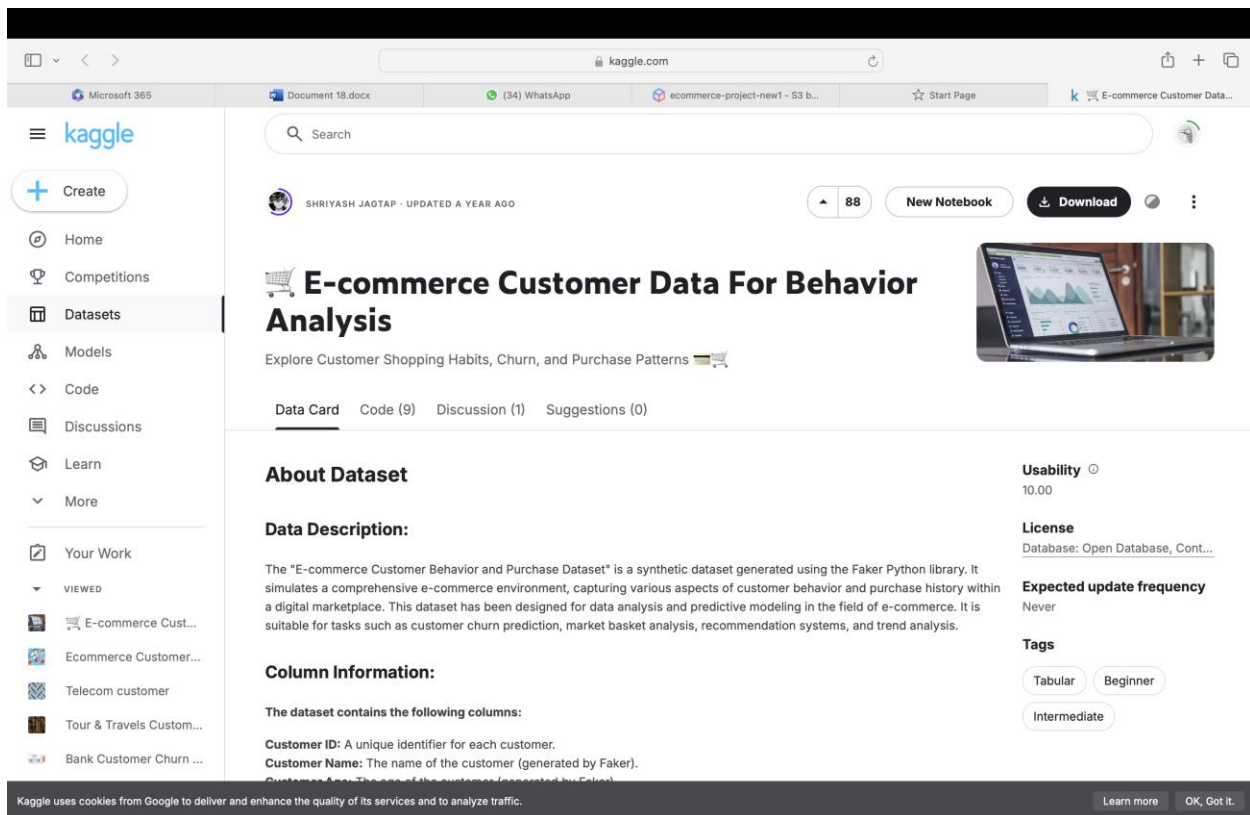
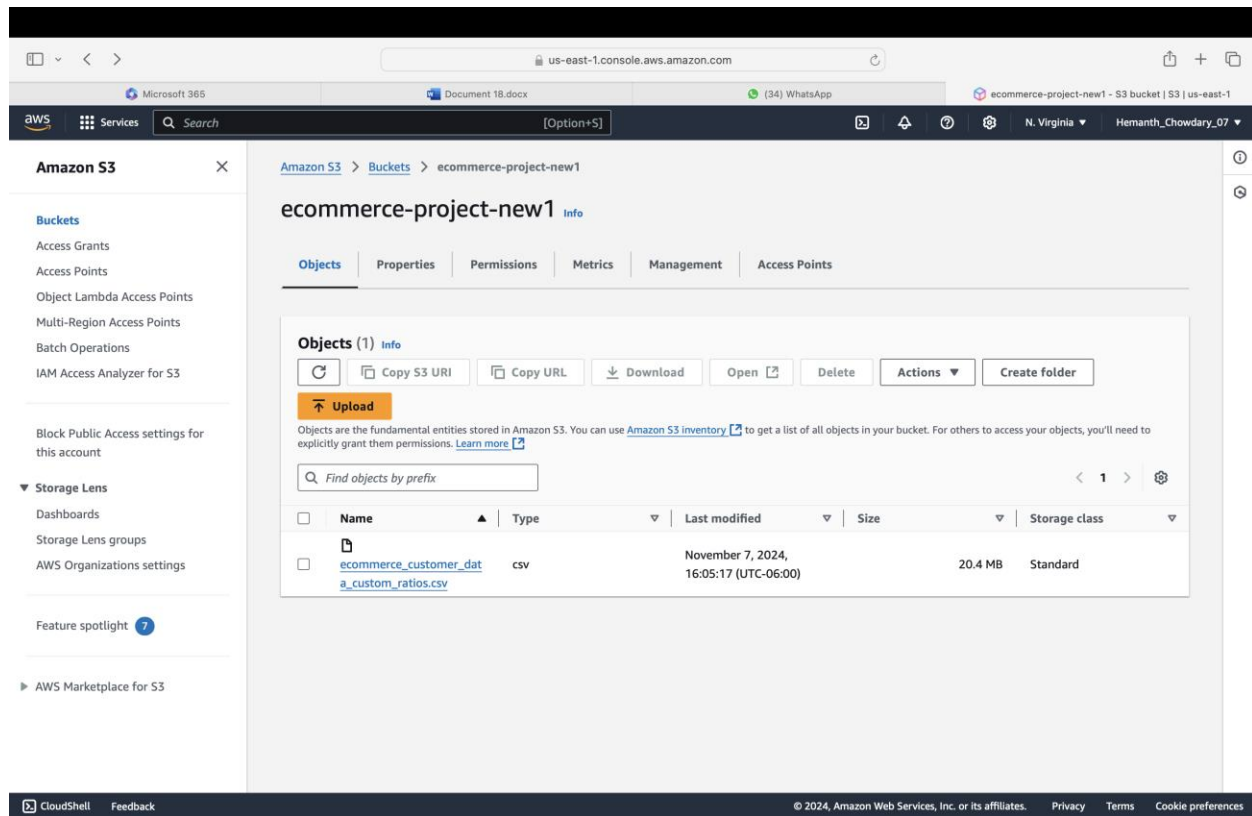


STEPS:

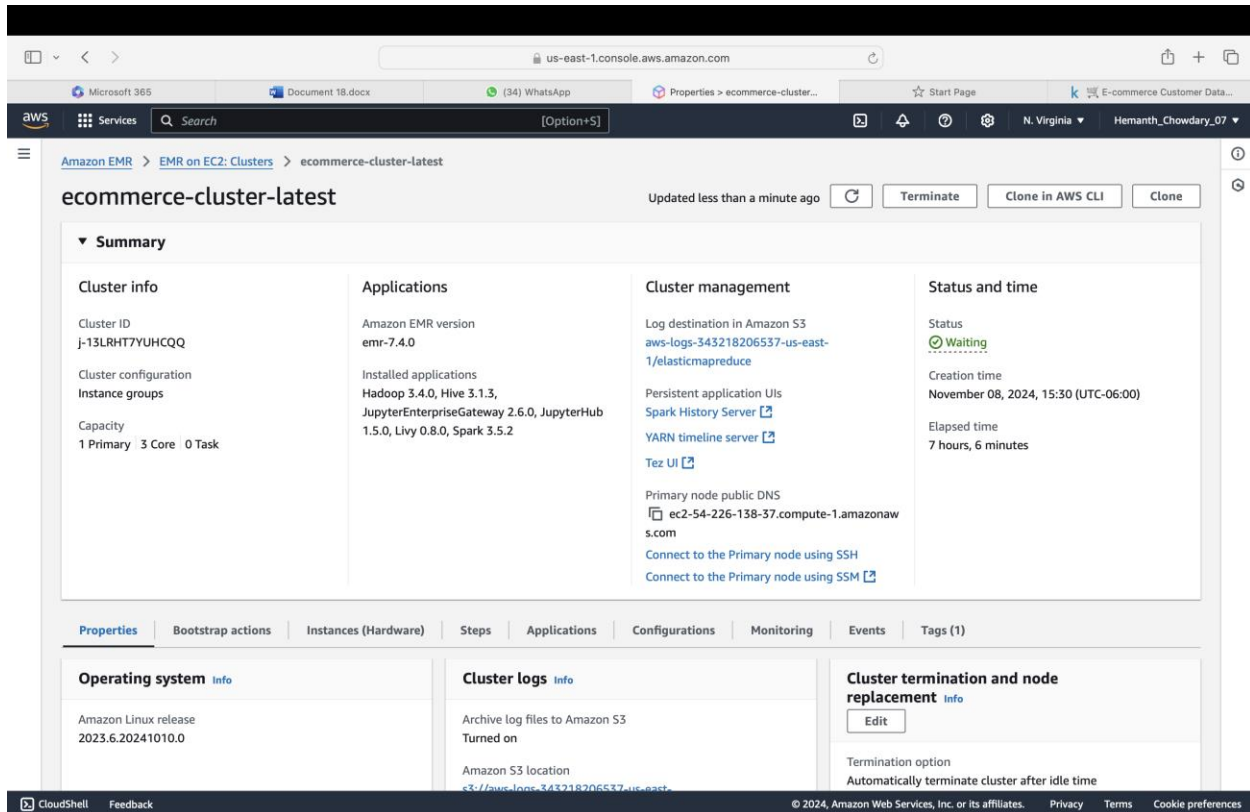
1. First we will download the data set from Kaggle and then we will store it on our system.



2. Then we will login into our AWS account and then create S3 bucket with project name and upload our data set init.



3. Then we will create a Cluster in EMR with all the permissions, key pairs and policies and run the cluster.



4. Then we will Copy the DNS and we will run EMR in our terminal using the DNS id and Keypair we used in EMR cluster.
5. Then we will start Hive in EMR.
6. Then we will Load the data set which is in S3 bucket by giving its path and create a table with that data.
7. Then we will analyze the data in HIVE by running some queries to know the properties and structure of our data.

```
[root@hadoop~]# hadoop@ip-172-31-77-105:-  
E:::EEEEEEEEE M:::M M:::R R:::  
EE:::EEEEEEEEE E M:::M M:::M R:::RRRRRR:::R  
E::::E EEEEE M:::M M:::M RR:::R R:::R  
E::::E EEEEE M:::M M:::M M:::M R:::R R:::R  
E::::EEEEEEEEE M:::M M:::M M:::M R:::RRRRRR:::R  
E::::EEEEEEEEE M:::M M:::M M:::M R:::RRRRRR:::R  
E::::EEEEEEEEE M:::M M:::M M:::M R:::RRRRRR:::R  
E::::E EEEEE M:::M M:::M R:::R R:::R  
E::::E EEEEE M:::M WWW M:::M R:::R R:::R  
EE:::EEEEEEEEE E M:::M M:::M R:::R R:::R  
E::::EEEEEEEEE E M:::M M:::M RR:::R R:::R  
EEEEEEEEEEEEEEEEEE WWWWWW WWWWRR RRRRRR
```

```
[hadoop@ip-172-31-77-105 ~]$ hive  
Hive Session ID = eaec179c-42c3-4b25-a7e6-f48bbc377b74  
  
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: true  
returns INT, customer_name STRING, age INT, gender STRING, churn STRING)row format delimited fields terminated by ',' stored as textfile;  
OK  
Time taken: 0.679 seconds  
hive> load data inpah 's3://hemant-task1/ecommerce_customer_data_custom_ratios.csv' into table customer_purchases;  
Loading data to table default.customer_purchases  
OK  
Time taken: 3.554 seconds  
hive> show tables;  
FAIL!D: ParseException line 2:5 extraneous input 'tables' expecting EOF near '<EOF>'  
hive> show tables;  
OK  
customer_purchases  
Time taken: 0.082 seconds, Fetched: 1 row(s)  
hive> describe customer_purchases;  
OK  
+-----+  
customer_id      string  
purchase_date    string  
product_category string  
product_price    double  
quantity         int  
total_purchase_amount double  
payment_method   string  
customer_age     int  
returns          int  
customer_name    string  
age              int  
gender           string  
churn            string  
Time taken: 0.04 seconds, Fetched: 13 row(s)  
hive>
```

```

hadoop@ip-172-31-77-105:~$
2738 2023-01-03 01:55:29 Clothing 42.0 1 1294.0 PayPal 22 0 Melissa Olson 22 Female 0
47364 2023-01-04 02:07:19 Home 122.0 2 2778.0 PayPal 63 NULL Jeremy Thomas 63 Female 0
23662 2023-04-04 10:07:34 Clothing 192.0 3 4159.0 PayPal 57 NULL Sean Spencer 57 Male 0
870 2020-05-13 01:56:05 Clothing 14.0 1 1588.0 Credit Card 66 0 Greg Jackson 66 Male 0
8730 2020-05-04 08:04:09 Electronics 22.0 4 1463.0 Credit Card 66 NULL Greg Jackson 66 Male 0
19836 2020-06-29 02:23:24 Electronics 287.0 3 2244.0 Credit Card 25 1 Jordan Hebert 25 Female 0
19836 2023-05-30 23:50:44 Home 418.0 4 3283.0 Credit Card 25 NULL Jordan Hebert 25 Female 0
5927 2023-08-21 03:36:59 Books 88.0 5 3228.0 PayPal 51 1 Edward Foster 51 Female 0
8511 2020-08-20 08:43:43 Electronics 43.0 3 2821.0 Credit Card 21 NULL Christopher Wood 21 Male 1
31680 2020-02-09 09:23:37 Home 77.0 1 4712.0 Credit Card 46 0 Lauren Ward 46 Male 0
31680 2020-08-13 16:26:31 Books 97.0 2 4863.0 PayPal 46 0 Lauren Ward 46 Male 0
44329 2020-07-22 18:36:27 Books 163.0 2 1891.0 PayPal 61 0 Jared Campbell 61 Male 0
228 2020-01-10 20:42:21 Home 162.0 5 4388.0 PayPal 22 0 Alyssa Watkins 22 Female 0
1247 2020-09-28 09:34:57 Clothing 236.0 2 1472.0 Credit Card 67 1 Ann Boyd 67 Female 0
28519 2020-12-23 16:32:45 Books 312.0 1 1147.0 PayPal 48 NULL Tara Thomas 48 Female 0
15196 2023-04-13 05:07:24 Electronics 227.0 4 3805.0 Credit Card 45 0 Andrew Patterson 45 Male 0
48895 2022-02-02 10:57:45 Electronics 170.0 2 2442.0 Credit Card 46 1 Chris Delgado 46 Female 0
9182 2023-04-01 21:07:58 Books 35.0 5 3263.0 Credit Card 21 0 Kristyn Lynn 21 Female 0
19531 2022-10-10 09:23:37 Clothing 5 1782.0 Credit Card 66 0 Stephen Vasquez 66 Female 0
19336 2022-11-20 19:59:29 Books 229.0 5 3474.0 Credit Card 56 0 Joshua Sanchez 56 Male 0
12886 2023-01-01 04:24:51 Clothing 193.0 1 2067.0 Cash 28 0 Mrs. Susan Monroe 28 Female 0
40179 2023-03-13 00:04:05 Electronics 63.0 4 3986.0 Credit Card 32 1 Dale Meadows 32 Male 1
906 2020-11-11 08:26:28 Clothing 171.0 3 2250.0 Credit Card 24 0 Margaret Wise 24 Female 0
33390 2022-11-19 22:28:23 Home 347.0 1 5089.0 PayPal 40 0 Rachel Martinez 40 Female 0
41832 2022-01-01 22:30:30 Home 333.0 1 2320.0 Credit Card 46 0 Matthew Orr 46 Male 0
31267 2022-10-24 14:59:34 Books 50.0 1 2308.0 Credit Card 48 0 Harry Wood 48 Male 1
892 2020-04-16 17:13:17 Electronics 321.0 1 2968.0 PayPal 58 1 Wendy Hernandez 58 Female 0
20049 2022-09-13 13:45:15 Clothing 81.0 2 3989.0 Credit Card 62 Lisa Johnson 62 Female 0
86415 2020-04-14 20:07:20 Home 147.0 5 1459.0 Cash 20 NULL Anthony McGrath 20 Male 1
6943 2023-09-01 00:12:41 Home 207.0 1 2968.0 Credit Card 58 1 Brian Barnes 58 Male 0
29504 2023-03-02 13:49:13 Home 124.0 5 2997.0 PayPal 58 1 Elizabeth Rhodes 58 Male 1
37382 2022-08-16 20:06:20 Books 406.0 2 4086.0 Credit Card 38 NULL Nancy Ortiz 38 Female 0
9093 2023-05-01 11:11:28 Clothing 111.0 4 4599.0 PayPal 36 0 Lisa Beard 36 Male 0
22888 2022-08-07 09:45:53 Electronics 169.0 1 4042.0 Credit Card 37 0 Melinda Wright 37 Female 0
25420 2020-07-15 04:55:39 Home 99.0 4 1665.0 Cash 26 1 Thomas Goodman 26 Male 0
23536 2021-09-17 02:50:11 Electronics 87.0 1 4073.0 Credit Card 70 NULL Heather Smith 70 Female 0
13612 2022-01-13 18:01:28 Books 116.0 4 4031.0 Cash 45 0 Katelyn Hebert 45 Female 0
977 2021-05-23 12:35:41 Books 24.0 1 2864.0 Credit Card 34 1 NULL Marcus Sanchez 34 Female 0
17165 2020-02-25 13:38:16 Clothing 230.0 4 3664.0 Credit Card 18 0 Dawn Perez 18 Male 0
45397 2022-02-18 04:18:18 Books 95.0 2 3397.0 Cash 54 NULL Scott Lindsey 54 Male 0
35965 2022-10-13 19:48:58 Home 83.0 5 1335.0 Credit Card 19 NULL Shelby Williams 19 Female 0
45410 2021-05-30 15:37:15 Home 311.0 2 3302.0 Credit Card 50 1 Johnny Riley 50 Male 0
13040 2020-03-13 09:59:55 Clothing 279.0 2 2187.0 PayPal 55 2 Michelle Flores 55 Female 1
48835 2021-11-23 01:30:42 Home 27.0 1 3615.0 Credit Card 42 1 Jeremy Rush 42 Male 1
21019 2020-07-02 14:04:48 Home 17.0 5 2466.0 Cash 41 0 Tina Craig 41 Male 0
49234 2020-12-30 02:02:40 Books 398.0 2 3668.0 Crypto 34 0 Jennifer Cooper 34 Female 1
10971 2021-03-13 16:28:39 Electronics 425.0 4 2370.0 Cash 36 1 Justin Lawson 36 Female 1
Time taken: 0.071 seconds, Fetched: 250001 row(s)
hive>

```

Select hadoop@ip-172-31-77-105~

```
customer_age > 45 THEN 'Above 45' END AS age_group, AVG(total_purchae_amount) AS avg_purchase_amount FROM customer_purchases GROUP BY CASE WHEN customer_age < 18 THEN 'Under 18' WHEN customer_age BETWEEN 18 AND 25 THEN '18-25' WHEN customer_age BETWEEN 26 AND 35 THEN '26-35' WHEN customer_age BETWEEN 36 AND 45 THEN '36-45' WHEN customer_age > 45 THEN 'Above 45' END ORDER BY age_group;
FAILED: ParseException line 13:0 missing EOF at 'SELECT' near 'age_group'
hive> SELECT CASE WHEN customer_age < 18 THEN 'Under 18' WHEN customer_age BETWEEN 18 AND 25 THEN '18-25' WHEN customer_age BETWEEN 26 AND 35 THEN '26-35' WHEN customer_age BETWEEN 36 AND 45 THEN '36-45' WHEN customer_age > 45 THEN 'Above 45' END AS age_group, AVG(total_purchae_amount) AS avg_purchase_amount FROM customer_purchases GROUP BY CASE WHEN customer_age < 18 THEN 'Under 18' WHEN customer_age BETWEEN 18 AND 25 THEN '18-25' WHEN customer_age BETWEEN 26 AND 35 THEN '26-35' WHEN customer_age BETWEEN 36 AND 45 THEN '36-45' WHEN customer_age > 45 THEN 'Above 45' END ORDER BY age_group;
Query ID = hadoop_20241107054052_e69d72c9-83f6-4936-9479-95aad53e2d2
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1730955833022_0002)
```

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	2	2	0	0	0	0
Reducer 2	container	SUCCEEDED	2	2	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0

VERTICES: 03/03 [=====] 100% ELAPSED TIME: 7.46 s

OK

18-25 2618.180526219719

26-35 2655.71349081141

36-45 2710.81235598186

Above 45 2793.973996036033

NULL NULL

Time taken: 7.884 seconds, Fetched: 5 row(s)

hive>

hadoop@ip-172-31-77-105~

```
VERTICES: 03/03 [=====] 100% ELAPSED TIME: 7.61 s
OK
Time taken: 14.31 seconds
hive>
> SELECT product_category, SUM(returns) AS total_returns
> FROM customer_purchases
> GROUP BY product_category
> ORDER BY total_returns DESC;
Query ID = hadoop_20241107060056_iddba25a-2f46-44f6-a4fb-0616c23a502f
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1730955833022_0003)
```

VERTICES	MODE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	2	2	0	0	0	0
Reducer 2	container	SUCCEEDED	2	2	0	0	0	0
Reducer 3	container	SUCCEEDED	1	1	0	0	0	0

VERTICES: 03/03 [=====] 100% ELAPSED TIME: 6.48 s

OK

Books 30229

Clothing 30122

Electronics 20271

Home 20147

Product Category NULL

Time taken: 6.734 seconds, Fetched: 5 row(s)

hive>

27°C
Light rain

Search

ENG
IN 23:51
06-11-2024

27°C
Light rain

Search

ENG
IN 00:01
07-11-2024

```

hadoop@ip-172-31-77-105~
ty, total_purchase_amount, payment_method, customer_age, returns, customer_name, age, gender, churn)
hive> select product_category, sum(total_purchase_amount) as total_sales from customer_purchases group by product_category order by total_sales desc;
Query ID = hadoop_20241107053631_7defa92a-4aa0-4258-a064-9d1fc0f7a35d
Total jobs = 1
Launching Job 1 out of 1
tez session was closed. Reopening...
Session re-established.
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1730955833022_0002)

-----
VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    2         2         0         0         0         0
Reducer 3 ..... container  SUCCEEDED    1         1         0         0         0         0
-----
VERTICES: 03/03 [=====] 100% ELAPSED TIME: 7.19 s
-----
OK
Books      2.04939601E8
Clothing   2.04532405E8
Electronics 1.36599467E8
Home       1.3527121E8
Product Category      NULL
Time taken: 12.966 seconds, Fetched: 5 row(s)
hive>

```

8. Now Data analysis is completed in HIVE. Then we will Preprocess the data and analyse, train and the test the data in Pyspark In aws using machine learning models.
9. So we will use the same bucket with our data set and run the cluster with required policies including Jupiter Notebook to run the Pyspark Code.
10. So we run the cluster and then we will go to studios and create a new workspace to run our pyspark queries in jupiter notebook.
11. Then we will quick launch our workspace and Jupiter notebook will be opened and now we can create a notebook with ecommerce name and we will run our queries by selecting Pyspark as kernel.

The screenshot displays the EMR Studio interface, which is divided into several sections:

- Left Sidebar (Navigation):**
 - EMR Studio** (with a close button)
 - Dashboard**
 - Workspaces** (highlighted)
 - Query editor** (with a 'New' button and 'Powered by Athena')
 - Serverless**
 - Applications
 - Clusters**
 - EMR on EC2
 - EMR on EKS
 - Service integrations**
 - SageMaker Data Wrangler (with a 'New' button and 'Low-code data prep and ML')
 - What's New**
 - Submit feedback**
 - Logout**
- Main Content Area (Workspaces):**
 - Header: **Workspaces** (with 'EMR Studio > Workspaces' breadcrumb)
 - Sub-header: **Studio: ecommerce-project**
 - Buttons: **Workspaces (2) Info**, **Actions** (dropdown), **Launch Workspace** (dropdown), and **Create Workspace** (orange button).
 - Search bar: **Find Workspaces by name, status, or last modified by**
 - Table of Workspaces:
- Bottom Section (Compute):**
 - Header: **Compute**
 - Text: "To run your notebooks, attach your workspace to a compute layer. All notebooks in the Workspace share the same compute. Make sure that your EMR Studio user role has appropriate permissions for your selected compute type. [Learn more](#)"
 - Compute type**
 - ☐ EMR Serverless application
 - ☒ EMR on EC2 cluster
 - ☐ EMR on EKS cluster
 - EMR on EC2 cluster** (with a dropdown menu showing 'ecommerce-cluster...') and a refresh button.
 - Runtime role** (with a dropdown menu and a note: "If you have configured your cluster for runtime roles, you must choose an IAM role for Amazon EMR to assume when it runs your")
- Bottom Panel (Jupyter Notebook):**
 - Tab: **ecommerce-project.ipynb**
 - Buttons: **Code** (dropdown), **Cluster attached...**, and **PySpark** (with a settings icon).
 - Code execution output:

	Workspace name	Status	Cluster ID	Creation time (UTC-05:00)	Last modified
<input type="radio"/>	ecommerce-project	Ready	j-1W92J8UKHFOUH	November 07, 2024, 18:00	ecommerce-pr
<input type="radio"/>	Studio_1_Workspace_1	Ready	-	November 07, 2024, 17:58	ecommerce-pr

```
Last executed at 2024-11-07 19:07:13 in 1m 0.06s

> Spark Job Progress

Gradient-Boosted Tree Churn Prediction Accuracy: 0.7993213669035859

•[11]: from pyspark.ml.classification import LinearSVC

# Initialize SVM model
svm = LinearSVC(labelCol="Churn", featuresCol="features", maxIter=10, regParam=0.1)
svm_model = svm.fit(train_df)
svm_predictions = svm_model.transform(test_df)

# Evaluate SVM model
svm_accuracy = evaluator.evaluate(svm_predictions)
print("SVM Churn Prediction Accuracy:", svm_accuracy)

Last executed at 2024-11-07 19:07:51 in 32.49s

> Spark Job Progress

SVM Churn Prediction Accuracy: 0.7993213669035859

•[12]: from pyspark.ml.classification import NaiveBayes

# Initialize Naive Bayes model
nb = NaiveBayes(labelCol="Churn", featuresCol="features", smoothing=1.0)
nb_model = nb.fit(train_df)
nb_predictions = nb_model.transform(test_df)
```

At the bottom of the interface, there is a status bar showing: **Simple** (toggle), **0**, **1**, **PySpark | Unknown**, **CodeWhisperer**, **Saving completed**, **Mode: Command**, **Ln 1, Col 1**, and **ecommerce-project.ipynb**.

The screenshot displays the Amazon EMR Studio interface, showing two notebooks in the 'ecommerce-project' workspace.

Top Notebook: ecommerce-project.ipynb

Code Cell [12]:

```
from pyspark.ml.classification import NaiveBayes

# Initialize Naive Bayes model
nb = NaiveBayes(labelCol="Churn", featuresCol="features", smoothing=1.0)
nb_model = nb.fit(train_df)
nb_predictions = nb_model.transform(test_df)

# Evaluate Naive Bayes model
nb_accuracy = evaluator.evaluate(nb_predictions)
print("Naive Bayes Churn Prediction Accuracy:", nb_accuracy)
```

Last executed at 2024-11-07 19:08:04 in 5.95s

Code Cell [13]:

```
from pyspark.ml.classification import MultilayerPerceptronClassifier

# Define network layers for MLP
layers = [len(feature_columns), 20, 10, 2] # Adjust based on your data

# Initialize MLP model
mlp = MultilayerPerceptronClassifier(labelCol="Churn", featuresCol="features", maxIter=100, layers=layers)
mlp_model = mlp.fit(train_df)
mlp_predictions = mlp_model.transform(test_df)

# Evaluate MLP model
mlp_accuracy = evaluator.evaluate(mlp_predictions)
print("MLP Churn Prediction Accuracy:", mlp_accuracy)
```

Bottom Notebook: ecommerce-project.ipynb

Code Cell [5]:

```
# 2. Churn Prediction (Predictive Modeling)
# Index categorical variables
indexer = StringIndexer(inputCols=["Product Category", "Payment Method", "Gender"], outputCols=["ProductCategoryIndex", "PaymentMethodIndex", "GenderIndex"])
df_indexed = indexer.fit(df).transform(df)

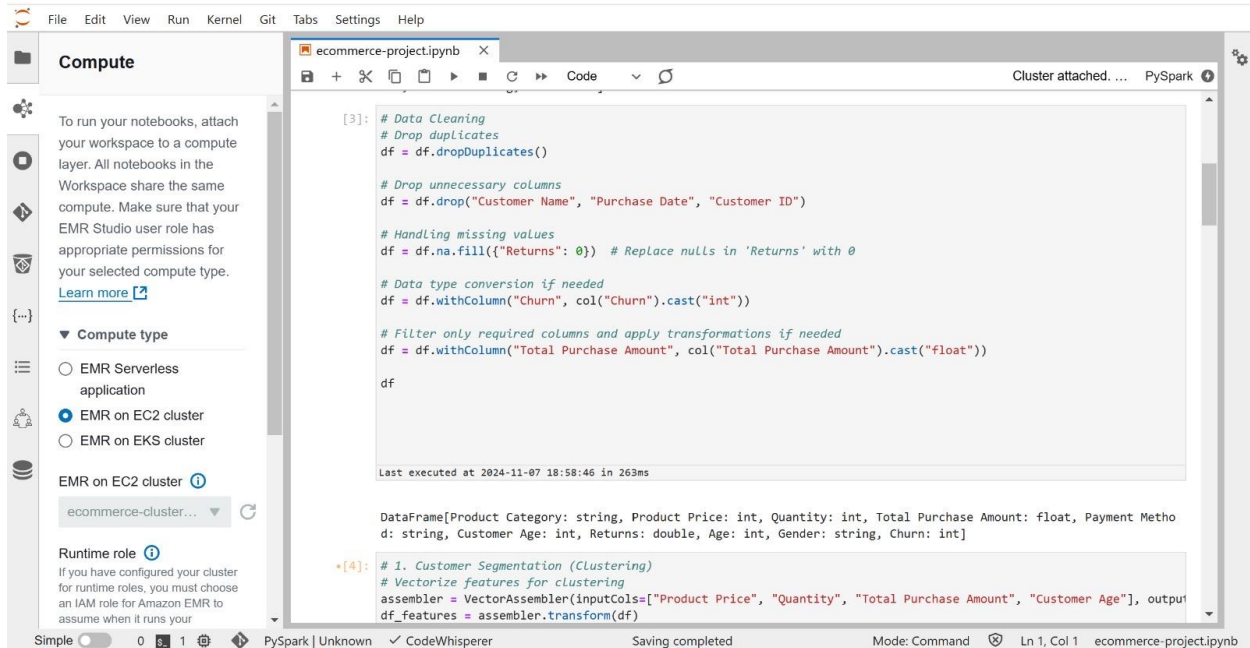
# Vectorize features for modeling
feature_columns = ["ProductCategoryIndex", "PaymentMethodIndex", "GenderIndex", "Product Price", "Quantity", "Total"]
assembler = VectorAssembler(inputCols=feature_columns, outputCol="features")
df_model = assembler.transform(df_indexed)

# Train-Test Split
train_df, test_df = df_model.randomSplit([0.8, 0.2], seed=42)

# Churn Prediction Model - Random Forest
rf = RandomForestClassifier(labelCol="Churn", featuresCol="features")
rf_model = rf.fit(train_df)
rf_predictions = rf_model.transform(test_df)

# Evaluation
evaluator = MulticlassClassificationEvaluator(labelCol="Churn", metricName="accuracy")
rf_accuracy = evaluator.evaluate(rf_predictions)
print("Random Forest Churn Prediction Accuracy:", rf_accuracy)

# Logistic Regression for comparison
lr = LogisticRegression(labelCol="Churn", featuresCol="features")
lr_model = lr.fit(train_df)
```

12. Then we will Visualize the data with our machine learning Models with accuracy and f1 score using pandas, matplotlib and seaborn libraries.

