

CSP 554 Big Data Technologies

Project Report

Online Retail Simulation and Analytics

Team Members:

Ravi Teja Batchala	(A20512513)
Rajesh Mavi	(A20481442)
Sadanand Satish Kolhe	(A20517969)
Rik Ganguli Biswas	(A20513718)

**Under The Supervision
of
Prof. Joseph Rosen**



**College Of Computing
Illinois Institute of Technology
Chicago, Illinois. December 2023**

Contents

1. Project Details

1.1 Project Topic

1.2 Overview

1.3 Introduction

1.4 Literature Review

1.5 Technologies

2. Data Operations and Analysis

2.1 Architecture and Data Flow

2.2 Data Source and Feature Analysis

2.3 Data Cleaning

2.4 Data Processing

2.5 Data Visualization

3. Difficulties Faced

4. Conclusion

5. Future Scope

6. References

I: PROJECT DETAILS

1.1 Project Topic

Apply a range of big data tools to explore some interesting data sets and derive insights from them. Ingest data, apply transformations, profile the data, summarize it, visualize it.

1.2 Overview

This project aims at crafting a lifelike simulation for an ecommerce platform's online shopping experience. The goal is to mimic user interactions, efficiently store data in both structured and unstructured formats, and analyze behaviors for valuable insights. Real-time processing, user segmentation, product analysis, and optional machine learning predictions are on the agenda.

1.3 Introduction

In the dynamic landscape of e-commerce, understanding user behavior and optimizing the online shopping experience is paramount for business success. This project endeavors to create a sophisticated and lifelike simulation of an e-commerce platform, aiming to replicate user interactions, streamline data storage, and extract valuable insights through robust real-time processing and analysis.

The primary objectives of this project encompass a comprehensive approach to user simulation, data management, and behavioral analysis. By generating synthetic data for various user actions, employing efficient storage mechanisms such as MySQL, CSV, and Parquet, and implementing a robust data pipeline using Apache Kafka, PySpark Streaming, and Apache Hive, the project aspires to deliver a holistic solution for e-commerce analytics. The project aims to improve decision-making by creating brief hourly reports and utilizing data visualization tools like Matplotlib and Plotly to display insightful charts and dashboards. This adaptability to changing user preferences and market trends is crucial to staying competitive in the dynamic field of online retail.

1.4 Literature Review

Citation	Author	Year	Key Findings
A study of preferences in a simulated online shopping experiment	Asle Fagerstrøm, Erik Arntzen, Gordon Foxall	2011	Study reveals consumer brand loyalty development influenced by environmental contingencies, showing preferences and reduced switching over time.
Effects of Pros and Cons of Applying Big Data Analytics to Consumers' Responses in an E-Commerce Context	Le, T. M, Liaw, S.-Y	2017	Positive effects of Big Data analytics on customer responses, emphasizing recommendation systems and dynamic pricing.
Analysis of the Impact of Big Data on E-Commerce in Cloud Computing Environment	Rongrui Yu, Chunqiong Wu, Bingwen Yan, Baoqin Yu, Xiukao Zhou, Yanliang Yu, Na Chen	2021	The study integrates MySQL, HBase, and cloud computing for e-commerce efficiency, leveraging big data for precision marketing and enhanced competitiveness.
Statistical Modeling and Simulation of Online Shopping Customer Loyalty Based on Machine Learning and Big Data Analysis	Jui-Chan Huang, Po-Chang Ko, Cher-Min Fong, Sn-Man Lai, Hsin-Hung Chen, Ching-Tang Hsieh	2021	Deleting null values and estimating replacements in commodity datasets improve data quality for effective credit risk assessment.
Mechanism of Big Data Analytics in Consumer Behavior on Online Shopping	Ruby Evangelin, Vasantha Shanmugam	2022	Big data variables (Capacity, Promptness, Data Range) impact online browsing behavior and influence purchasing decisions.

1.5 Technologies

- 1. Apache Spark:** Apache Spark is an open-source, distributed computing system that provides a fast and general-purpose cluster-computing framework for big data processing, enabling scalable and efficient data analytics and machine learning applications.
- 2. Apache PySpark Streaming:** Spark Streaming is a micro-batch processing framework for real-time data. It enables the processing of data streams in near real-time, making it suitable for applications requiring low-latency analytics.
- 3. Apache Kafka:** Apache Kafka is a distributed event streaming platform designed for building real-time data pipelines and streaming applications. It is capable of handling high-throughput, fault-tolerant, and scalable data streaming. Kafka provides a publish-subscribe model for exchanging records between producers and consumers.
- 4. Apache Hive:** Apache Hive is a data warehousing and SQL-like query language system built on top of Hadoop. It provides a high-level interface for analyzing and querying large datasets stored in Hadoop Distributed File System (HDFS) or other compatible file systems.
- 5. MySQL:** MySQL is an open-source relational database management system widely used for storing and retrieving structured data. Renowned for its performance, scalability, and ease of use,
- 6. Matplotlib/Plotly:** Matplotlib, a versatile Python plotting library, empowers our data analysis by offering a comprehensive suite of visualization tools. Its diverse range of customizable plots, from line graphs to histograms, enables us to effectively communicate and explore complex datasets.
- 7. Hadoop:** Hadoop is an open-source framework for distributed storage and processing of large datasets across clusters of computers, using a programming model called MapReduce
- 8. AWS S3:** Amazon S3, part of Amazon Web Services (AWS), is a scalable object storage service facilitating the storage and retrieval of data from any web location, providing seamless scalability.
- 9. Zookeeper:** Apache ZooKeeper is an open-source coordination service for distributed applications. It provides a centralized infrastructure for maintaining configuration information, naming, providing distributed synchronization, and group services.
- 10. Jupyter Notebook**

II: DATA OPERATIONS AND ANALYSIS

2.1 Architecture and Data Flow

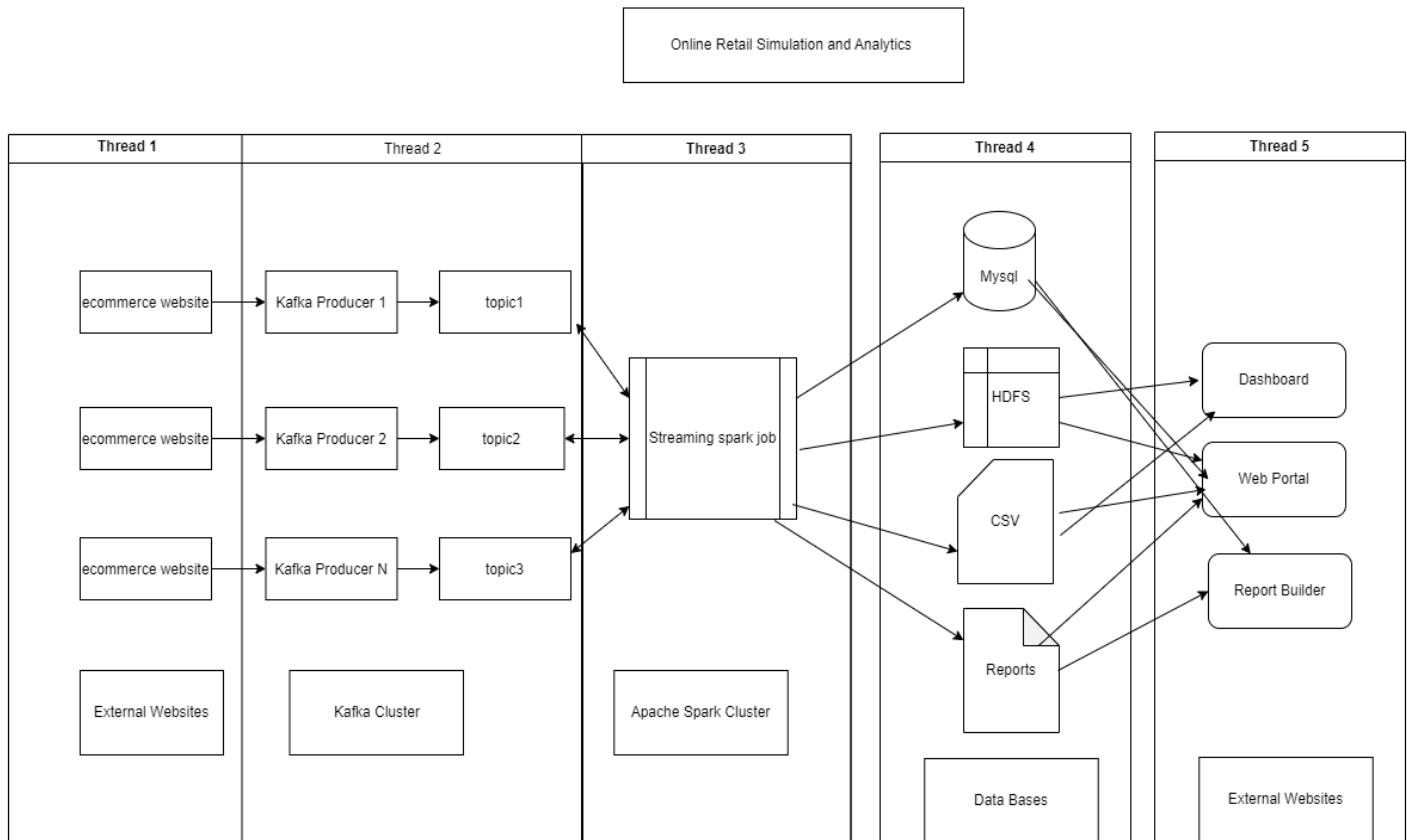


Fig 1: Architectural Diagram of Online Retail Simulation and Analytics System

1. Overview:

- System designed for real-time e-commerce system with Kafka producers, brokers, and PySpark streaming job. Workflow: ingest data, publish to Kafka, process with PySpark, generate reports, and store in a database.

2. Data Flow:

- E-commerce data sent to Kafka producers and Kafka producers publish data to Kafka topics.
- PySpark streaming job consumes data from Kafka topics in real-time.

3. Processing:

- PySpark streaming job executes intricate transformations on the incoming data.
- Processed data forms the basis for generating detailed reports.

4. Storage:

- Reports are systematically stored in a structured database.
- Database acts as a repository for future analytical exploration.

2.2 Data Source and Feature analysis

UserProfilesGenerator

Utilized UserProfilesGenerator to synthesize user profiles by considering factors such as demographic information, preferences, and historical behavior. The generated **UserProfiles_data.csv** includes attributes such as User_ID, Name, Age, Gender, Income, Contact, and Location, offering a comprehensive collection of user characteristics within the dataset.

ProductDataGenerator

Executed ProductDataGenerator to create synthetic product data by considering features such as Product_ID, Product_Name, Category, and Price. The resulting **Product_Data.csv** dataset encapsulates a varied assortment of products, mimicking a diverse product ecosystem.

UserInteractionsGenerator

Simulated user interactions with products using the tool, incorporating parameters such as user preferences, behavior patterns, and interaction types. The resulting dataset, **UserInteractions_data.csv**, includes attributes such as Interaction_ID, Timestamp, User_ID, Interaction_Type, Product_ID, Product_Name, Quantity, and Payment_Method, providing a diverse set of user engagement scenarios.

Feature Analysis

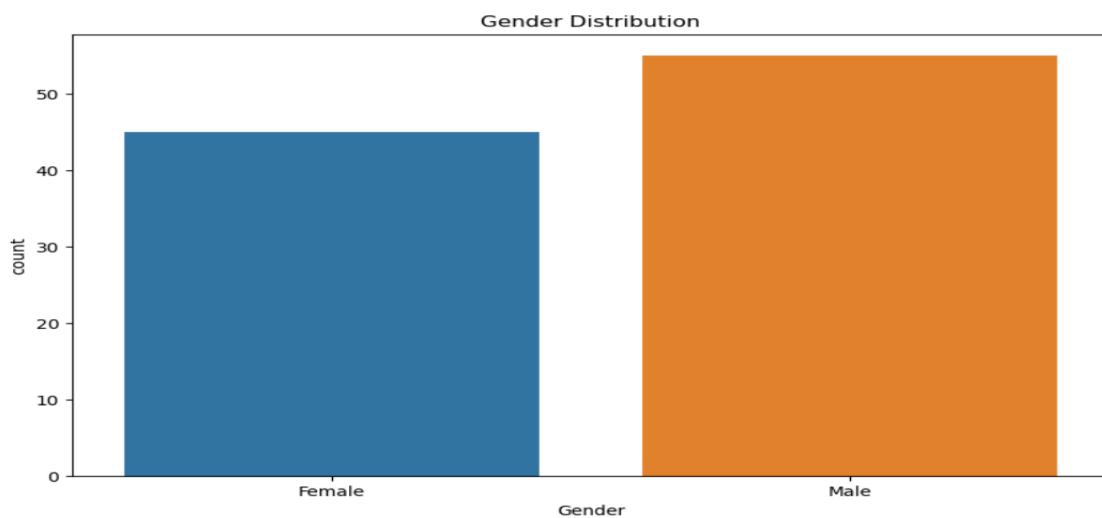
1) UserProfiles_data

Attribute	Description	Data Type
User_ID	User identifier	Numeric
Name	User's name	String
Age	User's age	Numeric
Gender	User's gender	String
Income	User's income	Numeric
Contact	User's contact info	String
Location	User's location	String

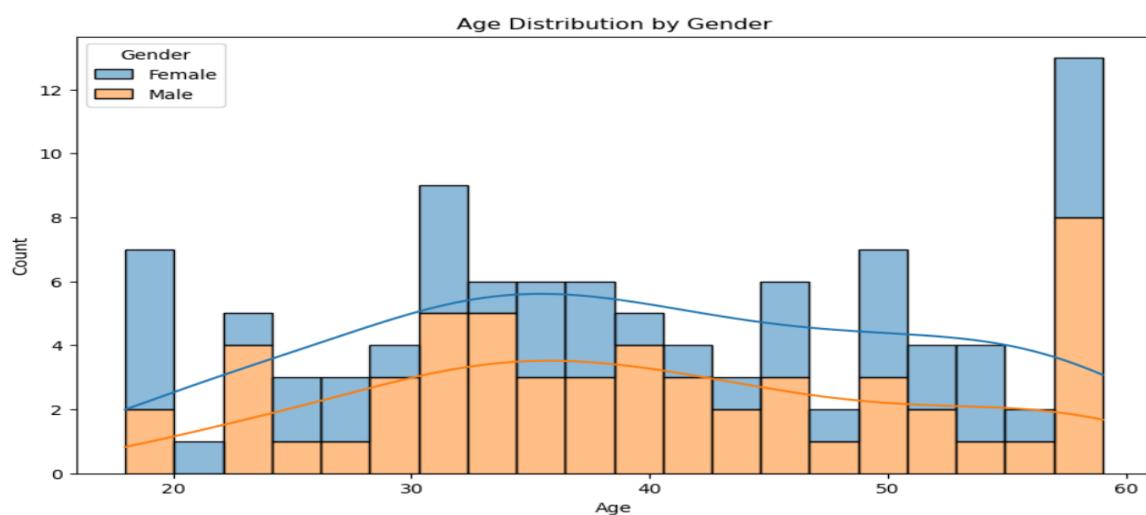
User_ID	Name	Age	Gender	Income	Contact	Location
1	Heather Burke	35	Female	68059	347-414-3911	61.389715, -7.198571
2	Margaret Brown	40	Male	73613	+1-632-865-1129	-61.234150, -9.520890
3	Pamela Crane	53	Male	75572	816.752.4431x56859	-68.287243, -135.834580
4	Erika Williamson	36	Male	63795	+1-652-669-8271	MDR37A, Basni, Nagaur Tehsil, Nagaur District, Rajasthan, 341021, India
5	Lori McIntosh	38	Male	31765	777-357-5255x67102	80.204909, -172.165824

Fig 2: Sample Data of UserProfile

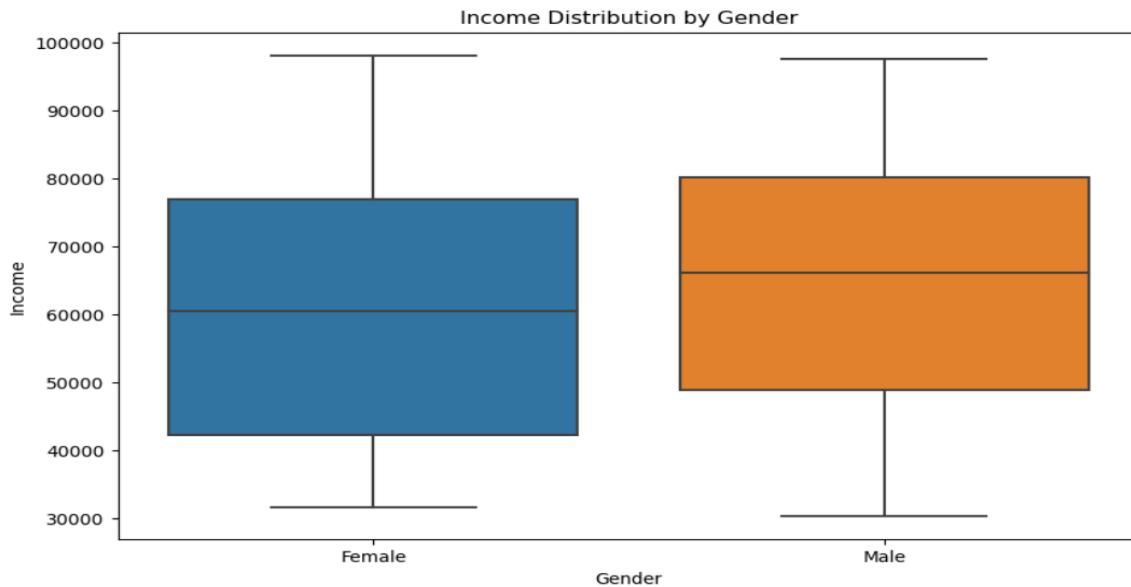
Gender Distribution: This provides a clear overview of the gender composition within a given dataset



Age distribution by gender: the x-axis represents different age groups, and the y-axis shows the count or percentage of individuals in each age group, segregated by male and female.



Income Distribution by gender: This visualization provides insights into income disparities or patterns across different gender demographics within a given population or dataset.



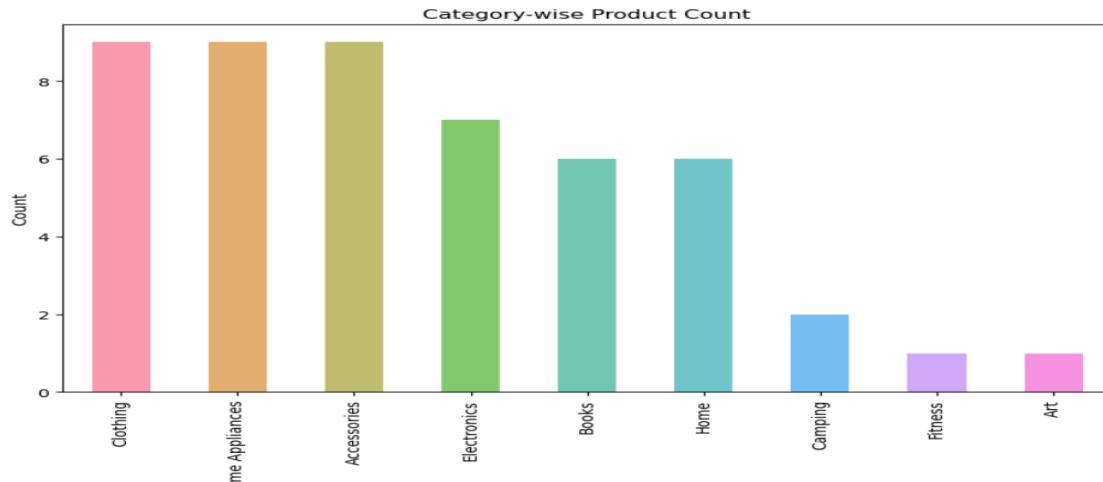
2) Products_data

Attribute	Description	Data Type
Product_ID	Product identifier	Numeric
Product_Name	Name of the product	String
Category	Product category	String
Price	Price of the product	Numeric

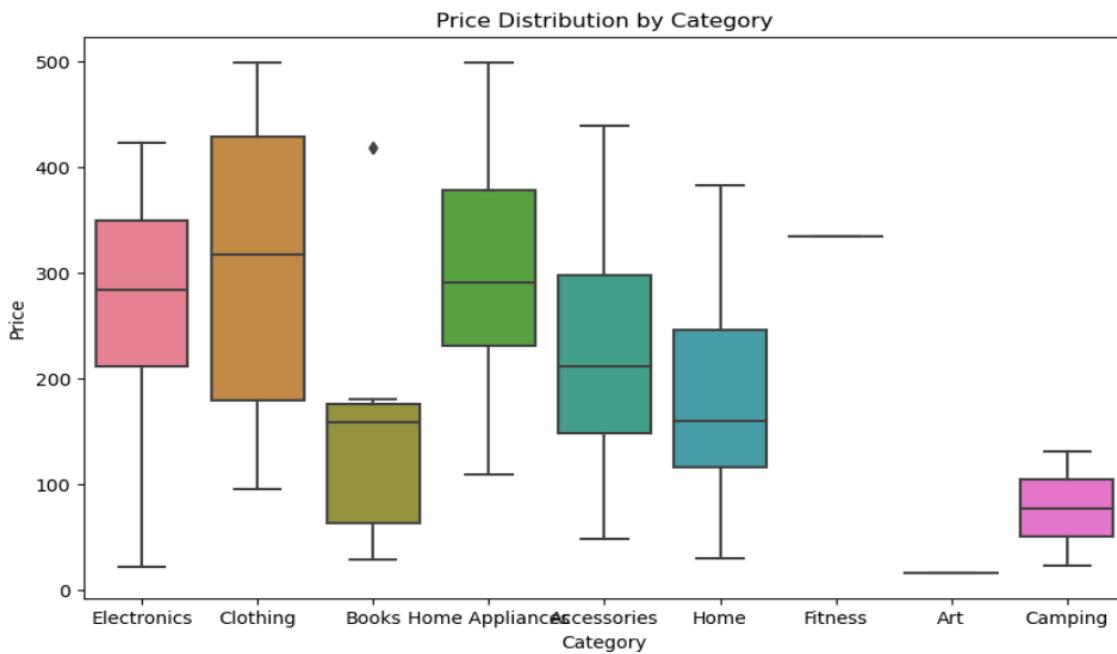
Product_ID	Product_Name	Category	Price
1	Laptop	Electronics	423
2	Smartphone	Electronics	174
3	Headphones	Electronics	406
4	Camera	Electronics	248
5	Television	Electronics	22

Fig 3: Sample Data of UserProfile

Category wise product count:-provides a clear overview of the distribution of products across various categories



Price Distribution by category:--The graph displays the distribution of prices for different product categories using a box plot. Provides a comparative analysis of price variations across different product categories



3) UserInteractions_data

Attribute	Description	Data Type
Interaction_ID	Interaction identifier	Numeric
Timestamp	Time of the interaction	Datetime
User_ID	User identifier	Numeric
Interaction_Type	Type of interaction	String
Product_ID	Product identifier	Numeric
Product_Name	Name of the product	String
Quantity	Quantity of the product	Numeric
Payment_Method	Payment method used	String

Interaction_ID	Timestamp	User_ID	Interaction_Type	Product_ID	Product_Name	Quantity	Payment_Method
1	2023-11-01 00:00:00	1	product_search	-	laptop	0.0	Not Applicable
2	2023-11-01 00:00:00	2	product_view	49	Camping Tent	0.0	Not Applicable
3	2023-11-01 00:00:00	3	product_search	-	shoes	0.0	Not Applicable
4	2023-11-01 00:00:00	4	product_search	-	electronics	0.0	Not Applicable
5	2023-11-01 00:00:00	5	product_search	-	phone	0.0	Not Applicable

Fig 4: Sample Data of UserInteractions_data

2.3 Data Cleaning

Removed outliers:- Removed outliers to enhance data quality and improve model performance by eliminating extreme values that could distort statistical analyses, ensuring a more accurate representation of the underlying data distribution.

```
def find_outliers(self, data):
    # Explain outliers for numerical attribute 'Price'
    outliers = data.filter(F.col('Price') > 1000)
    print("Outliers in Price:")
    outliers.show()

def find_outliers(self, data):
    # Explain outliers for numerical attribute 'Age'
    outliers = data.filter(F.col('Age') > 100)
    print("Outliers in Age:")
    outliers.show()
```

Null/empty/hyphen values:- Handled null values by replacing them with the mean, mitigating bias in analyses, maintaining data integrity, and ensuring consistency in data processing. This approach enhances the robustness of statistical computations and facilitates a more accurate representation of the dataset.

```
def preprocess_data(self, data):
    # Handle null values for numerical attribute 'Price'
    data = data.na.fill({'Price': data.select(F.mean('Price')).collect()[0][0]})
    return data

def preprocess_data(self, data):
    # Handle null values for numerical attribute 'Age'
    data = data.na.fill({'Age': data.select(F.mean('Age')).collect()[0][0]})
    return data

def preprocess_data(self, data: DataFrame) -> DataFrame:
    # Additional handling for null, empty, hyphen, or negative Product_ID
    return data.withColumn('Product_ID', F.when((F.col('Product_ID').isNull()) | (
        F.col('Product_ID') == '') | (F.col('Product_ID') == '-')) |
        (F.col('Product_ID') < 0), 0).otherwise(F.col('Product_ID')))

    df['Product_ID'].replace('-', np.nan, inplace=True)
```

2.4 Data Processing

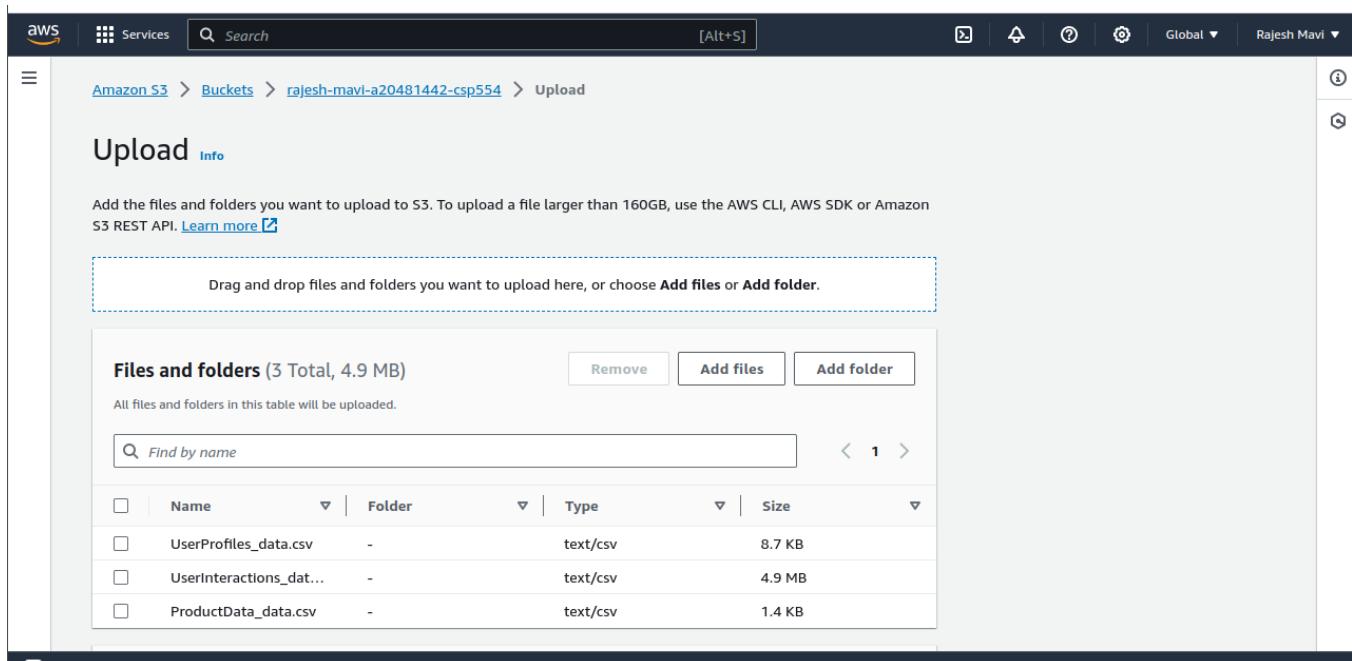


Fig 5: Creating a Bucket and added the data files.

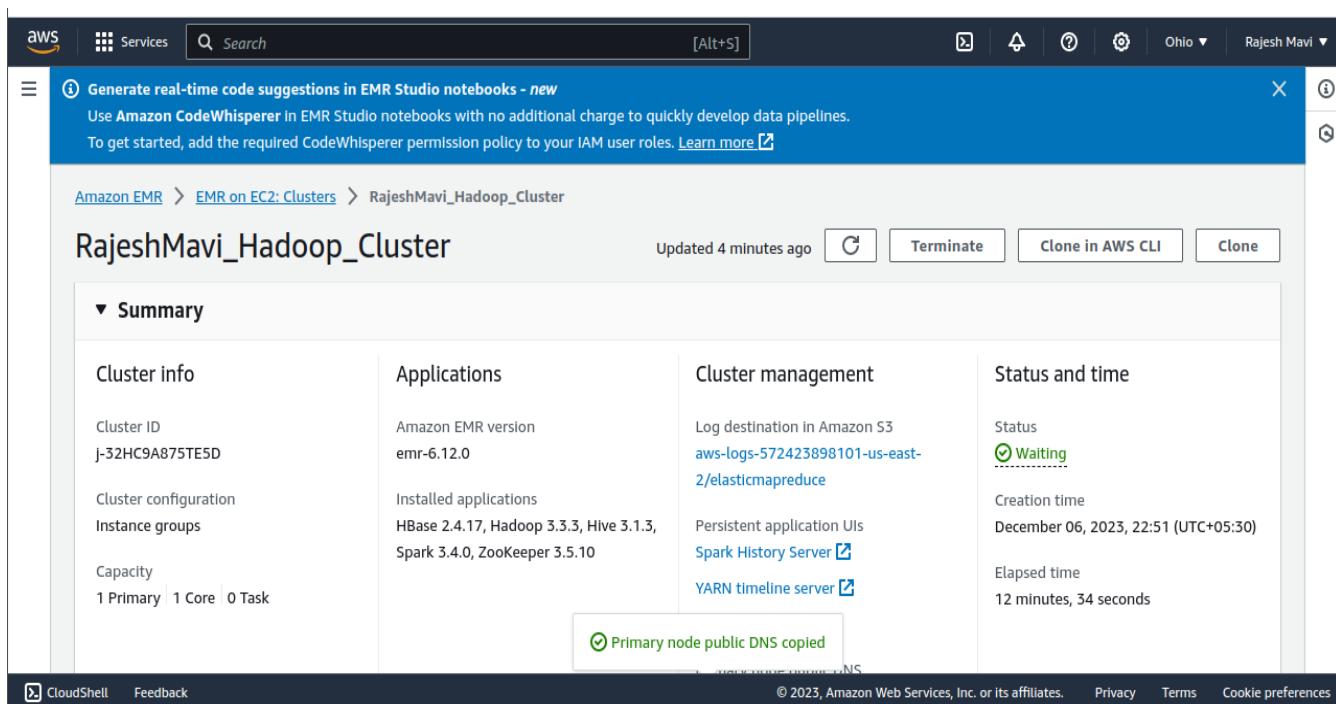


Fig 6: Started an EMR cluster to run jobs

Fig 7: Creating kafka topic user_Interaction for stream processing

```
[+]
hadoop@ip-172-31-8-158:~ hadoopuser@raja-Lenovo-V310-14ISK:~ hadoop@ip-172-31-8-158:~ hadoop@ip-172-31-8-158:~ hadoopuser@raja-Lenovo-V310-14I... ~

raja@raja-Lenovo-V310-14ISK:~ x      hadoop@ip-172-31-8-158:~/kafka_... x      hadoop@ip-172-31-8-158:~ x      hadoopuser@raja-Lenovo-V310-14I... x

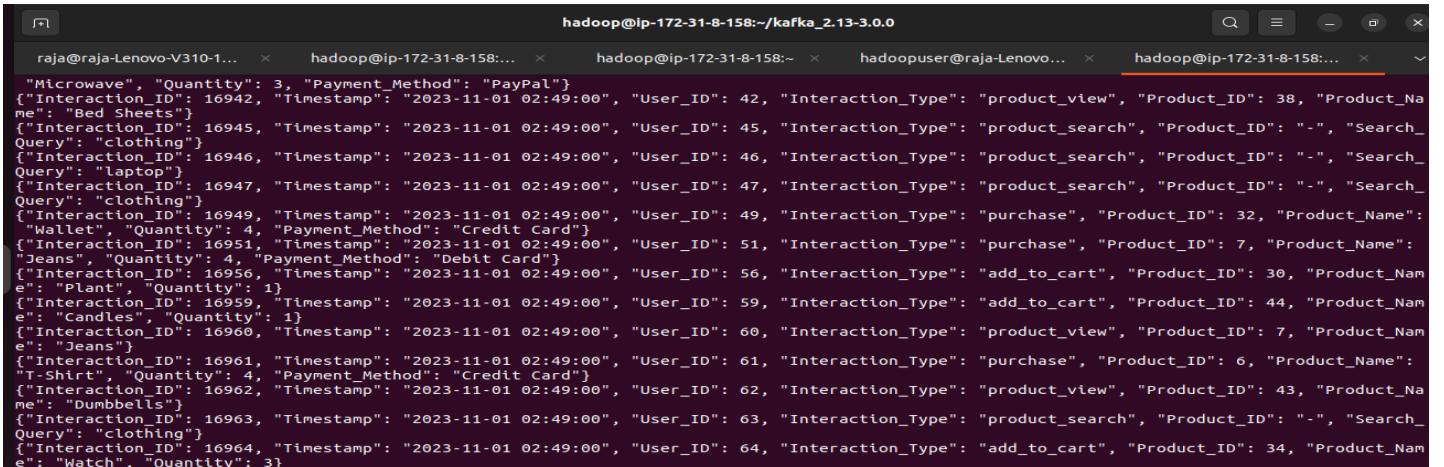
E:::::::::::::E M :::::M          M::::::M RR::::R          R:::::R
EEEEEEEEE::::: EEEEEE MMMMMMM          MMMMM RRRRRRR          RRRRR

[hadoop@ip-172-31-8-158 ~]$ ls -lthr
total 83M
-rw-rw-r-- 1 hadoop hadoop 83M Dec  6 17:39 kafka_2.13-3.0.0.tgz
drwxr-xr-x 8 hadoop hadoop 117 Dec  6 17:49 kafka_2.13-3.0.0
[hadoop@ip-172-31-8-158 ~]$ vim UserInteractionsProducerAws.py
[hadoop@ip-172-31-8-158 ~]$ pip install pandas numpy kafka-python boto3
Defaulting to user installation because normal site-packages is not writeable
Collecting pandas
  Downloading pandas-1.3.5-cp37-cp37m-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (11.3 MB)
    [██████████] 11.3 MB 2.6 MB/s
Requirement already satisfied: numpy in /usr/local/lib64/python3.7/site-packages (1.20.0)
Requirement already satisfied: kafka-python in ./local/lib/python3.7/site-packages (2.0.2)
Collecting boto3
  Downloading boto3-1.33.8-py3-none-any.whl (139 kB)
    [██████████] 139 kB 52.0 MB/s
Collecting python-dateutil>=2.7.3
  Downloading python_dateutil-2.8.2-py2.py3-none-any.whl (247 kB)
    [██████████] 247 kB 56.1 MB/s
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/site-packages (from pandas) (2023.3)
Collecting botocore<1.34.0,>=1.33.8
  Downloading botocore-1.33.8-py3-none-any.whl (11.8 MB)
    [██████████] 11.8 MB 50.3 MB/s
Collecting s3transfer<0.9.0,>=0.8.2
  Downloading s3transfer-0.8.2-py3-none-any.whl (82 kB)
    [██████████] 82 kB 277 kB/s
Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /usr/local/lib/python3.7/site-packages (from boto3) (1.0.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/site-packages (from python-dateutil>=2.7.3->pandas) (1.13.0)
Collecting urllib3<1.27,>=1.25.4; python_version < "3.10"
  Downloading urllib3-1.26.18-py3-none-any.whl (143 kB)
    [██████████] 143 kB 48.1 MB/s
Installing collected packages: python-dateutil, pandas, urllib3, botocore, s3transfer, boto3
Successfully installed boto3-1.33.8 botocore-1.33.8 pandas-1.3.5 python-dateutil-2.8.2 s3transfer-0.8.2 urllib3-1.26.18
[hadoop@ip-172-31-8-158 ~]$
```

Fig 8: Installing pandas, numpy, kafka-python, boto3

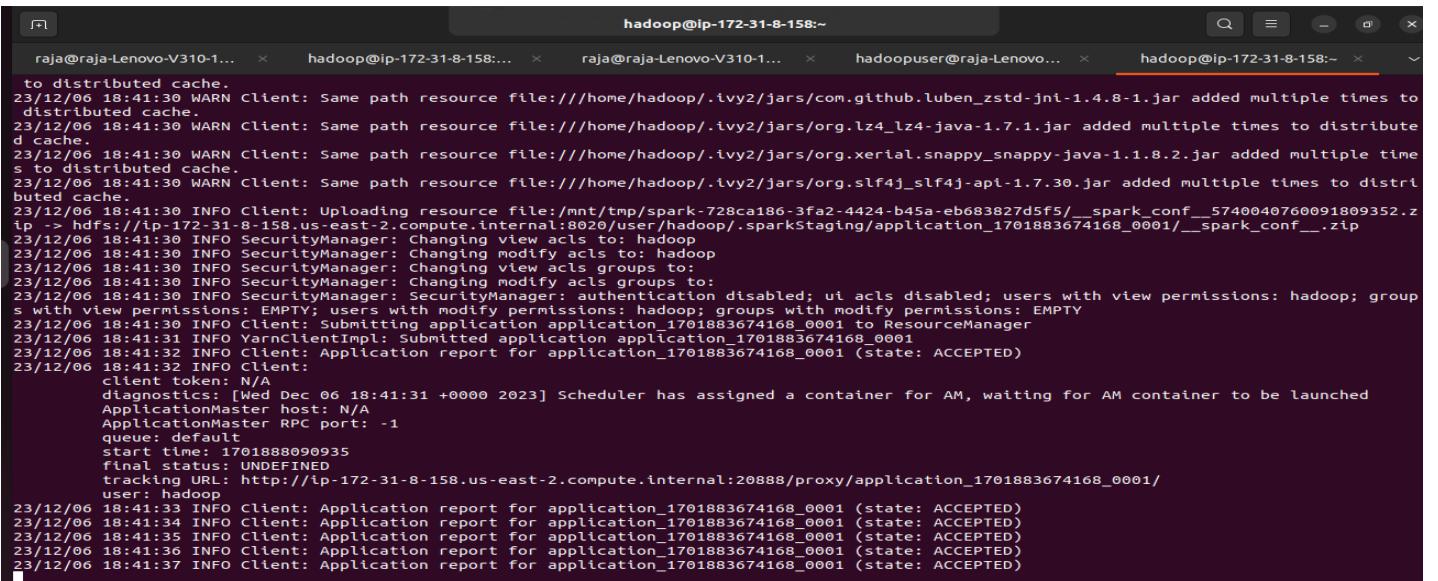
```
[hadoop@ip-172-31-8-158 ~]$ ls -ltrh
total 88M
-rw-rw-r-- 1 hadoop hadoop 83M Dec  6 17:39 kafka_2.13-3.0.0.tgz
drwxr-xr-x 8 hadoop hadoop 117 Dec  6 17:49 kafka_2.13-3.0.0
-rw-rw-r-- 1 hadoop hadoop 3.4K Dec  6 18:14 UserInteractionsProducerAws.py
-rw-rw-r-- 1 hadoop hadoop 8.8K Dec  6 18:22 UserProfiles_data.csv
-rw-rw-r-- 1 hadoop hadoop 1.4K Dec  6 18:23 ProductData_data.csv
-rw-rw-r-- 1 hadoop hadoop 5.0M Dec  6 18:23 UserInteractions_data.csv
[hadoop@ip-172-31-8-158 ~]$ pwd
/home/hadoop
[hadoop@ip-172-31-8-158 ~]$ rm -rf UserInteractionsProducerAws.py
[hadoop@ip-172-31-8-158 ~]$ vim UserInteractionsProducerAws.py
[hadoop@ip-172-31-8-158 ~]$ pw
-bash: pw: command not found
[hadoop@ip-172-31-8-158 ~]$ pwd
/home/hadoop
[hadoop@ip-172-31-8-158 ~]$ rm -rf UserInteractionsProducerAws.py
[hadoop@ip-172-31-8-158 ~]$ vim UserInteractionsProducerAws.py
[hadoop@ip-172-31-8-158 ~]$ python UserInteractionsProducerAws.py
```

Fig 9: Running Kafka producer UserInteractionsProducerAWS



```
hadoop@ip-172-31-8-158:~/kafka_2.13-3.0.0
raja@raja-Lenovo-V310-1...      hadoop@ip-172-31-8-158:...      hadoop@ip-172-31-8-158:~      hadoopuser@raja-Lenovo...      hadoop@ip-172-31-8-158:...
"Microwave", "Quantity": 3, "Payment_Method": "PayPal"} {"Interaction_ID": 16942, "Timestamp": "2023-11-01 02:49:00", "User_ID": 42, "Interaction_Type": "product_view", "Product_ID": 38, "Product_Name": "Bed Sheets"} {"Interaction_ID": 16945, "Timestamp": "2023-11-01 02:49:00", "User_ID": 45, "Interaction_Type": "product_search", "Product_ID": "-", "Search_Query": "Clothing"} {"Interaction_ID": 16946, "Timestamp": "2023-11-01 02:49:00", "User_ID": 46, "Interaction_Type": "product_search", "Product_ID": "-", "Search_Query": "laptop"} {"Interaction_ID": 16947, "Timestamp": "2023-11-01 02:49:00", "User_ID": 47, "Interaction_Type": "product_search", "Product_ID": "-", "Search_Query": "clothing"} {"Interaction_ID": 16949, "Timestamp": "2023-11-01 02:49:00", "User_ID": 49, "Interaction_Type": "purchase", "Product_ID": 32, "Product_Name": "Walls", "Quantity": 4, "Payment_Method": "Credit Card"} {"Interaction_ID": 16951, "Timestamp": "2023-11-01 02:49:00", "User_ID": 51, "Interaction_Type": "purchase", "Product_ID": 7, "Product_Name": "Jeans", "Quantity": 4, "Payment_Method": "Debit Card"} {"Interaction_ID": 16956, "Timestamp": "2023-11-01 02:49:00", "User_ID": 56, "Interaction_Type": "add_to_cart", "Product_ID": 30, "Product_Name": "Plant", "Quantity": 1} {"Interaction_ID": 16959, "Timestamp": "2023-11-01 02:49:00", "User_ID": 59, "Interaction_Type": "add_to_cart", "Product_ID": 44, "Product_Name": "Candles", "Quantity": 1} {"Interaction_ID": 16960, "Timestamp": "2023-11-01 02:49:00", "User_ID": 60, "Interaction_Type": "product_view", "Product_ID": 7, "Product_Name": "Jeans"} {"Interaction_ID": 16961, "Timestamp": "2023-11-01 02:49:00", "User_ID": 61, "Interaction_Type": "purchase", "Product_ID": 6, "Product_Name": "T-Shirt", "Quantity": 4, "Payment_Method": "Credit Card"} {"Interaction_ID": 16962, "Timestamp": "2023-11-01 02:49:00", "User_ID": 62, "Interaction_Type": "product_view", "Product_ID": 43, "Product_Name": "Dumbbells"} {"Interaction_ID": 16963, "Timestamp": "2023-11-01 02:49:00", "User_ID": 63, "Interaction_Type": "product_search", "Product_ID": "-", "Search_Query": "clothing"} {"Interaction_ID": 16964, "Timestamp": "2023-11-01 02:49:00", "User_ID": 64, "Interaction_Type": "add_to_cart", "Product_ID": 34, "Product_Name": "Watch", "Quantity": 3}
```

Fig 10: User interaction kafka topic data



```
hadoop@ip-172-31-8-158:~
raja@raja-Lenovo-V310-1...      hadoop@ip-172-31-8-158:...      raja@raja-Lenovo-V310-1...      hadoopuser@raja-Lenovo...      hadoop@ip-172-31-8-158:~
to distributed cache.
23/12/06 18:41:30 WARN Client: Same path resource file:///home/hadoop/.ivy2/jars/com.github.luben_zstd-jni-1.4.8-1.jar added multiple times to distributed cache.
23/12/06 18:41:30 WARN Client: Same path resource file:///home/hadoop/.ivy2/jars/org.lz4_lz4-java-1.7.1.jar added multiple times to distributed cache.
23/12/06 18:41:30 WARN Client: Same path resource file:///home/hadoop/.ivy2/jars/org.xerial.snappy_snappy-java-1.1.8.2.jar added multiple times to distributed cache.
23/12/06 18:41:30 WARN Client: Same path resource file:///home/hadoop/.ivy2/jars/org.slf4j_slf4j-api-1.7.30.jar added multiple times to distributed cache.
23/12/06 18:41:30 INFO Client: Uploading resource file:/mnt/tmp/spark-728ca186-3fa2-4424-b45a-eb683827df5f/_spark_conf_5740040760091809352.zip ip-> hdfs://ip-172-31-8-158.us-east-2.compute.internal:8020/user/hadoop/.sparkStaging/application_1701883674168_0001/_spark_conf_.zip
23/12/06 18:41:30 INFO SecurityManager: Changing view acls to: hadoop
23/12/06 18:41:30 INFO SecurityManager: Changing modify acls to: hadoop
23/12/06 18:41:30 INFO SecurityManager: Changing view acls groups to:
23/12/06 18:41:30 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: hadoop; groups with view permissions: EMPTY; users with modify permissions: hadoop; groups with modify permissions: EMPTY
23/12/06 18:41:30 INFO Client: Submitting application application_1701883674168_0001 to ResourceManager
23/12/06 18:41:31 INFO YarnClientImpl: Submitted application application_1701883674168_0001
23/12/06 18:41:32 INFO Client: Application report for application_1701883674168_0001 (State: ACCEPTED)
23/12/06 18:41:32 INFO Client: client token: N/A
diagnostics: [Wed Dec 06 18:41:31 +0000 2023] Scheduler has assigned a container for AM, waiting for AM container to be launched
ApplicationMaster host: N/A
ApplicationMaster RPC port: -1
queue: default
start time: 1701888090935
final status: UNDEFINED
tracking URL: http://ip-172-31-8-158.us-east-2.compute.internal:20888/proxy/application_1701883674168_0001/
user: hadoop
23/12/06 18:41:33 INFO Client: Application report for application_1701883674168_0001 (state: ACCEPTED)
23/12/06 18:41:34 INFO Client: Application report for application_1701883674168_0001 (state: ACCEPTED)
23/12/06 18:41:35 INFO Client: Application report for application_1701883674168_0001 (state: ACCEPTED)
23/12/06 18:41:36 INFO Client: Application report for application_1701883674168_0001 (state: ACCEPTED)
23/12/06 18:41:37 INFO Client: Application report for application_1701883674168_0001 (state: ACCEPTED)
```

Fig 11: Streaming kafkasparkstreamingCSVWrite job

```

hadoop@ip-172-31-8-158:~          raja@raja-Lenovo-V310-1...      raja@raja-Lenovo-V310-1...      hadoopuser@raja-Lenovo...      hadoop@ip-172-31-8-158:~

:: loading settings :: url = jar:file:/usr/lib/spark/jars/ivy-2.5.1.jar!/org/apache/ivy/core/settings/ivysettings.xml
Ivy Default Cache set to: /home/hadoop/.ivy2/cache
The jars for the packages stored in: /home/hadoop/.ivy2/jars
org.apache.spark#spark-sql-kafka-0-10_2.12 added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent-d8aec0ed-34e6-43e3-8495-485d3d99d424;1.0
  confs: [default]
    found org.apache.spark#spark-sql-kafka-0-10_2.12;3.1.2 in central
    found org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.1.2 in central
    found org.apache.kafka#kafka-clients;2.6.0 in central
    found com.github.luben#zstd-jni;1.4.8-1 in central
    found org.lz4#lz4-java;1.7.1 in central
    found org.xerial.snappy#snappy-java;1.1.8.2 in central
    found org.slf4j#slf4j-api;1.7.30 in central
    found org.spark-project.spark#unused;1.0.0 in central
    found org.apache.commons#commons-pool2;2.6.2 in central
:: resolution report :: resolve 926ms :: artifacts dl 22ms
  :: modules in use:
  com.github.luben#zstd-jni;1.4.8-1 from central in [default]
  org.apache.commons#commons-pool2;2.6.2 from central in [default]
  org.apache.kafka#kafka-clients;2.6.0 from central in [default]
  org.apache.spark#spark-sql-kafka-0-10_2.12;3.1.2 from central in [default]
  org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.1.2 from central in [default]
  org.lz4#lz4-java;1.7.1 from central in [default]
  org.slf4j#slf4j-api;1.7.30 from central in [default]
  org.spark-project.spark#unused;1.0.0 from central in [default]
  org.xerial.snappy#snappy-java;1.1.8.2 from central in [default]
-----
|           |           modules      ||   artifacts   | | | | |
|       conf    |   number| search|dwlded|evicted||   number|dwlded|
|       default |       9  |   0   |   0   |   0   ||   9   |   0   |
-----
:: retrieving :: org.apache.spark#spark-submit-parent-d8aec0ed-34e6-43e3-8495-485d3d99d424
  confs: [default]
  0 artifacts copied, 9 already retrieved (0kB/17ms)

```

Fig 12: Job starting

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	Start Time	End Time
application_1701883674168_0003	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:18:25 +0550 2023	Thu Dec 7 00:20:00 +0550 2023
application_1701883674168_0002	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:13:58 +0550 2023	Thu Dec 7 00:20:00 +0550 2023
application_1701883674168_0001	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:11:30 +0550 2023	Thu Dec 7 00:20:00 +0550 2023

Fig 13: Hadoop Yarn UI to track the spark streaming job

Spark Jobs (?)

User: hadoop
Total Uptime:
Scheduling Mode: FIFO
Active Jobs: 1

► Event Timeline
▼ Active Jobs (1)

Job Id (Job Group) ▾	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0 (f40fa7ad-eefb-4f37-ab31-45d9e3f6f66f)	id = b1ef9355-4f45-49c2-a5c1-d2ddf258a564 runld = f40fa7ad-eefb-4f37-ab31-45d9e3f6f66f ... start at NativeMethodAccessorImpl.java:0	2023/12/06 18:53:52	2.3 min	0/2	0/200

Fig 14: List of tasks spark job is assigned to finish

```

ion 20, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 14.0 in stage 1.0 (TID 14) in 291 ms on ip-172-31-7-27.us-east-2.compute.internal (executor 1) (13/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 21.0 in stage 1.0 (TID 21) (ip-172-31-9-149.us-east-2.compute.internal, executor 2, partition 21, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 10.0 in stage 1.0 (TID 10) in 616 ms on ip-172-31-9-149.us-east-2.compute.internal (executor 2) (14/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 22.0 in stage 1.0 (TID 22) (ip-172-31-9-149.us-east-2.compute.internal, executor 2, partition 22, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 8.0 in stage 1.0 (TID 8) in 635 ms on ip-172-31-9-149.us-east-2.compute.internal (executor 2) (15/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 23.0 in stage 1.0 (TID 23) (ip-172-31-9-149.us-east-2.compute.internal, executor 2, partition 23, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 9.0 in stage 1.0 (TID 9) in 638 ms on ip-172-31-9-149.us-east-2.compute.internal (executor 2) (16/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 24.0 in stage 1.0 (TID 24) (ip-172-31-7-27.us-east-2.compute.internal, executor 1, partition 24, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 17.0 in stage 1.0 (TID 17) in 175 ms on ip-172-31-7-27.us-east-2.compute.internal (executor 1) (17/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 25.0 in stage 1.0 (TID 25) (ip-172-31-9-149.us-east-2.compute.internal, executor 2, partition 25, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 21.0 in stage 1.0 (TID 21) in 127 ms on ip-172-31-9-149.us-east-2.compute.internal (executor 2) (18/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 26.0 in stage 1.0 (TID 26) (ip-172-31-7-27.us-east-2.compute.internal, executor 1, partition 26, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 18.0 in stage 1.0 (TID 18) in 179 ms on ip-172-31-7-27.us-east-2.compute.internal (executor 1) (19/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 27.0 in stage 1.0 (TID 27) (ip-172-31-9-149.us-east-2.compute.internal, executor 2, partition 27, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 13.0 in stage 1.0 (TID 13) in 485 ms on ip-172-31-9-149.us-east-2.compute.internal (executor 2) (20/200)
23/12/06 19:17:58 INFO TaskSetManager: Starting task 28.0 in stage 1.0 (TID 28) (ip-172-31-7-27.us-east-2.compute.internal, executor 1, partition 28, PROCESS_LOCAL, 7374 bytes)
23/12/06 19:17:58 INFO TaskSetManager: Finished task 20.0 in stage 1.0 (TID 20) in 195 ms on ip-172-31-7-27.us-east-2.compute.internal (executor 1) (21/200)

```

Fig 15: Logs of spark jobs showing completion of the job

```

raja@raja-Lenovo-V310-1...      hadoop@ip-172-31-8-158:~ 
ssl.truststore.password = null
ssl.truststore.type = JKS

23/12/06 19:18:04 WARN AdminClientConfig: The configuration 'key_deserializer' was supplied but isn't a known config.
23/12/06 19:18:04 WARN AdminClientConfig: The configuration 'value_deserializer' was supplied but isn't a known config.
23/12/06 19:18:04 WARN AdminClientConfig: The configuration 'enable.auto.commit' was supplied but isn't a known config.
23/12/06 19:18:04 WARN AdminClientConfig: The configuration 'max.poll.records' was supplied but isn't a known config.
23/12/06 19:18:04 WARN AdminClientConfig: The configuration 'auto.offset.reset' was supplied but isn't a known config.
23/12/06 19:18:04 INFO AppInfoParser: Kafka version: 2.6.0
23/12/06 19:18:04 INFO AppInfoParser: Kafka commitId: 62abed01bee039651
23/12/06 19:18:04 INFO AppInfoParser: Kafka startTimeMs: 1701890284030
23/12/06 19:18:04 WARN NetworkClient: [AdminClient clientId=adminclient-1] Connection to node -1 (localhost/127.0.0.1:9092) could not be established. Broker may not be available.
23/12/06 19:18:04 WARN NetworkClient: [AdminClient clientId=adminclient-1] Connection to node -1 (localhost/127.0.0.1:9092) could not be established. Broker may not be available.
23/12/06 19:18:04 WARN NetworkClient: [AdminClient clientId=adminclient-1] Connection to node -1 (localhost/127.0.0.1:9092) could not be established. Broker may not be available.
23/12/06 19:18:04 WARN NetworkClient: [AdminClient clientId=adminclient-1] Connection to node -1 (localhost/127.0.0.1:9092) could not be established. Broker may not be available.
23/12/06 19:18:04 WARN NetworkClient: [AdminClient clientId=adminclient-1] Connection to node -1 (localhost/127.0.0.1:9092) could not be established. Broker may not be available.

```

Fig 16: Showing the connected nodes

The screenshot shows the Hadoop Application History interface. On the left, there's a sidebar with 'Application History' (selected), 'About Applications' (with FINISHED, FAILED, KILLED options), and 'Tools'. The main area displays application details for 'application_1701883674168_0004'. The 'Application Overview' section includes:

- User: hadoop
- Name: KafkaSparkStreamingCSVWriter
- Application Type: SPARK
- Application Tags: (empty)
- Application Priority: 0 (Higher Integer value indicates higher priority)
- YarnApplicationState: FINISHED
- Queue: default
- FinalStatus Reported by AM: SUCCEEDED
- Started: Wed Dec 06 18:53:32 +0000 2023
- Launched: Wed Dec 06 18:53:32 +0000 2023
- Finished: Wed Dec 06 19:11:22 +0000 2023
- Elapsed: 17mins, 50sec
- Tracking URL: History
- Diagnostics: (empty)
- Unmanaged Application: false
- Application Node Label expression: <Not set>
- AM container Node Label expression: <DEFAULT_PARTITION>

Below this is a table for 'Attempt ID' showing one entry: 'appattempt_1701883674168_0004_000001' with 'Started' at 'Thu Dec 7 00:23:32 +0550 2023' and 'Logs' at 'http://ip-172-31-7-27.us-east-2.compute.internal:8042'. The bottom of the table says 'Showing 1 to 1 of 1 entries'.

Fig 17: Complete information about the spark job

```

raja@raja-Lenovo-V310-1...      hadoop@ip-172-31-8-158:~      raja@raja-Lenovo-V310-1...      hadoopuser@raja-Lenovo...      hadoop@ip-172-31-8-158:~ 
r
23/12/06 19:23:38 INFO ServerInfo: Adding filter to /static/sql: org.apache.hadoop.yarn.server.webproxy.amfilter.AmIpFilter
23/12/06 19:23:38 INFO ResolveWriteToStream: Checkpoint root /path/to/checkpoint resolved to hdfs://ip-172-31-8-158.us-east-2.compute.internal:8020/path/to/checkpoint.
23/12/06 19:23:38 WARN ResolveWriteToStream: spark.sql.adaptive.enabled is not supported in streaming DataFrames/Datasets and will be disabled.
23/12/06 19:23:38 INFO MicroBatchExecution: Starting Tid = bief9355-4f45-49c2-a5c1-d2ddf258a564, runId = e9869d2f-ec30-470f-81a5-519a36276c14]
23/12/06 19:23:38 INFO UserHdfs: [/ip-172-31-8-158.us-east-2.compute.internal:8020/path/to/checkpoint].
23/12/06 19:23:38 INFO MicroBatchExecution: Reading table [org.apache.spark.sql.KafkaSourceProvider$KafkaTable@3f8e1464] from DataSou
rcceV2 named 'kafka' [org.apache.spark.sql.KafkaSourceProvider$KafkaTable@48a4fa32]
23/12/06 19:23:38 INFO OffsetSeqLog: BatchIds found from listing: 0
23/12/06 19:23:38 INFO OffsetSeqLog: Getting latest batch 0
23/12/06 19:23:38 INFO OffsetSeqLog: BatchIds found from listing: 0
23/12/06 19:23:38 INFO OffsetSeqLog: Getting latest batch 0
23/12/06 19:23:38 INFO CommitLog: BatchIds found from listing: 0
23/12/06 19:23:38 INFO CommitLog: Getting latest batch 0
23/12/06 19:23:38 INFO MicroBatchExecution: Resuming at batch 1 with committed offsets {KafkaV2[Subscribe[User_Interactions]]: {"User_Interacti
ons": [{"2": "463200", "1": "463178", "0": "462599"}]} and available offsets {KafkaV2[Subscribe[User_Interactions]]: {"User_Interactions": {"2": "463200", "1": "4
63178", "0": "462599"}]}
23/12/06 19:23:38 INFO MicroBatchExecution: Stream started from {KafkaV2[Subscribe[User_Interactions]]: {"User_Interactions": {"2": "463200", "1": "4
63178", "0": "462599"}}}
23/12/06 19:23:38 INFO AdminClientConfig: AdminClientConfig values:
  bootstrap.servers = [localhost:9092]
  client.dns.lookup = use_all_dns_ips
  client.id =
  connections.max.idle.ms = 300000
  default.apt.timeout.ms = 600000
  metadata.max.age.ms = 300000
  metric.reporters = []
  metrics.num.samples = 2
  metrics.recording.level = INFO
  metrics.sample.window.ms = 30000
  receive.buffer.bytes = 65536
  reconnect.backoff.ms = 1000
  reconnect.backoff.ms = 50
  request.timeout.ms = 30000

```

Fig 18: Reading the data from the Kafka topic by the streaming spark job

Spark Jobs (1)

User: hadoop
Total Uptime: 2.7 min
Scheduling Mode: FIFO
Completed Jobs: 1

► Event Timeline
▼ Completed Jobs (1)

Job Id (Job Group) ▾	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0 (f0ae766c-12fd-4db4-842b-96acb36d9090)	id = b1ef9355-4f45-49c2-a5c1-d2ddf258a564 runid = f0ae766c-12fd-4db4-842b-96acb36d90... start at NativeMethodAccessorImpl.java:0	2023/12/06 19:17:53	9 s	1/1 (1 skipped)	200/200

Fig 19: Tracking the list of tasks that has executed successfully.

Details for Stage 1 (Attempt 0)

Resource Profile Id: 0
Total Time Across All Tasks: 47 s
Locality Level Summary: Process local: 200
Associated Job Ids: 0

► DAG Visualization
► Show Additional Metrics
► Event Timeline

Summary Metrics for 200 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	23.0 ms	43.0 ms	63.0 ms	0.4 s	2 s
GC Time	0.0 ms	0.0 ms	0.0 ms	0.0 ms	0.2 s

Aggregated Metrics by Executor

Tasks (200)

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Errors
0	0	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
1	1	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
2	2	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
3	3	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
4	4	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
5	5	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
6	6	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
7	7	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
8	8	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:57	0.5 s		
9	9	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr	2023-12-07	0.6 s		

Fig 20: Details of spark job stages

Tasks (200)

Index	Task ID	Attempt	Status	Locality level	Executor ID	Host	Logs	Launch Time	Duration	GC Time	Errors
0	0	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
1	1	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
2	2	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
3	3	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
4	4	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
5	5	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
6	6	0	SUCCESS	PROCESS_LOCAL	1	ip-172-31-7-27.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
7	7	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:54	2 s	0.2 s	
8	8	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr stdout	2023-12-07 00:47:57	0.5 s		
9	9	0	SUCCESS	PROCESS_LOCAL	2	ip-172-31-9-149.us-east-2.compute.internal	stderr	2023-12-07	0.6 s		

Fig 21: All the details of number of tasks executed on a node

The screenshot shows the Hadoop application history interface. On the left, there's a sidebar with a 'hadoop' logo, a 'Application History' section (with 'About Applications' dropdown and 'FINISHED', 'FAILED', 'KILLED' buttons), and a 'Tools' section. The main area is titled 'All Applications' and displays a table of application logs. The table has columns: ID, User, Name, Application Type, Application Tags, Queue, Application Priority, Start Time, and End Time. There are seven entries listed, all from the user 'hadoop' and of type 'SPARK', all named 'KafkaSparkStreamingCSVWriter'. The start times range from Dec 7, 2023, at 00:59:08 to 00:11:30.

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	Start Time	End Time
application_1701883674168_0007	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:59:08 +0550 2023	Thu Dec 7 00:59:20 +0550 2023
application_1701883674168_0006	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:53:25 +0550 2023	Thu Dec 7 00:53:26 +0550 2023
application_1701883674168_0005	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:47:31 +0550 2023	Thu Dec 7 00:47:32 +0550 2023
application_1701883674168_0004	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:23:32 +0550 2023	Thu Dec 7 00:23:33 +0550 2023
application_1701883674168_0003	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:18:25 +0550 2023	Thu Dec 7 00:18:26 +0550 2023
application_1701883674168_0002	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:13:58 +0550 2023	Thu Dec 7 00:13:59 +0550 2023
application_1701883674168_0001	hadoop	KafkaSparkStreamingCSVWriter	SPARK		default	0	Thu Dec 7 00:11:30 +0550 2023	Thu Dec 7 00:11:31 +0550 2023

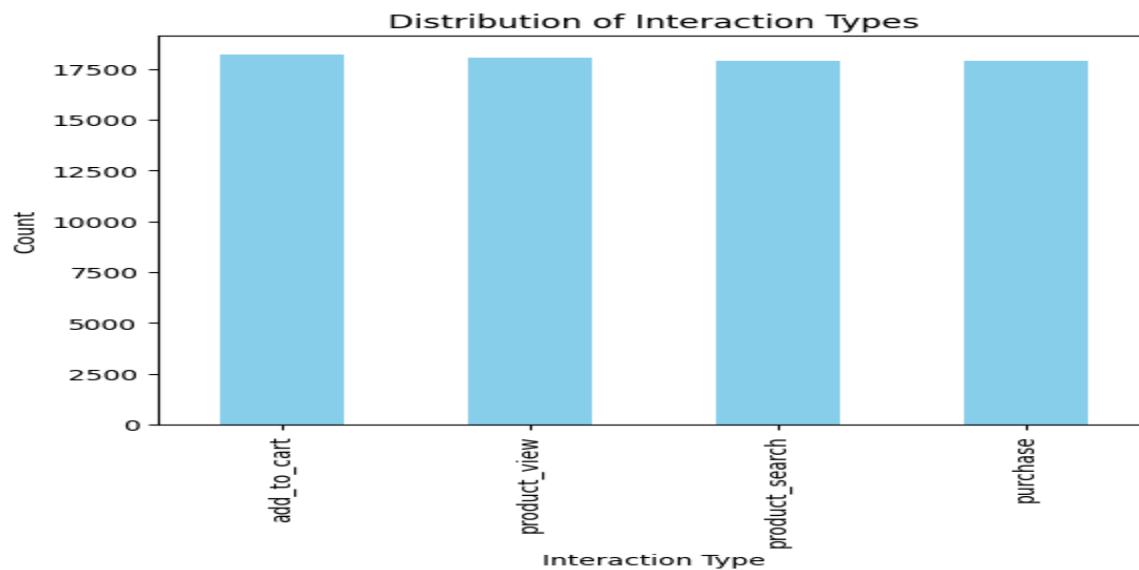
Fig 22: History of the all streaming spark jobs

```
-rw-rw-r-- 1 hadoopuser hadoopuser 5.9K Dec 7 03:48 UserInteractions_data_20231207_0348.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.7K Dec 7 03:49 UserInteractions_data_20231207_0349.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.7K Dec 7 03:50 UserInteractions_data_20231207_0350.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 5.9K Dec 7 03:51 UserInteractions_data_20231207_0351.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.7K Dec 7 03:52 UserInteractions_data_20231207_0352.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 5.9K Dec 7 03:53 UserInteractions_data_20231207_0353.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 03:54 UserInteractions_data_20231207_0354.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 03:55 UserInteractions_data_20231207_0355.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 5.9K Dec 7 03:56 UserInteractions_data_20231207_0356.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 03:57 UserInteractions_data_20231207_0357.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 03:58 UserInteractions_data_20231207_0358.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 03:59 UserInteractions_data_20231207_0359.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.2K Dec 7 04:00 UserInteractions_data_20231207_0400.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:01 UserInteractions_data_20231207_0401.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 04:02 UserInteractions_data_20231207_0402.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.7K Dec 7 04:03 UserInteractions_data_20231207_0403.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 04:04 UserInteractions_data_20231207_0404.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 04:05 UserInteractions_data_20231207_0405.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:06 UserInteractions_data_20231207_0406.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:07 UserInteractions_data_20231207_0407.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 04:08 UserInteractions_data_20231207_0408.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 04:09 UserInteractions_data_20231207_0409.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:10 UserInteractions_data_20231207_0410.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 04:11 UserInteractions_data_20231207_0411.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.9K Dec 7 04:12 UserInteractions_data_20231207_0412.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.9K Dec 7 04:13 UserInteractions_data_20231207_0413.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 04:14 UserInteractions_data_20231207_0414.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:15 UserInteractions_data_20231207_0415.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:16 UserInteractions_data_20231207_0416.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.2K Dec 7 04:17 UserInteractions_data_20231207_0417.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 04:18 UserInteractions_data_20231207_0418.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 04:19 UserInteractions_data_20231207_0419.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.1K Dec 7 04:20 UserInteractions_data_20231207_0420.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.8K Dec 7 04:21 UserInteractions_data_20231207_0421.csv
-rw-rw-r-- 1 hadoopuser hadoopuser 6.0K Dec 7 04:31 UserInteractions_data_20231207_0431.csv
```

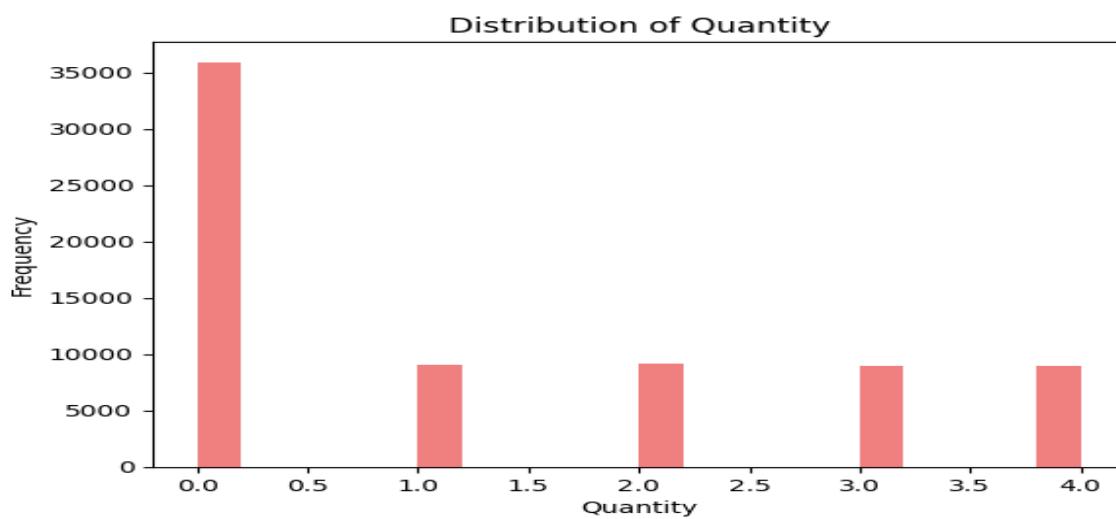
Fig 23: List of all files generated

2.5 Data Visualization

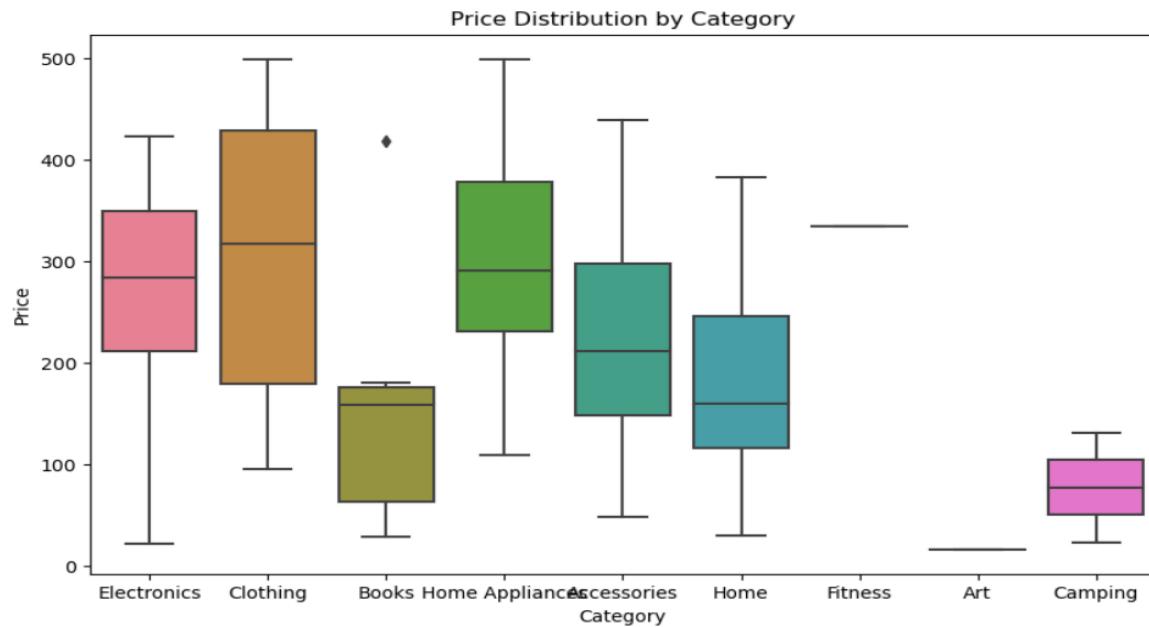
Distribution of interaction types: The graph illustrates the distribution of user interactions, including "Add to Cart," "Product View," "Product Search," and "Purchase," with the count represented on the y-axis. This visual representation provides insights into the varying engagement levels across these key interaction types, offering a comprehensive overview of user behavior within the given context.



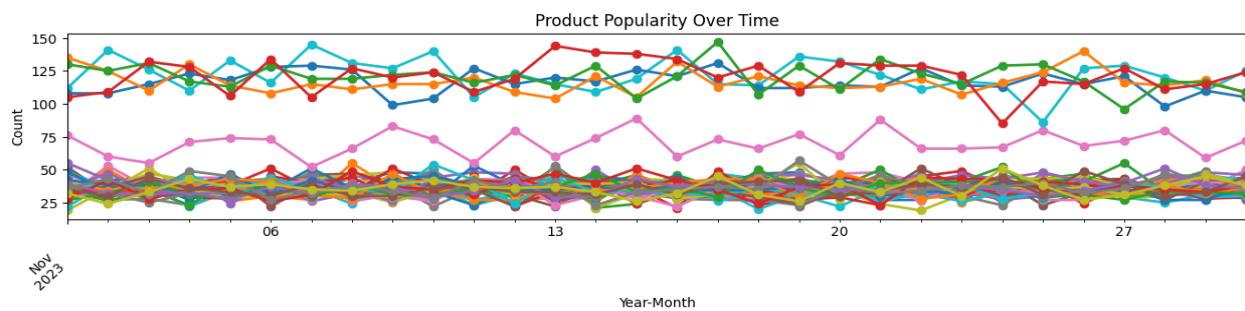
Distribution of quantity: The graph depicts a frequency distribution based on quantities, with the x-axis representing different quantity values and the y-axis indicating the frequency of occurrences. This visual representation offers insights into the distribution pattern and prevalence of specific quantity values within the analyzed dataset.



Quantity by interaction type: The box plot graph illustrates the distribution of quantities for different interaction types, emphasizing a focus on "Purchase" and "Add to Cart." The plot highlights the central tendency, spread, and potential outliers in the quantity data for these specific interactions, with a notable emphasis on the absence of quantities for "Product Search" and "Product View," marked as zero. This visual analysis provides a understanding of the variation in quantities associated with distinct user interactions, particularly emphasizing the transactional nature of "Purchase" and "Add to Cart."



Products Popularity over time: The graph showcases the popularity of products over a one-month period, utilizing the x-axis for time intervals and the y-axis to represent order counts. This visual representation allows for a temporal analysis of how various products have performed in terms of sales or orders, offering insights into their relative popularity dynamics over the specified timeframe.



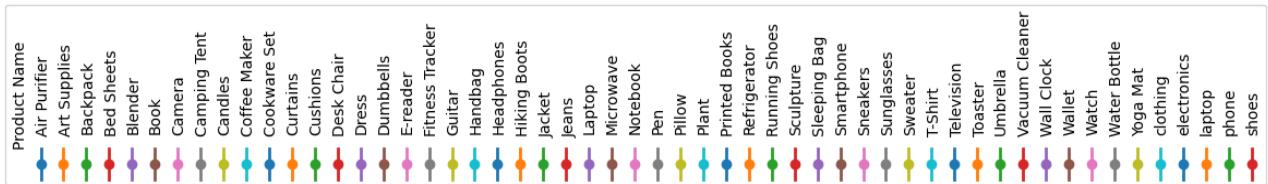
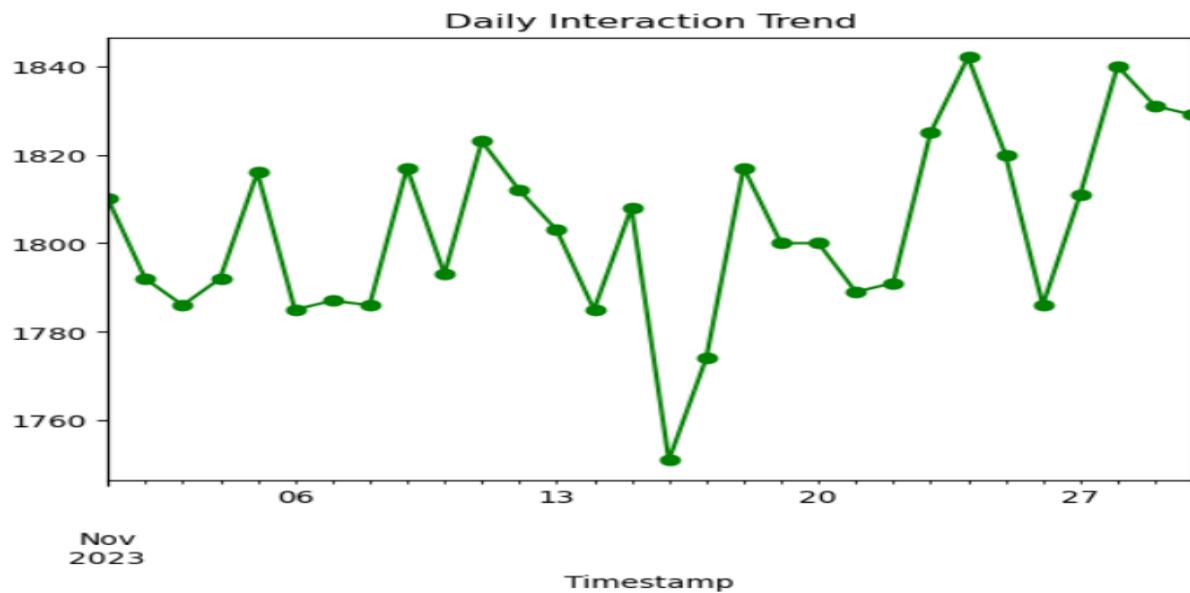


Fig 24: Legend for Products Popularity over time

Daily Interaction trend: This graph aims to depict how user interactions have fluctuated on a daily basis over the past month, offering a visual representation of the overall trend. Visualizing daily interaction trends is useful for businesses to track user engagement patterns, optimize marketing strategies, and make informed decisions about product performance and resource allocation.



Product Popularity Forecasting (MACHINE LEARNING): Used the ARIMA (AutoRegressive Integrated Moving Average) time series forecasting model from the statsmodel library to forecast the popularity of three products (laptop, phone, shoes) based on historical data, plots the historical data along with the forecasted values, and visualizes the product popularity trends over time. This visualization provides insights into the predicted trends and variations in the popularity of the specified products, aiding in decision-making and planning based on anticipated future product demand.

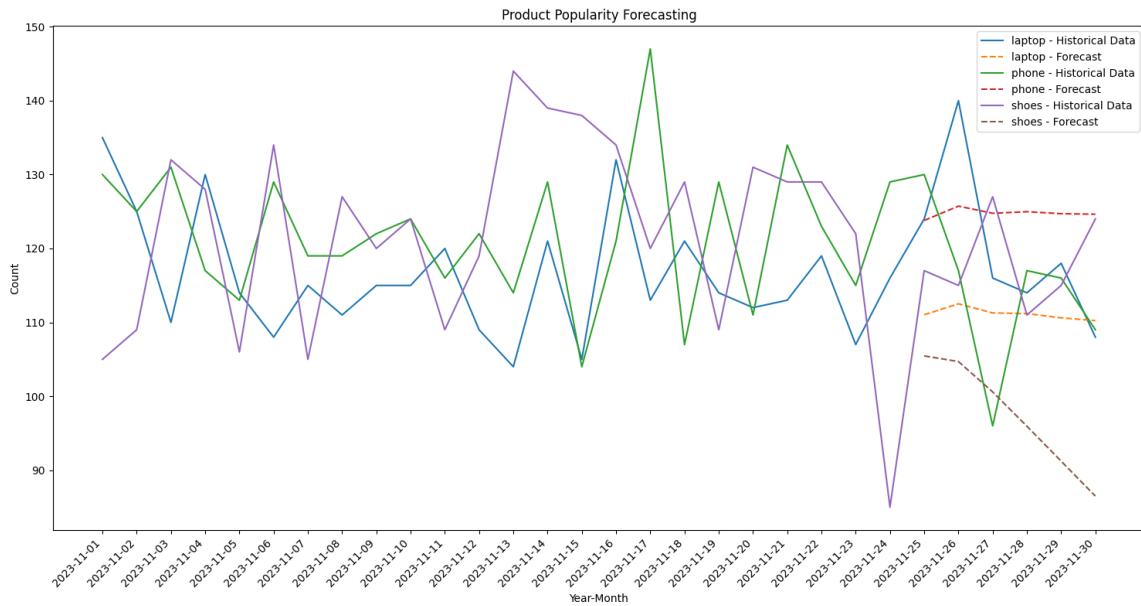


Fig 25: Product Popularity Forecasting

III: DIFFICULTIES FACED

1. Managing real-time data ingestion from S3, navigating challenges related to data consistency, latency, and optimizing transfer protocols.
2. Architecting Spark streaming jobs for AWS, addressing complexities in scalability, fault tolerance, and resource allocation across heterogeneous nodes.
3. Customizing an EMR cluster involves intricate configuration, dependency management, and the optimization of resource utilization.
4. Encountered challenges while configuring and publishing JSON-formatted messages to a Kafka topic. Additionally, faced difficulties in reading these messages within a PySpark Streaming job for the purpose of creating minute-wise aggregations.
5. Conducting a literature review poses challenges in identifying pertinent research papers on streaming data processing, custom cluster design, and live data from S3 for a comprehensive scholarly evaluation.
6. Considered various tools for distributed and real-time processing, including Hadoop MapReduce, Apache Flink, and Apache Storm. After meticulous evaluation, the chosen tool is PySpark Streaming with Kafka, leveraging PySpark's simplicity for real-time processing and Kafka's efficiency for event streaming.

IV: CONCLUSION

In conclusion, our project has successfully delivered a fully functional simulation that not only accommodates user interactions but also ensures efficient storage of both structured and unstructured data.

The implementation of real-time data processing has empowered continuous user behavior analysis, offering valuable insights into patterns, preferences, and product performance. With user segmentation, we gain a detailed understanding of different user types, while visualizations and reports provide actionable information for decision-making. Scalable system design positions us well to handle increasing data volumes in the future. Overall, this project establishes a robust foundation for comprehensive data-driven insights and user-focused strategies.

V: FUTURE SCOPE

- 1. Advanced Analytics:** Explore and implement advanced analytics techniques, such as machine learning algorithms, for deeper insights into user behavior, predictive analytics, or anomaly detection.
- 2. Integration with External Data Systems:** Extend the system's capabilities by integrating with external data sources or third-party APIs to enrich the dataset and provide a more comprehensive analysis.
- 3. Real-time Decision Support:** Strengthen the decision-making process by incorporating real-time decision support mechanisms, allowing stakeholders to make informed decisions on the fly.
- 4. Behavioral Economics Analysis:** Incorporate principles of behavioral economics into the analysis to understand how psychological factors influence user decisions, leading to more effective design and engagement strategies.

VI: REFERENCES

- [1] [Python](#)
 - [2] [Apache Hive](#)
 - [3] [Matplotlib](#)
 - [4] [Apache PySpark](#)
 - [5] [Apache Kafka](#)
 - [6] [Hadoop](#)
 - [7] [Zookeeper](#)
 - [8] [AWS](#)
- [9] Le, T. M., & Liaw, S.-Y. (2017). *Effects of Pros and Cons of Applying Big Data Analytics to Consumers' Responses in an E-Commerce Context*. Sustainability, 9(5), 798. [Mdpi](#)
- [10] Evangelin, R., & Shanmugam, V. (2022). *Mechanism of Big Data Analytics in Consumer Behavior on Online Shopping*. May 2022. [ResearchGate](#)
- [11] Rongrui Yu, Chunqiong Wu, Bingwen Yan, Baoqin Yu, Xiukao Zhou, Yanliang Yu, Na Chen, (2021). "Analysis of the Impact of Big Data on E-Commerce in Cloud Computing Environment", Article ID 5613599. [Hindawi](#)
- [12] Jui-Chan Huang, Po-Chang Ko, Cher-Min Fong, Sn-Man Lai, Hsin-Hung Chen, Ching-Tang Hsieh,(2021) "Statistical Modeling and Simulation of Online Shopping Customer Loyalty Based on Machine Learning and Big Data Analysis". [Hindawi](#)
- [13] Asle Fagerstrøm, Erik Arntzen, Gordon Foxall, (2011). "A study of preferences in a simulated online shopping experiment". [ResearchGate](#)