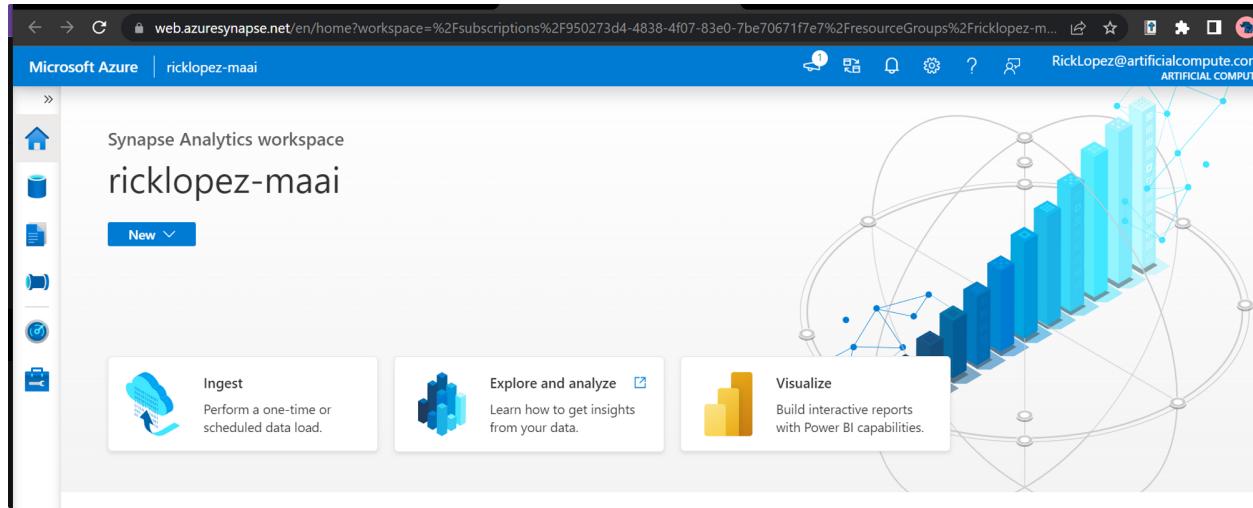


Assignment 7 Screenshots - Richard Lopez

The goal of this assignment is to utilize Azure Synapse integrated Machine Learning capabilities to predict an outcome.



Download and access the dataset

PC > Blade 14 (C:) > Users > RickL > Dropbox > msaa > BigData > wk7

Name	Date modified	Type	Size
NYCTripSmall.parquet	2/26/2023 11:27 PM	PARQUET File	5,047 KB

Load the dataset into Azure Synapse

A screenshot of the Microsoft Azure Synapse Data workspace interface. The top navigation bar shows the URL 'web.azuresynapse.net/en/authoring/explore/linked/storageaccounts/ricklopez-maai-WorkspaceDefaultStorage-ricklopezmaai%2Fricklopezmaai%2F...'. The left sidebar shows 'Data' with 'Linked' selected, and a list of storage accounts: 'Azure Data Lake Storage Gen2' (2 items) and 'ricklopez-maai (Primary - ricklopezmaai Primary)' (1 item). The main area shows a 'Sales Bulk Upload' section with options like 'New SQL script', 'New data flow', 'New integration dataset', 'Upload', 'Download', and 'New folder'. A file browser window is open over the interface, showing a folder structure under 'ricklopezmaai': '_rejectedrows' (Folder, Last Modified: 2/25/2023, 11:16:43 PM) and 'SalesInfo.csv' (File, Last Modified: 2/18/2023, 4:20:22 PM).

Microsoft Azure | ricklopez-maai

Synapse live | Validate all | Publish all | Search

Data

- Workspace
- Linked

Filter resources by name

- Azure Data Lake Storage Gen2 2
 - ricklopez-maai (Primary - ricklopezmaai)
 - ricklopezmaai (Primary)
- (Attached Containers) ...

Sales Bulk Upload ricklopezmaai

New SQL script New data flow New integration dataset

ricklopezmaai

Name Last Modified

_rejectedrows	2/25/2023, 11:16:43 PM
SalesInfo.csv	2/18/2023, 4:20:22 PM

File Upload NYCTripSmall.parquet

Overwrite existing files

File name	Size	Action
NYCTripSmall.parquet	4.93 MB	Remove

Showing 1 to 2 of 2 cached items

Upload Cancel

ARTIFICIAL COMPUTE

Synapse live | Validate all | Publish all | Search

Data

- Workspace
- Linked

Filter resources by name

- Azure Data Lake Storage Gen2 2
 - ricklopez-maai (Primary - ricklopezmaai)
 - ricklopezmaai (Primary)
- (Attached Containers) ...

Sales Bulk Upload ricklopezmaai

New SQL script New data flow New integration dataset

ricklopezmaai

Name Last Modified Content Type Size

_rejectedrows	2/25/2023, 11:16:43 PM	Folder	
NYCTripSmall.parquet	2/26/2023, 11:51:54 PM		4.9 MB
SalesInfo.csv	2/18/2023, 4:20:22 PM		364 B

Upload completed

1 files uploaded to 'ricklopezmaai/'. View detail

Synapse Analytics | ricklopez-maai

live | Validate all | Publish all | Search

Data

- Bulk Upload
- script 1
- script 2
- script 3
- script 4
- script 5
- script 6
- script 7
- script 8
- script 9
- script 10
- script 11

Sales Bulk Upload ricklopezmaai SQL script 11

Run Undo Publish Query plan Connect to Built-in Use database master

```

1 -- This is auto-generated code
2 SELECT
3     TOP 100 *
4 FROM
5 OPENROWSET(
6     BULK 'https://ricklopezmaai.dfs.core.windows.net/ricklopezmaai/NYCTripSmall.parquet',
7     FORMAT = 'PARQUET'
8 ) AS [result]
9

```

Properties

General Related (0)

Name * SQL script 11

Description

Type .sql script

Size 213 bytes

Results settings per query
 First 5000 rows (default)
 All rows

Results Messages

View Table Chart Export results

Search

DateID	MedallionID	HackneyLicens...	PickupTimeID	DropoffTimeID	PickupGeogra...	DropoffGeogr...	Picku...
20131231	10624	35117	63879	64392	172348	219532	40.81...
20131231	7263	11597	61733	62114	87083	137768	40.73...
20131231	1899	5301	59488	59815	82160	146431	40.76...
20131231	9171	28727	60002	60518	276312	12408	40.75...

Microsoft Azure | Synapse Analytics | ricklopez-maai

Synapse live | Validate all | Publish all | Search

Develop

- SQL scripts 12

Filter resources by name

- Sales Bulk Upload
- SQL script 1
- SQL script 2
- SQL script 3
- SQL script 4
- SQL script 5
- SQL script 6
- SQL script 7

Sales Bulk Upload ricklopezmaai SQL script 11

New SQL script New notebook New data flow New integration dataset

Load to DataFrame

Name Last Modified Content Type Size

New Spark table			
_rejectedrows	2/25/2023, 11:16:43 PM	Folder	
NYCTripSmall.parquet	2/26/2023, 11:51:54 PM		4.9 MB
SalesInfo.csv	2/18/2023, 4:20:22 PM		364 B

Microsoft Azure

Home > Microsoft.Azure.SynapseAnalytics.SparkPool-202302270005342506 | Overview

Deployment

Search Delete Cancel Redeploy Download Refresh

Your deployment is complete

Deployment name: Microsoft.Azure.SynapseAnalytics.SparkPool-20... Start time: 2/27/2023, 12:05:32 AM

Subscription: Microsoft Azure Sponsorship Correlation ID: ccbdfc25-ae07-421d-a896-f006c128ddef

Resource group: ricklopez-maai

Deployment details Next steps Go to resource

Microsoft Azure | Synapse Analytics > ricklopez-maai

Synapse live Validate all Publish all

Data Workspace Linked Filter resources by name

Azure Data Lake Storage Gen2 ricklopez-maai (Primary - ricklo... ricklopezmaai (Primary) (Attached Containers)

Sales Bulk Upload ricklopezmaai SQL script 11 Notebook 1 Run all Undo Publish Outline Attach to ricklopezmaai Language PySpark (Python) Variables

[1] Ready

```
1 %pyspark
2 df = spark.read.load('abfss://ricklopezmaai@ricklopezmaai.dfs.core.windows.net/NYCTripSmall.parquet', format='parquet')
3 display(df.limit(10))
```

3 min 31 sec - Apache Spark session started in 2 min 53 sec 939 ms. Command executed in 37 sec 648 ms by RickLopez on 12:10:17 AM, 2/27/23

Job execution Succeeded Spark 2 executors 8 cores View in monitoring Open Spark UI

DateID	MedallionID	HackneyLicenseID	PickupTimeID	DropoffTimeID
20131231	10624	35117	63879	64392
20131231	7263	11597	61733	62114
20131231	1899	5301	59488	59815
20131231	9471	28727	69093	69548
20131231	2556	16124	51720	52440
20131231	1433	33809	31260	31100

💡 See Full notebook here: <https://ricklopez.github.io/AI-540-01-SP23-Applications-of-AI/wk7/richard-lopez.html>

ricklopez.github.io

<https://ricklopez.github.io/AI-540-01-SP23-Applications-of-AI/wk7/richard-lopez.html>

Assignment 7.1: Hands-on with Azure Synapse

Richard Lopez

```
In [1]: %%pyspark
df = spark.read.load('abfss://ricklopezmaai@ricklopezmaai.dfs.core.windows.net/NYCTripSmall.parquet', format='parquet')
display(df.limit(10))
```

StatementMeta(ricklopezmaai, 0, 1, Finished, Available)

SynapseWidget(Synapse.DataFrame, dcc47129-9079-4be1-8a5b-ab146bf895fc)

Import Libs

```
In [3]: import matplotlib.pyplot as plt
from datetime import datetime
from dateutil import parser
from pyspark.sql.functions import unix_timestamp, date_format, col, when
from pyspark.ml import Pipeline
from pyspark.ml import PipelineModel
from pyspark.ml.feature import RFormula
from pyspark.ml.feature import OneHotEncoder, StringIndexer, VectorIndexer
from pyspark.ml.classification import LogisticRegression
from pyspark.mllib.evaluation import BinaryClassificationMetrics
from pyspark.ml.evaluation import BinaryClassificationEvaluator
```

```
from pyspark.mllib.evaluation import BinaryClassificationMetrics
from pyspark.ml.evaluation import BinaryClassificationEvaluator
```

StatementMeta(ricklopezmaai, 0, 3, Finished, Available)

Conduct feature preprocessing (if needed)

```
[4]: # build DF
from azureml.opendatasets import NycTlcYellow
```

```
end_date = parser.parse('2018-06-06')
start_date = parser.parse('2018-05-01')
nyc_tlc = NycTlcYellow(start_date=start_date, end_date=end_date)
filtered_df = nyc_tlc.to_spark_dataframe()
```

StatementMeta(ricklopezmaai, 0, 4, Finished, Available)

```
[6]: # Sample data for faster initial development.
sampled_taxi_df = filtered_df.sample(True, 0.002, seed=1234)
```

StatementMeta(ricklopezmaai, 0, 6, Finished, Available)

```
[7]: display(sampled_taxi_df)
```

StatementMeta(ricklopezmaai, 0, 7, Finished, Available)

SynapseWidget(Synapse.DataFrame, bee7e0ad-d3f5-49cd-8a3d-4d7deaf502c2)

Complete Feature Processing

```
In [8]: taxi_df = sampled_taxi_df.select('totalAmount', 'fareAmount', 'tipAmount', 'paymentType', 'rateCodeId', 'passengerCount'\
, 'tripDistance', 'tpepPickupDateTime', 'tpepDropoffDateTime'\
, date_format('tpepPickupDateTime', 'hh').alias('pickupHour')\
, date_format('tpepDropoffDateTime', 'EEE').alias('weekdayString')\
, (unix_timestamp(col('tpepDropoffDateTime')) - unix_timestamp(col('tpepPickupDateTime'))).alias('tripTimeSecs')\
, (when(col('tipAmount') > 0, 1).otherwise(0)).alias('tipped')\
)\\
.filter((sampled_taxi_df.passengerCount > 0) & (sampled_taxi_df.passengerCount < 8) \
& (sampled_taxi_df.tipAmount > 0) & (sampled_taxi_df.tipAmount <= 25) \
& (sampled_taxi_df.fareAmount >= 1) & (sampled_taxi_df.fareAmount <= 250) \
& (sampled_taxi_df.tipAmount < sampled_taxi_df.fareAmount) \
& (sampled_taxi_df.tripDistance > 0) & (sampled_taxi_df.tripDistance <= 100) \
& (sampled_taxi_df.rateCodeId <= 5) \
& (sampled_taxi_df.paymentType.isin(["1", "2"])) \
)

StatementMeta(ricklopezmaai, 0, 8, Finished, Available)

In [9]: taxi_featurised_df = taxi_df.select('totalAmount', 'fareAmount', 'tipAmount', 'paymentType', 'passengerCount'\
, 'tripDistance', 'weekdayString', 'pickupHour', 'tripTimeSecs', 'tipped'\
, when((taxi_df.pickupHour <= 6) | (taxi_df.pickupHour >= 20), "Night")\
, when((taxi_df.pickupHour >= 7) & (taxi_df.pickupHour <= 10), "AMRush")\
, when((taxi_df.pickupHour >= 11) & (taxi_df.pickupHour <= 15), "Afternoon")\
, when((taxi_df.pickupHour >= 16) & (taxi_df.pickupHour <= 19), "PMRush")\
, otherwise(0).alias("trafficTimeBins")\
)\\
.filter((taxi_df.tripTimeSecs >= 30) & (taxi_df.tripTimeSecs <= 7200))

StatementMeta(ricklopezmaai, 0, 9, Finished, Available)
```

Define the model

```
In [10]: # convert the categorical columns into numbers
sI1 = StringIndexer(inputCol="trafficTimeBins", outputCol="trafficTimeBinsIndex")
en1 = OneHotEncoder(dropLast=False, inputCol="trafficTimeBinsIndex", outputCol="trafficTimeBinsVec")
sI2 = StringIndexer(inputCol="weekdayString", outputCol="weekdayIndex")
en2 = OneHotEncoder(dropLast=False, inputCol="weekdayIndex", outputCol="weekdayVec")

encoded_final_df = Pipeline(stages=[sI1, en1, sI2, en2]).fit(taxi_featurised_df).transform(taxi_featurised_df)

StatementMeta(ricklopezmaai, 0, 10, Finished, Available)
```

Build the pipeline and train the model

```
In [11]: # Create training and testing data split from the DataFrame
trainingFraction = 0.7
testingFraction = (1-trainingFraction)
seed = 1234

train_data_df, test_data_df = encoded_final_df.randomSplit([trainingFraction, testingFraction], seed=seed)

StatementMeta(ricklopezmaai, 0, 11, Finished, Available)

In [12]: # Logistic regression object for the model
logReg = LogisticRegression(maxIter=10, regParam=0.3, labelCol = 'tipped')

classFormula = RFormula(formula="tipped ~ pickupHour + weekdayVec + passengerCount + tripTimeSecs + tripDistance + fareAmount + paymentType+ trafficTimeBinsIndex + weekdayIndex")

# training pipeline
lrModel = Pipeline(stages=[classFormula, logReg]).fit(train_data_df)

datestamp = datetime.now().strftime('%m-%d-%Y-%S')
fileName = "lrModel_" + datestamp
logRegDirfilename = fileName
lrModel.save(logRegDirfilename)

StatementMeta(ricklopezmaai, 0, 12, Finished, Available)

In [17]: predictions = lrModel.transform(test_data_df)

StatementMeta(ricklopezmaai, 0, 17, Finished, Available)

In [19]: from pyspark.ml.evaluation import MulticlassClassificationEvaluator
evaluator = MulticlassClassificationEvaluator( \
    labelCol="tipped", \
    predictionCol="prediction", \
    metricName="accuracy")

StatementMeta(ricklopezmaai, 0, 19, Finished, Available)

In [20]: lr_accuracy = evaluator.evaluate(predictions)
print("Accuracy of LogisticRegression is %%" % lr_accuracy)
print("Test Error of LogisticRegression = %%" % (1.0 - lr_accuracy))

StatementMeta(ricklopezmaai, 0, 20, Finished, Available)
Accuracy of LogisticRegression is 0.970428
Test Error of LogisticRegression = 0.0295725
```

Make predictions and evaluate model performance

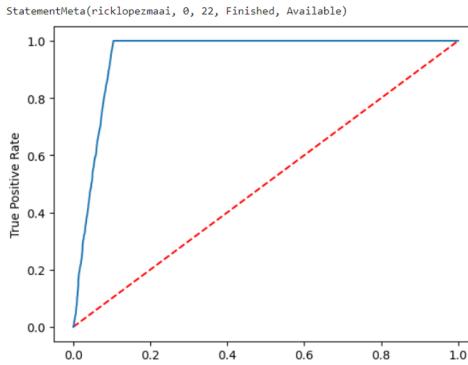
```
In [21]: predictionAndLabels = predictions.select("label","prediction").rdd
metrics = BinaryClassificationMetrics(predictionAndLabels)
print("Area under ROC = %s" % metrics.areaUnderROC)

StatementMeta(ricklopezmaai, 0, 21, Finished, Available)
/opt/spark/python/lib/pyspark.zip/pyspark/sql/context.py:157: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
Area under ROC = 0.9786493385936412
```

```
Arced under ROC = 0.9786493385936412

In [22]: ## Plot the ROC curve one of the ways to reason about the results of this Model test
modelSummary = lrModel.stages[1].summary

plt.plot([0, 1], [0, 1], 'r--')
plt.plot(modelSummary.roc.select('FPR').collect(),
         modelSummary.roc.select('TPR').collect())
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.show()
```



References

- Chakraborty, P. (2021, January 12). 9 Classification Methods From Spark MLlib We Should Know. Medium; Medium. <https://cprosenjit.medium.com/9-classification-methods-from-spark-mllib-we-should-know-c41f55c0425>
- AAI-540-01-SP23 - Applications of AI Module 7 Videos (2023). Sandiego.edu. https://ole.sandiego.edu/ultra/courses/_107930_1/cl/outline
- saveenr. (2022, November 18). Quickstart: Get started analyzing with Spark - Azure Synapse Analytics. Microsoft.com. <https://learn.microsoft.com/en-us/azure/synapse-analytics/get-started-analyze-spark>