(x.split()), drop_cols))
45 req_cols = list(map(lambda x: "_".join(x.split()), req_cols))
46
47 oldColumns = df.schema.names
48 df = reduce(lambda df, idx: df.withColumnRenamed(oldColumns[idx], "_".join(oldColumns[idx].
49     split()))), range(len(oldColumns)), df)
50
51 df.drop(*drop_cols)
52
53 string_cols = req_cols.copy()
54 string_cols.remove("Violation_Code")
55 string_cols.remove("Street_Code1")
56 string_cols.remove("Street_Code2")
57 string_cols.remove("Street_Code3")
58 string_cols.remove("Vehicle_Year")
59 string_cols.remove("Feet_From_Curb")
60
61 #Remove Extra two columns
62 string_cols.remove("Issue_Date")
63 string_cols.remove("Violation_Time")
64
65 #Defining two udfs
66 @F.udf(returnType=IntegerType())
67 def dateParse(dstr):
68     BrStr = dstr.split("/")
69     NewBrStr = [BrStr[2], BrStr[0], BrStr[1]]
70     DateAsNo = int("".join(NewBrStr))
71     return DateAsNo
72
73 @F.udf(returnType=IntegerType())
74 def timeParse(tstr):
75     try:
76         NoPart = tstr[:-1]
77         dayLight = tstr[-1]
78         if(dayLight == 'A'):
79             return(int(NoPart))
80         else:
81             MM = NoPart[-2:]
82             HH = str(int(NoPart[:-2])%12+12)
83             TimeAsNo = int("".join([HH,MM]))
84             return TimeAsNo
85     except Exception:
86         return None
87
88 #Pass into both udf's
89 df = df.withColumn("Issue_Date", dateParse("Issue_Date"))
```

```

80 df = df.withColumn("Violation_Time", timeParse("Violation_Time"))
81
82 #Assign the Class Label column as label
83 indexer = StringIndexer(inputCol="Violation_County", outputCol="Violation_County_ind")
84 df = indexer.fit(df).transform(df)
85
86 df = df.withColumn("label", df["Violation_County_ind"])
87 df.select(F.col("Violation_Code").cast('int').alias("Violation_Code"))
88 df.select(F.col("Street_Code1").cast('int').alias("Street_Code1"))
89 df.select(F.col("Street_Code2").cast('int').alias("Street_Code2"))
90 df.select(F.col("Street_Code3").cast('int').alias("Street_Code3"))
91 df.select(F.col("Vehicle_Year").cast('int').alias("Vehicle_Year"))
92 df.select(F.col("Feet_From_Curb").cast('float').alias("Feet_From_Curb"))
93
94 # Label Encoding of Classes and removing the original column
95 indexers = [StringIndexer(inputCol=column, outputCol=column+"_ind").setHandleInvalid(value="
    skip") for column in string_cols]
96 req_cols = [col+"_ind" for col in string_cols] + [col for col in req_cols if col not in
    string_cols]
97
98 imputer = Imputer()
99 imputer.setStrategy("mode")
100 imputer.setInputCols(req_cols)
101 imputer.setOutputCols([col+"_imp" for col in req_cols])
102 req_cols = [col+"_imp" for col in req_cols]
103
104 assembler1 = VectorAssembler(inputCols = req_cols, outputCol = "features_old")
105 scaler = StandardScaler(inputCol="features_old", outputCol="scaledFeatures",withStd=True,
    withMean=True)
106 assembler2 = VectorAssembler(inputCols = ["scaledFeatures"], outputCol = "features")
107
108 # Construct a new NaiveBayes object and fit the training data.
109 nb = NaiveBayes(modelType="gaussian")
110 #pipeline
111 pipe = Pipeline(stages = indexers + [imputer,assembler1,scaler,assembler2, nb])
112
113 # Create a grid of multiple values of the hyper-parameter regParam
114 paramGrid = ParamGridBuilder().addGrid(nb.smoothing,[0.1,0.5,1.0]).build()
115
116 #Create a CrossValidator Object
117 obj = TrainValidationSplit(estimator=pipe,
118     estimatorParamMaps=paramGrid,
119     evaluator=MulticlassClassificationEvaluator(metricName = '
        accuracy'),
120     trainRatio=0.9,
121     seed = 2021)
122
123 #Train the model with the CrossValidator Object
124 trained_model = obj.fit(df)
125
126 #SAVING...
127 trained_model.bestModel.save("gs://bigdata_lab7/bestModel_NaiveBayes")
128
129 # Acquire and print the best model details from the CrossValidator object
130 print("\033[1mValidation Results\033[0m")
131 val = list(zip(trained_model.validationMetrics, trained_model.getEstimatorParamMaps()))
132 print(val)

```

Listing 3: NaiveBayes using Pyspark

5.4 Code for Kafka Producer

```
1 from kafka import KafkaProducer
2 from google.cloud import storage
3
4 #kafka details
5 IP_ = "10.128.0.44:9092"
6 topic_name = 'PROJECT'
7 producer = KafkaProducer(bootstrap_servers = [IP_])
8
9 #download data for gcs
10 client = storage.Client()
11 bucket = client.get_bucket("bd12021_final_project")
12 blobs_all = list(bucket.list_blobs(prefix="nyc_tickets_train.csv/"))
13
14
15 i = 0
16 for blob in blobs_all[2:]:          #ignore the first 2 irrelevant files
17
18     #download data from GCS
19     content = blob.download_as_string()
20     content = content.decode('utf-8')
21
22     #split and iterate over each line
23     lines = content.split('\n')
24     for line in lines[1:]:          #ignore title line
25         if line != '':
26             line = bytes(line,'utf-8')
27             producer.send(topic_name,line)
28             producer.flush()
29
30     print("{} lines from part no: {} written to Kafka".format(len(lines),i))
31     i += 1
```

Listing 4: Kafka Producer

5.5 Code for Kafka Subscriber

```
1 #Importing Requirements
2 from __future__ import print_function
3 from pyspark.context import SparkContext
4 from pyspark.ml.linalg import Vectors
5 from pyspark.ml.classification import RandomForestClassifier
6 from pyspark.sql.session import SparkSession
7 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
8 from pyspark.sql.types import DoubleType
9 from pyspark.ml.feature import StringIndexer, StringIndexerModel, Imputer
10 from pyspark.ml import Pipeline, PipelineModel
11 from pyspark.ml.feature import VectorAssembler
12 from pyspark.ml.feature import StandardScaler
13 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
14 from pyspark.sql import functions as F
15
16 from pyspark.sql.types import IntegerType
17 from functools import reduce
18
19
20 spark = SparkSession.builder.appName("Project_NYC_tickets").getOrCreate()
21 spark.sparkContext.setLogLevel('WARN')
22
23 #address of Kafka server
24 IP_ = "10.128.0.48:9092"
25
26 #read from kafka
27 df = spark.readStream \
28     .format("kafka") \
29     .option("kafka.bootstrap.servers", IP_) \
30     .option("subscribe", "PROJECT") \
31     .load()
32
33
34 #All features
35 features = ['Summons Number', 'Plate ID', 'Registration State', 'Plate Type', 'Issue Date', '
36     Violation Code', 'Vehicle Body Type', 'Vehicle Make',
37     'Issuing Agency', 'Street Code1', 'Street Code2', 'Street Code3', 'Vehicle
38     Expiration Date', 'Issuer Code', 'Issuer Command', 'Issuer Squad',
39     'Violation Time', 'Time First Observed', 'Violation County', 'Violation In Front Of
40     Or Opposite', 'House Number', 'Street Name', 'Intersecting Street',
41     'Date First Observed', 'Law Section', 'Sub Division', 'Violation Legal Code', 'Days
42     Parking In Effect', 'From Hours In Effect', 'To Hours In Effect',
43     'Vehicle Color', 'Unregistered Vehicle?', 'Vehicle Year', 'Meter Number', 'Feet From
44     Curb', 'Violation Post Code', 'Violation Description',
45     'No Standing or Stopping Violation', 'Hydrant Violation', 'Double Parking
46     Violation', 'Latitude', 'Longitude', 'Community Board', 'Community Council',
47     'Census Tract', 'BIN', 'BBL', 'NTA']
48
49 features = list(map(lambda x: "_".join(x.split()), features))
50
51 #add feature columns to DataFrame
52 columns = F.split(df.value, ',')
53 for i in range(len(features)):
54     df = df.withColumn(features[i], columns[i].cast('string'))
55
56
57 #Defining two udfs
58 @F.udf(returnType=IntegerType())
59 def dateParse(dstr):
60     BrStr = dstr.split("/")
61     NewBrStr = [BrStr[2], BrStr[0], BrStr[1]]
62     DateAsNo = int("_".join(NewBrStr))
63     return DateAsNo
64
65 @F.udf(returnType=IntegerType())
66 def timeParse(tstr):
67     try:
68         NoPart = tstr[:-1]
69         dayLight = tstr[-1]
70         if(dayLight == 'A'):
71             return(int(NoPart))
72         else:
73             MM = NoPart[-2:]
74             HH = str(int(NoPart[:-2])%12+12)
75             TimeAsNo = int("_".join([HH, MM]))
76             return TimeAsNo
77     except Exception:
78         return None
79
80 #Pass into both udf's
81 df = df.withColumn("Issue_Date", dateParse("Issue_Date"))
82 df = df.withColumn("Violation_Time", timeParse("Violation_Time"))
83
84 #cast columns to appropriate dtype
85 df = df.withColumn("Violation_Code", F.col("Violation_Code").cast('int'))
86 df = df.withColumn("Street_Code1", F.col("Street_Code1").cast('int'))
87 df = df.withColumn("Street_Code2", F.col("Street_Code2").cast('int'))
88 df = df.withColumn("Street_Code3", F.col("Street_Code3").cast('int'))
```

```

84 df = df.withColumn("Vehicle_Year", F.col("Vehicle_Year").cast('int'))
85 df = df.withColumn("Feet_From_Curb", F.col("Feet_From_Curb").cast('float'))
86
87
88 #load label indexer and transform df
89 indexer = StringIndexerModel.load('gs://avinashbagali/Label_Indexer/')
90 df = indexer.transform(df)
91 df = df.withColumn("label", df["Violation_County_ind"])
92
93
94 #load best model and transform df
95 best_model = PipelineModel.load("gs://avinashbagali/bestModel_RandomForest/")
96 df = best_model.transform(df)
97
98 #result df -> ["Violation_County","label","prediction"]
99 result_df = df[["Violation_County","label","prediction"]]
100
101
102 # function to use with foreachBatch to compute accuracy and F1-score
103 def Metrics(df, epoch_id):
104     evaluator_acc = MulticlassClassificationEvaluator(predictionCol = "prediction",
105                                                         labelCol="label",
106                                                         metricName="accuracy")
107
108     evaluator_f1 = MulticlassClassificationEvaluator(predictionCol = "prediction",
109                                                         labelCol="label",
110                                                         metricName="f1")
111
112     acc = evaluator_acc.evaluate(df)
113     f1 = evaluator_f1.evaluate(df)
114
115
116     print("+-----+")
117     print("|Number of rows in batch {} = {}|".format(epoch_id,df.count()))
118     print("| Accuracy on batch {} is {} |".format(epoch_id,round(acc,3)))
119     print("| F1 score on batch {} is {} |".format(epoch_id,round(f1,3)))
120     print("+-----+")
121
122
123 #print result_df -> ["Violation_County","label","prediction"]
124 query1 = result_df \
125     .writeStream \
126     .outputMode("append") \
127     .format("console") \
128     .start()
129
130
131 #print Accuracy and F1-score
132 query2 = result_df \
133     .writeStream \
134     .format("console") \
135     .foreachBatch(Metrics) \
136     .start()
137
138
139 query1.awaitTermination()
140 query2.awaitTermination()

```

Listing 5: Kafka Subscriber