E-commerce Customer Behavior Analysis Report

1. Introduction

The dataset used for this analysis can be accessed from the following Kaggle URL: E-commerce Customer Behavior Dataset. It contains a wealth of information that allows us to explore various aspects of customer behavior in the e-commerce domain.

To facilitate the analysis, we set up a Hadoop Docker cluster. The Hadoop cluster played a crucial role in storing and preprocessing the large datasets of user logs and transaction data. Additionally, we utilized Apache Spark, a powerful distributed computing framework, to perform machine learning algorithms for analyzing customer behavior and predicting future buying patterns.

2. Methodology

The analysis process involved several steps, as outlined below:

```
nswar@norhanswar-HP-Laptop-15s-fq5xxx:~/Career/Master/BigData/Labs/Lab2Material-20231203/docker-spark-cluster$ sudo docker ps
  [sudo] password for norhanswar:
CONTAINER ID IMAGE
                                                                                                                                                                                                                                    CREATED
                                                                                                                                                                        NAMES
"/bin/bash /start-sp..." 6
080/tcp db50e401bc52_spar
"/entrypoint.sh /run..." 6
datanode
  a5a702552e2b cluster-apache-spark:3.0.2
77/tcp, :::7077->7077/tcp, 7000/tcp, 0.0.0.0:8080->8080/tcp, :::8080->808
775576b3838 bde2020/hadoop-datanode:2.0.0-hadoop3.2.1-java8 "
54/tcp, :::9864->9864/tcp
odbb635916db bde2020/hadoop-nodemanager:2.0.0-hadoop3.2.1-java8 "
                                                                                                                                                                                                                                 6 hours ago Up 5 hours
                                                                                                                                                                                                                                                                                                                           0.0.0.0:7077->70
                                                                                                                                                                                                                                   6 hours ago Up 5 hours (healthy) 0.0.0.0:9864->98
                                                                                                                                                                         "/entrypoint.sh /run..."
a2d82f6b9ee5 bde2020/hadoop-historyserver:2.0.0-hadoop3.2.1-java8
                                                                                                                                                                                                                                 6 hours ago Up 5 hours (healthy)
                                                                                                                                                                                                                                                                                                                           8188/tcp
11d378fb5a32 bde2020/hadoop-resourcemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run...
resourcemanager:2.0.0-hadoop3.2.1-java8 "/entrypoint.sh /run..."
                                                                                                                                                                                                                                6 hours ago Up 5 hours (healthy) 8088/tcp
ad858515a992 bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8 "resourcemanager ad858515a992 bde2020/hadoop-namenode:2.0.0-hadoop3.2.1-java8 "/entrypoint.s.h /run..." 6 hours ago Up 5 hours 70/tcp, :::9870->9870/tcp, 0.0.0.0:9010->9000/tcp, :::9010->9000/tcp namenode "/bin/bash /start-sp..." 4 days ago Up 5 hours 00/tcp, :::7000->7000/tcp, 77/tcp, 0.0.0.0:9091->8080/tcp, :::9091->8080/tcp docker-spark-cluster-spark-worker-a-1 a93ac807d60a cluster-apache-sparksi3.0.2 "/bin/bash /start-sp..." 4 days ago Up 5 hours 0:7001->7000/tcp, :::7001->7000/tcp, 0.0.0.0:9092->8080/tcp, :::9092->8080/tcp docker-spark-cluster-spark-worker-b-1 ef41aa42d125 postgres:11.7-alpine "docker-entrypoint.s..." 4 days ago Up 5 hours 32/tcp, :::5432->5432/tcp docker-spark-ruster-spark-worker-b-1 morbanswarengen answaren - HP-laptop-15s-fg5xxx:-/career/Master/8lg0ata/Labs/Lab2Material-20231233/docker-spark-cluster-$
                                                                                                                                                                                                                                                                                                                           0.0.0.0:7000->70
                                                                                                                                                                                                                                                                                                                           7077/tcp. 0.0.0.
                                                                                                                                                                                                                                                                                                                           0.0.0.0:5432->54
```

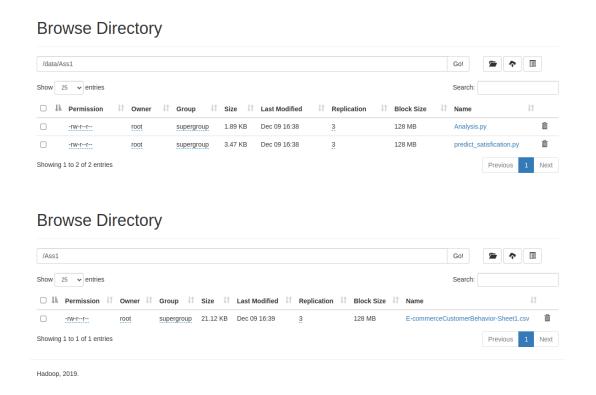
.

2.1. Setting up the Hadoop Docker Cluster

To begin with, we ran the Docker Compose file to set up the Hadoop Docker cluster.

2.2. Storing the Dataset

Once the Hadoop cluster was up and running, we uploaded and stored the dataset in the Hadoop Distributed File System (HDFS). This step ensured that the dataset was readily available for analysis and processing.



2.3. Analysis and Preprocessing with PySpark

Using PySpark, we performed analysis and preprocessing on the dataset. The analysis focused on answering the following key questions:

Segment customers based on their demographic information (Age, Gender, City) and shopping behaviors (Total Spend, Number of Items Purchased, Membership Type)

Silver Between 35-50 Los Angeles 288 814.65 Bronze Between 35-50 Chicago 546 499.8827586206897 Gold Under 35 New York 901 1165.0355932203393 Silver Under 35 Los Angeles 401 799.2114285714289 Bronze Between 35-50 Houston 425 447.6482142857141 Silver Under 35 Miami 675 690.3896551724141	Membership Type	Age Group	City	sum(Items Purchased)	avg(Total Spend)
Gold Under 35 San Francisco 1141 1460.4561403508774 Gold Between 35-50 San Francisco 19 1420.8	Bronze Gold Silver Bronze Silver Gold	Between 35-50 Under 35 Under 35 Between 35-50 Under 35 Under 35	Chicago New York Los Angeles Houston Miami San Francisco	546 901 401 425 675 1141	499.8827586206897 1165.0355932203393 799.2114285714289 447.6482142857141 690.3896551724141 1460.4561403508774

```
| City|sum(Items Purchased)| avg(Total Spend)|
| Los Angeles| 689| 805.4915254237288|
|San Francisco| 1160|1459.7724137931039|
| Chicago| 546| 499.8827586206897|
| Houston| 425| 447.6482142857141|
| Miami| 675| 690.3896551724141|
| New York| 901|1165.0355932203393|
```

Which customers are at risk of not making future purchases based on their Days Since Last Purchase and Satisfaction Level

			eGenerator: Code							
-+	ID Gender #	Age	City Membership	Type Total	Spend Items	Purchased Average	Rating Discount	Applied Days Sin	ce Last Purchase Satisf	action Leve
-+ 	103 Female				510.75	9	3.4	true	42	Unsatisfie
:	105 Male	27	Miami S	ilver	720.4	13	4.0	true	55	Unsatisfie
a 	109 Female	41 Ch	icago B	ronze ·	495.25	10	3.6	true	40	Unsatisfie
a 	111 Male	32	Miami S	ilver	690.3	11	3.8	true	34	Unsatisfie
d 	115 Female	42 Ch	icago B	ronze	530.4	9	3.5	true	38	Unsatisfie
d 	117 Male	26	Miami S	ilver	700.6	12	3.7	true	48	Unsatisfie
d 	121 Female	43 Ch	icago B	ronze	505.75	10	3.3	true	41	Unsatisfie
d 	123 Male	27	Miami S	ilver	710.4	13	4.1	true	54	Unsatisfie
d	127 Female	41 Ch	icago B	ronze	485.25	9	3.6	true	39	Unsatisfie
d 	1291 Malel	321	Miamil S	ilverl	670.31	101	3.81	truel	331	Unsatisfie



Spark Master at spark://a5a702552e2b:7077

URL: spark://a5a702552e2b:7077 Alive Workers: 2 Cores in use: 2 Total, 0 Used Memory in use: 2.0 GiB Total, 0.0 B Used Resources in use: Applications: 0 Running, 1 Completed Drivers: 0 Running, 0 Completed Status: ALIVE

→ Workers (2)

Worker Id	Address	State	Cores	Memory	Resources
worker-20231209143209-172.19.0.10-7000	172.19.0.10:7000	ALIVE	1 (0 Used)	1024.0 MiB (0.0 B Used)	
worker-20231209143209-172.19.0.9-7000	172.19.0.9:7000	ALIVE	1 (0 Used)	1024.0 MiB (0.0 B Used)	

→ Running Applications (0)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration

▼ Completed Applications (1)

Application ID	Name	Cores	Memory per Executor	Resources Per Executor	Submitted Time	User	State	Duration
app-20231209144049-0000	PySparkShell	2	1024.0 MiB		2023/12/09 14:40:49	root	FINISHED	10 s