



INSTITUTE OF BUSINESS ADMINISTRATION
KARACHI

PROJECT REPORT

(Big Data Analytics)

[Big Data Analysis of English Tweets
& Comparison of Big Data Tools]

| S.NO | Name | ERP |
|------|-----------------------|-------|
| 1 | Bilawal Ali | 17278 |
| 2 | Muhammed Imran Magsi | 17670 |
| 3 | Muhammad Omar Qureshi | 07133 |

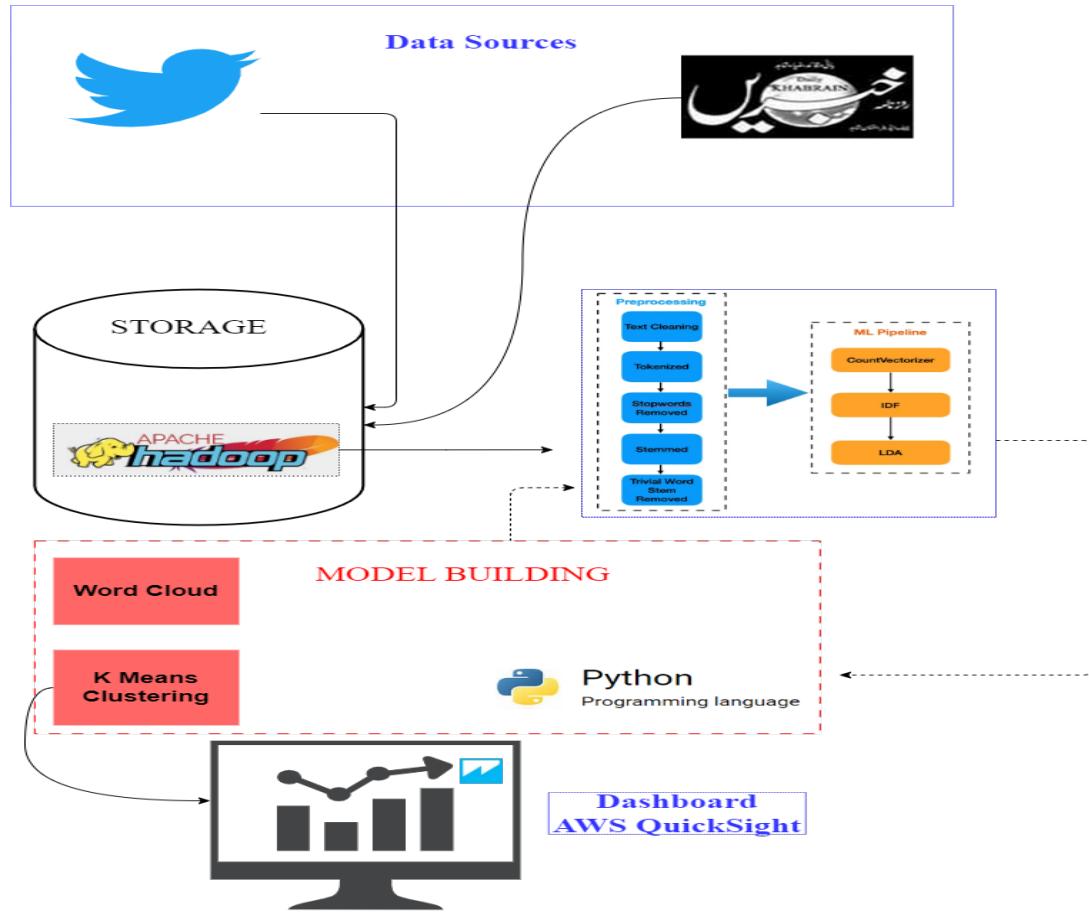
TABLE OF CONTENTS

| | | |
|-----|--|----|
| 1 | Project # 01:..... | 3 |
| 2 | Primary Project: Urdu Tweet Analysis | 3 |
| 2.1 | Pipeline..... | 3 |
| 2.2 | Project Abstract:..... | 3 |
| 2.3 | Preprocessing news raw text..... | 4 |
| 2.4 | TF-IDF..... | 4 |
| 2.5 | K-Means Clustering..... | 4 |
| 2.6 | Knee Point | 5 |
| 2.7 | Word Cloud..... | 5 |
| | Interface:..... | 6 |
| 3 | Project #02:..... | 7 |
| 4 | Secondary Project: Apache Hive VS Apache Spark VS Apache Drill | 7 |
| 4.1 | Data Pipeline | 7 |
| 4.2 | Dataset Introduction | 7 |
| 4.3 | Spark Introduction..... | 8 |
| 4.4 | Spark on Docker | 8 |
| 4.5 | Spark Commands & Analytics..... | 10 |
| 5 | Apache Hive: Bilawal..... | 16 |
| 5.1 | What is Hive..... | 16 |
| 5.2 | Steps..... | 17 |

1 PROJECT #01:

2 PRIMARY PROJECT: URDU TWEET ANALYSIS

2.1 PIPELINE



2.2 PROJECT ABSTRACT:

We plan to build a dockerized application that first fetches live tweets from Twitter and RSS feeds of different news articles/sources. We are saving all these tweets in .csv file format. Once we have the data, we start Pre-processing files. During Pre-processing we are removing punctuations, English texts, numbers and special characters then we start removing Stop words. Once we are done with Pre-processing, we do lemmatization and tokenization of Urdu words. Then we apply Term-Frequency and Inverse document frequency to evaluate how relevant a word is to a document in a collection of documents. Now comes the modelling, we are applying K-means Algorithm to cluster the tweets on the basis of TF-IDF vectors, here

Knee Locator finds the optimal value of K. Then comes the Principal Component Analysis (PCA) that we use for denoising and data compression, and finally we get the clusters. We are creating Tag Clouds for all clusters formed during K-means also showing Tweets of all clusters on a separate page. There are three ways a user can cluster or fetch data. First, they can see the default clustering the application does by itself (Live data fetching) secondly users can upload the file, thirdly they can also search a keyword and all tweets containing the keyword will be fetched live and all the modeling will start from the beginning.

2.3 PREPROCESSING NEWS RAW TEXT

We have used UrduHack library to preprocess our Urdu text data, and have removed any emojis, numbers, English alphabets, links and stop-words before using the data for any visual analysis. Here are some functions which we have used from `urduhack.preprocessing` module while processing data.

- `urduhack_normalize`
- `normalize_whitespace`
- `remove_punctuation`
- `replace_urls`
- `replace_emails`
- `replace_numbers`
- `remove_english_alphabets`
- `remove_diacritics`

Besides, we have used python regular expressions and lemmatization techniques from UrduHack library to convert the words to their basic form.

2.4 TF-IDF

We have utilized Term Frequency-Inverse Document Frequency as a text vectorizer to transform the text into a usable vector. This vector is then used to generate clusters.

2.5 K-MEANS CLUSTERING

We used the K-means clustering algorithm for generating clusters. K-Means aims to partition a set of observations into several clusters.

2.6 KNEE POINT

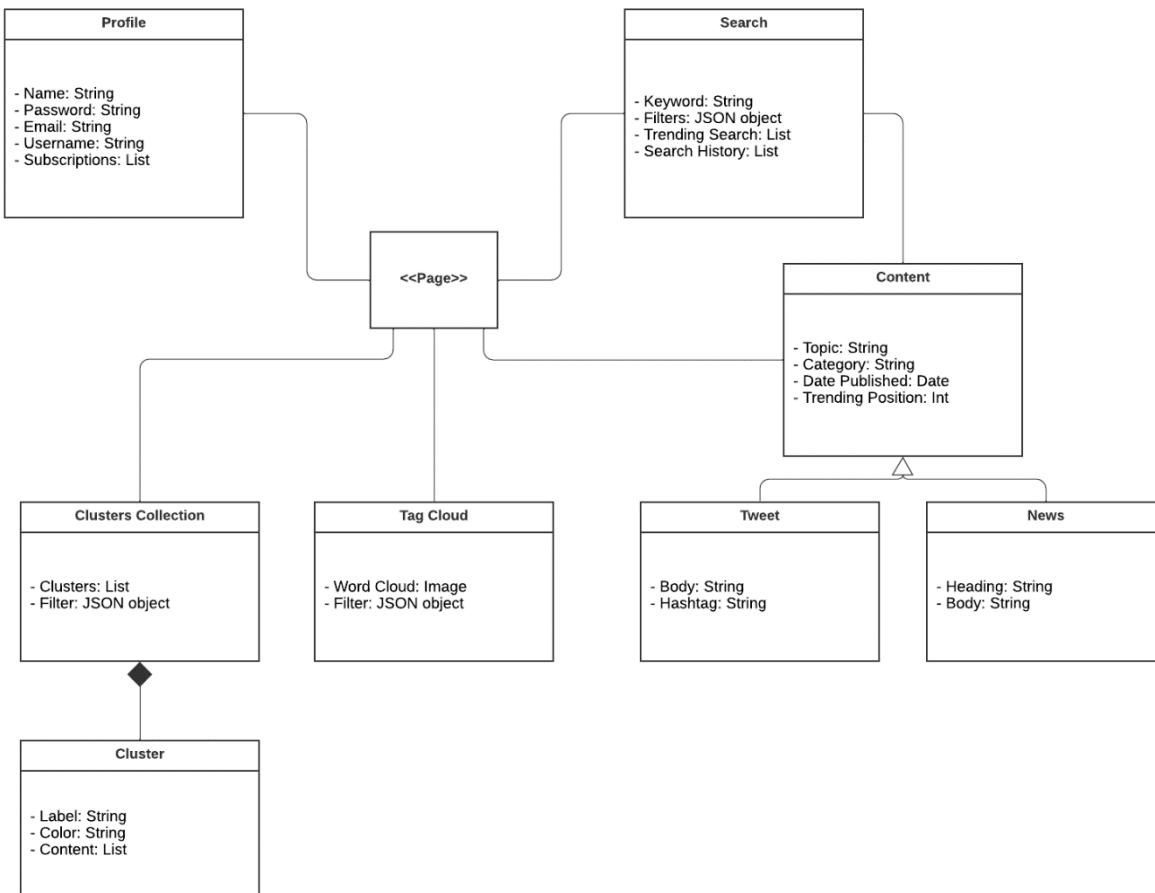
To find the optimal number of clusters, we are using Knee Point Locator which determines the suitable number of clusters using SSE inertia values.

2.7 WORD CLOUD

Given a textual data, the word-cloud function returns a dictionary of objects with unique words as keys, and their font size (directly proportional to the frequency) in word-cloud as value. This dictionary object is used in node.js to create and visualize the cluster.

Note: **Other details of this project is in .ipynb notebook**

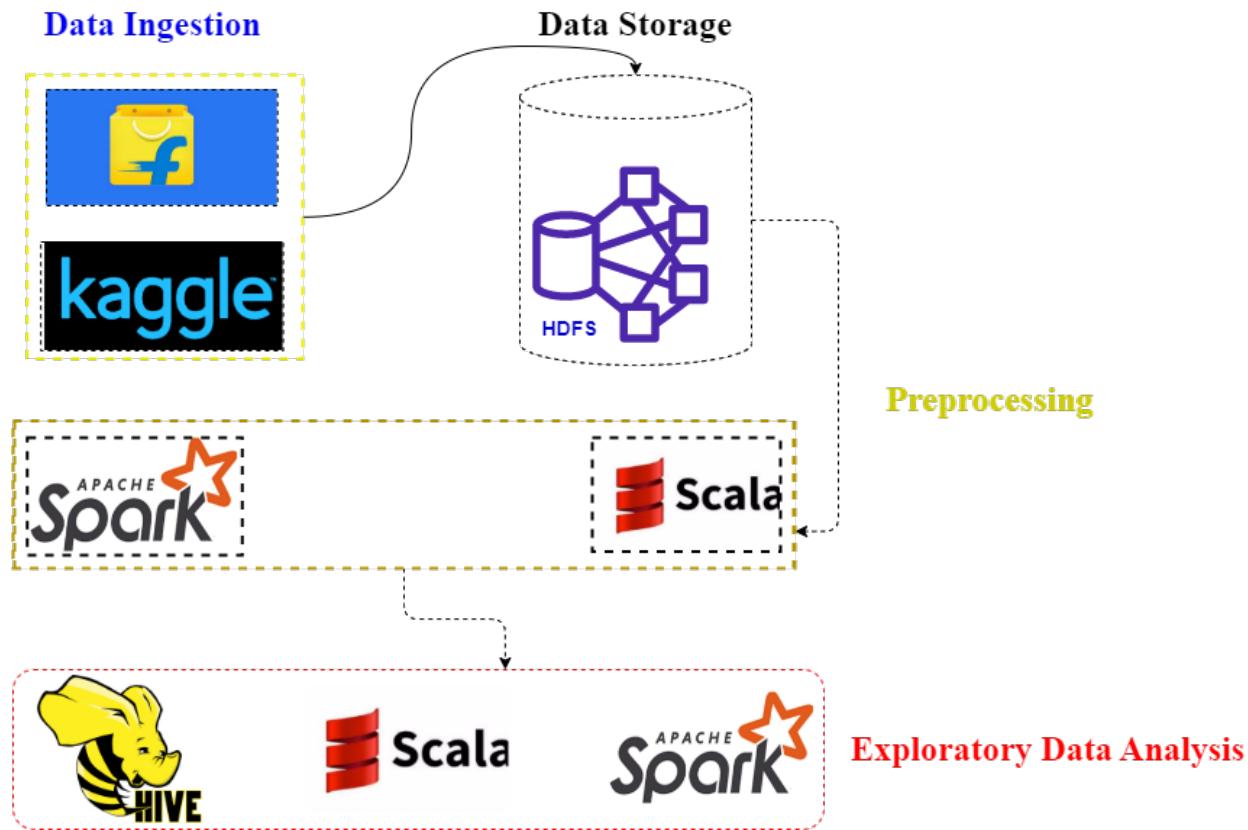
INTERFACE:



3 PROJECT #02:

4 SECONDARY PROJECT: APACHE HIVE VS APACHE SPARK VS APACHE DRILL

4.1 DATA PIPELINE



4.2 DATASET INTRODUCTION

Size: 3GB

Mobile prices extracted from flipkart a day before big billion days.

Categories

Mobiles, Laptops, Tvs, Tablets

Columns

- name - Name of the item
- offer_price - Current price with offer
- original_price - Orginal price of the item without any offer
- off_now - % of off on item

- total_ratings - Number of ratings given
- total_reviews - Number of reviews written
- rating - Final rating
- description - Description of the product
- created_at - Dataset created time
- u_id - Unique id
- item_link - link of the particular item

4.3 SPARK INTRODUCTION

Spark is an open-source distributed computing system designed to process large-scale data sets and perform data analysis at high speeds. It provides a unified and flexible framework for data processing, enabling developers to efficiently process data across various data sources and execute complex data transformations and computations.

One of the key features of Spark is its ability to handle data processing tasks in-memory, which significantly enhances the overall performance compared to traditional disk-based processing systems. Spark's in-memory computing allows for iterative and interactive data analysis, making it well-suited for real-time analytics, machine learning, and graph processing.

Spark offers a rich set of APIs and libraries that support various programming languages, including Scala, Java, Python, and R. This versatility enables developers to work with Spark using their preferred programming language and leverage the vast ecosystem of tools and libraries available in the Spark ecosystem.

With its distributed computing model, Spark enables parallel processing of data across clusters of machines, making it highly scalable and capable of handling large-scale data processing tasks. It provides fault tolerance and resilience by automatically recovering from failures, ensuring the reliability of data processing workflows.

Furthermore, Spark provides a wide range of high-level abstractions, such as DataFrames and Datasets, which simplify the development of complex data processing pipelines. These abstractions offer a more intuitive and declarative programming model, allowing developers to focus on the logic of their data analysis tasks rather than dealing with low-level details.

Overall, Spark has revolutionized big data processing and analytics, empowering organizations to efficiently extract valuable insights from massive datasets. Its speed, scalability, and rich feature set make it a popular choice for data scientists, engineers, and analysts seeking to harness the power of distributed computing for their data-driven applications.

4.4 SPARK ON DOCKER

1. Install Docker for Desktop from (<https://docs.docker.com/desktop/windows/install/>)
2. Clone the Spark using the command line: `git clone https://github.com/big-data-europe/docker-spark.git`
3. Navigate to the cloned folder docker-spark and Execute: `docker-compose up`

4. Copy the directory (here tvs) from local machine to the Spark using : `docker cp D:\big.data\docker-hadoop\spark\docker-spark\tvs 76a4c7ff984b:home` (where **76a4c7ff984b** is the Spark container's address running on the Docker)
5. Bash into the spark-master: `docker exec -it spark-master /bin/bash`
6. Use the Spark shell (*Scala*) by executing : `/spark/bin/spark-shell --master=local spark://spark-master:7077`

```
D:\big.data\docker-hadoop\spark\docker-spark>docker cp D:\big.data\docker-hadoop\spark\docker-spark\tvs 76a4c7ff984b:home

D:\big.data\docker-hadoop\spark\docker-spark>docker cp D:\big.data\docker-hadoop\spark\docker-spark\mobiles 76a4c7ff984b:home

D:\big.data\docker-hadoop\spark\docker-spark>docker cp D:\big.data\docker-hadoop\spark\docker-spark\other_electronics 76a4c7ff984b:home

D:\big.data\docker-hadoop\spark\docker-spark> docker exec -it spark-master /bin/bash
bash-5.0# cd home
bash-5.0# ls
laptops  mobiles  other_electronics  tablet  tvs
bash-5.0#
```

Dataset directory : tablets

```
bash-5.0# /spark/bin/spark-shell --master=local spark://spark-master:7077
```

Welcome to



Using Scala version 2.12.15 (OpenJDK 64-Bit Server VM, Java 1.8.0_275)
Type in expressions to have them evaluated.
Type :help for more information.

```
scala> val laptopsData = spark.read.format("csv").option("header", "true").load("/home/laptops")
laptopsData: org.apache.spark.sql.DataFrame = [u_id: string, name: string ... 9 more fields]

scala> val tvsData = spark.read.format("csv").option("header", "true").load("/home/tvs")
tvsData: org.apache.spark.sql.DataFrame = [u_id: string, name: string ... 9 more fields]
```

```
scala> val tabletData = spark.read.format("csv").option("header", "true").load("/home/tablet")
tabletData: org.apache.spark.sql.DataFrame = [u_id: string, name: string ... 9 more fields]

scala> val laptopsData = spark.read.format("csv").option("header", "true").load("/home/laptops")
laptopsData: org.apache.spark.sql.DataFrame = [u_id: string, name: string ... 9 more fields]
```

4.5 SPARK COMMANDS & ANALYTICS

```
scala> val data = spark.read.format("csv")
data: org.apache.spark.sql.DataFrameReader = org.apache.spark.sql.DataFrameReader@5b33446d

scala> val data = spark.read.format("csv").option("header", "true").option("inferSchema", "true").load("/home/tablet")
data: org.apache.spark.sql.DataFrame = [u_id: string, name: string ... 9 more fields]

scala> data.printSchema()
root
 |-- u_id: string (nullable = true)
 |-- name: string (nullable = true)
 |-- offer_price: integer (nullable = true)
 |-- original_price: integer (nullable = true)
 |-- off_now: string (nullable = true)
 |-- total_ratings: integer (nullable = true)
 |-- total_reviews: integer (nullable = true)
 |-- rating: double (nullable = true)
 |-- description: string (nullable = true)
 |-- item_link: string (nullable = true)
 |-- created_at: string (nullable = true)
```

The detailed description of the dataset

```
scala> data.show()
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| u_id | name|offer_price|original_price|off_now|total_ratings|total_reviews|rating| description| item_link| cre
+-----+-----+-----+-----+-----+-----+-----+-----+
| SEYTLMU5|Oppo Pad Air 4 GB...| 16999| 29999|43% off| 1908| 238| 4.4|['4 GB RAM | 64 G...|https://www.flipk...|2022-11-03 22
| 33:...| 0WWIN9FG|Oppo Pad Air 4 GB...| 19999| 34999|42% off| 1908| 238| 4.4|['4 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| 33:...| 2BE7C106|Redmi Pad 6 GB RA...| 19999| 33999|41% off| 503| 80| 4.4|['6 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| 33:...| A97WTC95|Redmi Pad 4 GB RA...| 17999| 28999|37% off| 503| 80| 4.4|['4 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| 33:...| KJYUQ9B|Redmi Pad 4 GB RA...| 17999| 28999|37% off| 503| 80| 4.4|['4 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| 33:...| 9F39HUEC|realme Pad X 6 GB...| 27999| 44999|37% off| 3263| 348| 4.4|['Compatible with...|https://www.flipk...|2022-11-03 22
| 33:...| ETTF77J10|realme Pad X 6 GB...| 27999| 44999|37% off| 3263| 348| 4.4|['Compatible with...|https://www.flipk...|2022-11-03 22
| 33:...| VKX7NOZS|SAMSUNG Galaxy Ta...| 24999| 30999|19% off| 10126| 838| 4.5|['4 GB RAM | 64 G...|https://www.flipk...|2022-11-03 22
| 33:...| U5RX4T8A|realme Pad 4 GB R...| 16849| 29999|43% off| 44194| 4830| 4.4|['4 GB RAM | 64 G...|https://www.flipk...|2022-11-03 22
| 33:...| U9FMK61|realme Pad 4 GB R...| 16899| 29999|43% off| 44194| 4830| 4.4|['4 GB RAM | 64 G...|https://www.flipk...|2022-11-03 22
| 33:...| 9N70MORV|SAMSUNG Galaxy Ta...| 13999| 21599|35% off| 3513| 339| 4.3|['3 GB RAM | 32 G...|https://www.flipk...|2022-11-03 22
| 33:...| V72CGWE3|realme Pad X 4 GB...| 19999| 29999|33% off| 3263| 348| 4.4|['Compatible with...|https://www.flipk...|2022-11-03 22
| 33:...| J8IWHIYJ|realme Pad Mini 4...| 14999| 23999|37% off| 4970| 484| 4.3|['18W Fast Chargi...|https://www.flipk...|2022-11-03 22
```

To show the number of rows in the dataset file

```
scala> val numRows = data.count()
numRows: Long = 7464

scala> println(s"Number of rows: $numRows")
Number of rows: 7464
```

Average ratings received for the tablets

```

scala> val avgRating = data.select(avg("rating")).first().getDouble(0)
avgRating: Double = 3.7867229367631134

scala> println(s"Average rating: $avgRating")
Average rating: 3.7867229367631134

```

High rated items: rating > 4.5

```

scala> val highRatedItems = data.filter($"rating" > 4.5)
highRatedItems: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [u_id: string, name: string ... 9 more fields]

scala> highRatedItems.show()
+-----+-----+-----+-----+-----+-----+-----+-----+
| u_id | name | offer_price | original_price | off_now | total_ratings | total_reviews | rating |
+-----+-----+-----+-----+-----+-----+-----+-----+
| UNKIHQQ | APPLE iPad (9th G...) | 30900 | 30900 | 0 | 10196 | 939 | 4.6 | ['64 GB ROM', '25...|https://www.flipk...|2022-11-03 22
| RVCVFUPK | SAMSUNG Galaxy Ta... | 51999 | 76999 | 32% off | 3220 | 344 | 4.6 | ['6 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| TA2F3Z1I | SAMSUNG Galaxy Ta... | 66999 | 85999 | 22% off | 160 | 10 | 4.7 | ['8 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| V0M008MQ | APPLE iPad (9th G...) | 44900 | 44900 | 0 | 10196 | 939 | 4.6 | ['256 GB ROM', '2...|https://www.flipk...|2022-11-03 22
| KSN551H | APPLE iPad Pro (2... | 70900 | 71900 | 1% off | 1961 | 167 | 4.7 | ['Apple M1 chip f...|https://www.flipk...|2022-11-03 22
| 8VCZR1SQ | SAMSUNG Galaxy Ta... | 66999 | 85999 | 22% off | 160 | 10 | 4.7 | ['8 GB RAM | 128 ...|https://www.flipk...|2022-11-03 22
| PSEWFZ2X | APPLE iPad Pro (2... | 79900 | 80900 | 1% off | 1961 | 167 | 4.7 | ['Apple M1 chip f...|https://www.flipk...|2022-11-03 22
| 9KH0YDSA | APPLE iPad mini (... | 60900 | 60900 | 0 | 409 | 49 | 4.6 | ['256 GB ROM', '2...|https://www.flipk...|2022-11-03 22
| 4TCZI9TN | APPLE iPad (9th G...) | 30900 | 30900 | 0 | 10196 | 939 | 4.6 | ['64 GB ROM', '25...|https://www.flipk...|2022-11-03 22
| 128P0388 | APPLE iPad mini (... | 60900 | 60900 | 0 | 409 | 49 | 4.6 | ['64 GB ROM', '21...|https://www.flipk...|2022-11-03 22

```

Compute the average rating and total reviews for each tablet:

```

scala> tabletAnalysis.show
+-----+-----+-----+
| name | avg_rating | total_reviews |
+-----+-----+-----+
| acer Spin 3 Core ... | 5.0 | 2.0 |
| acer Spin 5 Core ... | 5.0 | 2.0 |
| ASUS ROG Flow X13... | 5.0 | 0.0 |
| HP Spectre x360 C... | 5.0 | 7.0 |
| Lenovo Yoga 6 Ryz... | 4.971428571428571 | 5.0 |
| APPLE iPad Pro (2... | 4.8 | 352.0 |
| APPLE iPad Pro (2... | 4.8 | 352.0 |
| APPLE iPad Pro (2... | 4.8 | 308.0 |
| APPLE iPad Pro (2... | 4.8 | 308.0 |
| APPLE iPad Pro (2... | 4.8 | 396.0 |
| APPLE iPad Pro (2... | 4.8 | 308.0 |
| APPLE iPad Pro (2... | 4.8 | 352.0 |
| Lenovo Yoga Slim ... | 4.8 | 48.0 |
| APPLE iPad Pro (2... | 4.8 | 352.0 |
| APPLE iPad Pro (2... | 4.8 | 308.0 |
| HP Envy 13 x360 I... | 4.8 | 16.0 |
| Lenovo Celeron Du... | 4.8 | 0.0 |
| APPLE iPad Pro (2... | 4.8 | 308.0 |
| APPLE iPad Pro (2... | 4.8 | 352.0 |
| APPLE iPad Pro (2... | 4.8 | 396.0 |
+-----+-----+-----+
only showing top 20 rows

```

Analyze the distribution of tablet ratings:

```
scala> ratingDistribution.show
+----+---+
|rating|count|
+----+---+
|    0|  735|
|  1.5|     4|
|  2.2|     4|
|  2.4|     7|
|  2.5|    21|
|  2.6|    22|
|  2.7|   104|
|  2.8|    88|
|  2.9|    31|
|    3|   131|
|  3.1|   178|
|  3.2|    92|
|  3.3|    93|
|  3.4|    78|
|  3.5|   115|
|  3.6|   209|
|  3.7|   214|
|  3.8|   240|
|  3.9|   131|
|    4|   225|
+----+---+
only showing top 20 rows
```

Lowest and highest rated product

```
scala> val highestRatedTablet = tabletsByRating.first()
highestRatedTablet: org.apache.spark.sql.Row = [HP Spectre x360 Core i5 10th Gen - (8 GB/512 GB SSD/Windows 10 Pro) 13-aw0211TU 2 in 1 Laptop,5.0]
scala> val lowestRatedTablet = tabletsByRating.orderBy("avg_rating").first()
lowestRatedTablet: org.apache.spark.sql.Row = [Domo Slate SLP9 OS11 4 GB RAM 64 GB ROM 10.1 inch with Wi-Fi+4G Tablet (Grey),0.0]
```

Analyze the distribution of ratings across different price ranges:

```
scala> val priceRangeAnalysis = tabletData.withColumn("price_range", when(col("offer_price") <= 10000, "Low").when(col("offer_price") <= 20000, "Medium").otherwise("High")).groupBy("price_range").agg(avg("rating").alias("avg_rating"), count("u_id").alias("count"))
priceRangeAnalysis: org.apache.spark.sql.DataFrame = [price_range: string, avg_rating: double ... 1 more field]

scala> priceRangeAnalysis.show
+-----+-----+-----+
|price_range|      avg_rating|count|
+-----+-----+-----+
|    High| 3.950926132559706| 4481|
|     Low| 3.2346350832266335| 1562|
| Medium| 3.8757916959887404| 1421|
+-----+-----+-----+
```

Determine the most popular brands based on the total number of ratings and reviews:

| brand | total_ratings | total_reviews |
|-----------|---------------|---------------|
| APPLE | 1.0337189E7 | 1032695.0 |
| realme | 2486790.0 | 268583.0 |
| SAMSUNG | 2132809.0 | 200068.0 |
| Lenovo | 1941314.0 | 245920.0 |
| Honor | 541250.0 | 67698.0 |
| Alcatel | 528099.0 | 94703.0 |
| I | 489911.0 | 62894.0 |
| Swipe | 165660.0 | 35173.0 |
| Huawei | 151229.0 | 20857.0 |
| HP | 138174.0 | 16051.0 |
| iball | 111902.0 | 15651.0 |
| MOTOROLA | 64304.0 | 8161.0 |
| Micromax | 58243.0 | 8093.0 |
| acer | 51631.0 | 5774.0 |
| Apple | 50967.0 | 7693.0 |
| Nokia | 44925.0 | 5829.0 |
| Panasonic | 33410.0 | 4466.0 |
| TCL | 29940.0 | 5547.0 |
| Oppo | 26596.0 | 3322.0 |
| Redmi | 22376.0 | 3596.0 |

Analyze the distribution of tablets across different brands:

| brand | count |
|------------|-------|
| APPLE | 2460 |
| Lenovo | 738 |
| SAMSUNG | 735 |
| HP | 483 |
| I | 407 |
| ASUS | 296 |
| Alcatel | 243 |
| Swipe | 232 |
| DELL | 190 |
| iball | 170 |
| realme | 155 |
| Maplin | 110 |
| DOMO | 108 |
| Wishitel | 106 |
| Honor | 85 |
| Redmi | 76 |
| Apple | 76 |
| LAVA | 67 |
| Smartbeats | 55 |
| IQOO | 54 |

Determine the top 5 tablets with the highest discount percentage:

| u_id | name | offer_price | original_price | off_now | total_ratings | total_reviews | rating | description | item_link | cre_ated_at |
|----------|-----------------------|-------------|----------------|---------|---------------|---------------|--------|--|---|-------------|
| 3SNS9Z4N | HP Chromebook Med... | 22499 | 24845 | 9% off | 5812 | 668 | 3.8 | ['MediaTek MediaT... https://www.flipk... 2022-11-10 22:58:... | '512 GB ROM' '26.67 cm (10.5 ... | |
| ZCSF014B | APPLE iPad Pro 51... | 72990 | 80500 | 9% off | 3215 | 392 | 4.7 | "[10.5 Retina di... https://www.flipk... 2022-11-10 22:0.5 ... | '512 GB ROM' '26.67 cm (10.5 ... | |
| KPQJ2QSJ | APPLE iPad Pro 51... | 72990 | 80500 | 9% off | 3215 | 392 | 4.7 | "[10.5 Retina di... https://www.flipk... 2022-11-10 22:0.5 ... | '512 GB ROM' '26.67 cm (10.5 ... | |
| HWURVWNG | MOTOROLA e32s (Mi...) | 7990 | 8799 | 9% off | 95 | 18 | 3.7 | [3 GB RAM 32 G... https://www.flipk... 2022-11-10 22:57:... | '3 GB RAM 32 G... https://www.flipk... 2022-11-10 22:59:... | |
| J9HUMMUX | APPLE iPad Pro 20... | 89900 | 98900 | 9% off | 819 | 84 | 4.7 | [6 GB RAM 512 ... https://www.flipk... 2022-11-10 22:59:... | '6 GB RAM 512 ... https://www.flipk... 2022-11-10 22:59:... | |

Maximum discount by each brand and its rating

| brand | max_discount | rating |
|----------|--------------|--------|
| APPLE | 9% off | 4.6 |
| Swipe | 9% off | 3.7 |
| Lenovo | 9% off | 4.3 |
| ankikrit | 9% off | 3.8 |
| iball | 9% off | 4 |
| HP | 9% off | 3.8 |
| MOTOROLA | 9% off | 4.1 |
| SAMSUNG | 9% off | 4.5 |
| Wishtel | 8% off | 3.8 |
| DELL | 8% off | 4.2 |
| ASUS | 8% off | 0 |
| Datawind | 8% off | 3.1 |
| Coolpad | 7% off | 3.4 |
| Elevn | 62% off | 3.7 |
| Contixo | 62% off | 2.7 |
| DOMO | 61% off | 3.7 |
| TCL | 60% off | 3.8 |
| FUSION5 | 60% off | 3.3 |
| VIZIO | 51% off | 0 |
| nuveck | 51% off | 2.4 |

AvgOriginalPriceByOfferPrice

| | |
|-------------|--------------------|
| offer_price | avg_original_price |
| 81900 | 84098.0 |
| 59990 | 66964.0 |
| 40900 | 40900.0 |
| 47900 | 57900.0 |
| 22990 | 28132.85714285714 |
| 2999 | 3306.6923076923076 |
| 162900 | 162900.0 |
| 91990 | 112990.0 |
| 85900 | 90390.625 |
| 80999 | 99999.0 |
| 89900 | 92900.0 |
| 11146 | 14999.0 |
| 16999 | 23974.4 |
| 63191 | 71900.0 |
| 82500 | 98200.0 |
| 4799 | 7999.0 |
| 217999 | 229599.0 |
| 38600 | 38600.0 |
| 130000 | 140990.0 |
| 134990 | 152990.0 |

Summary of the datasets

```
scala> summaryData.show()
+-----+-----+
|Category|Count|      AverageRating|
+-----+-----+
| Laptops|16848| 3.214693732193721|
| Mobiles|17760| 4.090196917808223|
|    TVs|16752|3.0908491107286307|
| Tablets| 7464|3.7867229367631134|
+-----+-----+
```

5 APACHE HIVE: BILAWAL

5.1 WHAT IS HIVE

Apache Hive is a data analysis tool that runs on top of Hadoop and serves as a data warehousing solution. It is designed for users familiar with SQL, providing a SQL-like language called HiveQL for managing and querying structured data. By using Hive, we can simplify the complexities associated with Hadoop, making it easier to work with and analyze data.

Our Working:

In our docker-compose.yml file, we have included important components like Hadoop and Hive as Docker images. These components are essential for our project as they allow us to perform queries on the Apache Hive platform. We will use Hadoop to store our large dataset, and then we will load this data into Apache Hive from a hdfs path.

5.2 STEPS.

- Run docker file and check if containers are running

```
bilawalali@tenpearlss-MacBook-Pro docker-hive % docker-compose up
[+] Running 6/0
  ● Container docker-hive-hive-server-1          Created                               0.0s
  ● Container docker-hive-namenode-1            Created                               0.0s
  ● Container docker-hive-presto-coordinator-1   Created                               0.0s
  ● Container docker-hive-datanode-1           Created                               0.0s
  ● Container docker-hive-hive-metastore-postgresql-1   Created                               0.0s
  ● Container docker-hive-hive-metastore-1        Created                               0.0s
Attaching to docker-hive-datanode-1, docker-hive-hive-metastore-1, docker-hive-hive-server-1, docker-hive-namenode-1, docker-hive-presto-coordinator-1
docker-hive-datanode-1      | Configuring core
                           | - Setting hadoop.proxyuser.hue.hosts=*
                           | - Setting fs.defaultFS=hdfs://namenode:8020
                           | - Setting hadoop.proxyuser.hue.groups=**
                           | - Setting hadoop.http.staticuser.user=root
                           | Configuring hdfs
                           | - Setting dfs.namenode.datanode.registration.ip-hostname-check=false
                           | - Setting dfs.datanode.data.dir=file:///hadoop/dfs/data
                           | - Setting dfs.permissions.enabled=false
                           | - Setting dfs.webhdfs.enabled=true
                           | Configuring yarn
                           | - Setting yarn.resourcemanager.fs.state-store.uri=/rmstate
                           | - Setting yarn.timeline-service.generic-application-history.enabled=true
                           | - Setting yarn.resourcemanager.recovery.enabled=true
                           | - Setting yarn.log-aggregation-enable=true
                           | Configuring core
                           | - Setting yarn.timeline-service.enabled=true
                           | Configuring core
                           | - Setting yarn.resourcemanager.store.class=org.apache.hadoop.yarn.server.resourcemanager.recovery.FileSystemRMStateStore
                           | - Setting hadoop.proxyuser.hue.hosts=*
                           | - Setting yarn.resourcemanager.system-metrics-publisher.enabled=true
```

| CONTAINER ID | IMAGE | COMMAND | CREATED | STATUS | PORTS |
|----------------|--|----------------------------------|--------------|-------------------------------|------------------------|
| 809ca313e43e | docker-bigdata-bigdata-cluster-nifi | "./scripts/start.sh_" | 19 hours ago | Exited (137) 5 minutes ago | |
| c9e3da0453f4 | docker-bigdata-bigdata-cluster-elasticsearch | "sh start-elasticsea..." | 19 hours ago | Exited (137) 5 minutes ago | |
| d28829b592a3 | bde2020/hadoop-base:2.0.0-hadoop3.2.1-jav8 | "/entrypoint.sh" | 19 hours ago | Restarting (0) 12 seconds ago | |
| 65d4c1f4244d | bde2020/hadoop-namenode:2.0.0-hadoop2.7.4-jav8 | "/entrypoint.sh /run..." | 19 hours ago | Up 4 minutes (healthy) | 0.0.0.0:50070->50070/ |
| tcp_2976a1143 | docker-hive-namenode-1 | docker-hive-namenode-1 | | | |
| 922c7e9705d7 | bde2020/hadoop-datanode:2.0.0-hadoop2.7.4-jav8 | "/entrypoint.sh /run..." | 19 hours ago | Up 4 minutes (healthy) | 0.0.0.0:50075->50075/ |
| tcp_10002/tcp | docker-hive-datanode-1 | docker-hive-datanode-1 | | | |
| 26ec1a86d055 | bde2020/hive:2.3.2-postgresql-metastore | "/entrypoint.sh /bin/..." | 19 hours ago | Up 4 minutes | 0.0.0.0:10000->10000/ |
| tcp_9983/tcp | bde2020/hive:2.3.2-postgresql-metastore | "/entrypoint.sh /opt/..." | 20 hours ago | Up 4 minutes | 10000/tcp, 0.0.0.0:90 |
| b70925180a20 | bde2020/hive-metastore-postgresql:2.3.0 | "/docker-entrypoint..." | 21 hours ago | Up 4 minutes | 5432/tcp |
| 1222ac4d9695 | shawnzhu/prestodb:0.181 | "/bin/launcher run" | 21 hours ago | Up 4 minutes | 0.0.0.0:8080->8080/tcp |
| p_26ec1a86d055 | docker-hive-presto-coordinator-1 | docker-hive-presto-coordinator-1 | | | |

- Copy data to container and bash into the container to see the data

```
bilawalali@tenpearlss-MBP docker-hive % ls
Dockerfile          README.md          docker-compose.yml    hadoop-hive.env      startup.sh
Makefile           conf                entrypoint.sh     my_data.csv
bilawalali@tenpearlss-MBP docker-hive % docker cp my_data.csv 26ec1a86d055:home
bilawalali@tenpearlss-MBP docker-hive % docker compose exec hive-server bash
bilawalali@tenpearlss-MBP docker-hive % docker compose exec hive-server bash
root@26ec1a86d055:/opt#
root@26ec1a86d055:/opt#
root@26ec1a86d055:/opt# cd ..
root@26ec1a86d055:/#
bin  dev   etc       home  lib64  mnt  proc  run   srv  tmp  var
boot entrypoint.sh  hadoop-data  lib   media  opt  root  sbin  sys  usr
root@26ec1a86d055:/# cd home
root@26ec1a86d055:/home#
root@26ec1a86d055:/home# ls
EcomData  my_data.csv
root@26ec1a86d055:/home#
```

3. Renaming the data so that we can read it properly

```
root@26ec1a86d055:/home# k
bash: k: command not found
root@26ec1a86d055:/home# ls
EcomData  my_data.csv
root@26ec1a86d055:/home# mv my_data.csv EcommerceData
root@26ec1a86d055:/home#
root@26ec1a86d055:/home# ls
EcomData  EcommerceData
root@26ec1a86d055:/home#
```

4. Saving Data to Hadoop

```
root@26ec1a86d055:/# hadoop fs -put -f /home /user/dataset
root@26ec1a86d055:/#
root@26ec1a86d055:/# hadoop fs -ls /user/dataset
Found 1 items
drwxr-xr-x  - root supergroup          0 2023-06-07 15:55 /user/dataset/home
root@26ec1a86d055:/#
```

5. Create the table and copy the values into it from our Data

```
root@026ec1a86d055:/# docker-hive --docker-compose - docker compose exec hive-server bash - 127x30
root@026ec1a86d055:/# /opt/hive/bin/beeline -u jdbc:hive2://localhost:10000
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/hive/lib/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/hadoop-2.7.4/share/hadoop/common/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLog
erBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Connecting to jdbc:hive2://localhost:10000
Connected to: Apache Hive (version 2.3.2)
Driver: Hive JDBC (version 2.3.2)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 2.3.2 by Apache Hive
0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000> CRETAE TABLE EcommerceDataTable (u_id String, name String, offer_price String, original_price
String, off_now String, total_ratings String, total_review String, rating String, description String) row format delimited file
ds terminated by ',';
Error: Error while compiling statement: FAILED: ParseException line 1:0 cannot recognize input near 'CRETAE' 'TABLE' 'Ecommerce
DataTable' (state=42000,code=40000)
0: jdbc:hive2://localhost:10000> CREATE TABLE EcommerceDataTable (u_id String, name String, offer_price String, original_price
String, off_now String, total_ratings String, total_review String, rating String, description String) row format delimited file
ds terminated by ',';
No rows affected (3.246 seconds)
0: jdbc:hive2://localhost:10000>
```

```
root@026ec1a86d055:/# docker-hive --docker-compose - docker compose exec hive-server bash - 138x30
0: jdbc:hive2://localhost:10000> LOAD DATA INPATH '/user/dataset/home/EcommerceData' OVERWRITE INTO TABLE EcommerceDataTable;
No rows affected (1.48 seconds)
0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000> select * from EcommerceDataTable limit 10;
+-----+-----+-----+-----+-----+
| ecommercedatatable.u_id | ecommercedatatable.name | ecommercedatatable.offer_price | ecommercedatatable ori
ginal_price | ecommercedatatable.off_now | ecommercedatatable.total_ratings | ecommercedatatable.total_review | ecommercedatatable.rat
ing | ecommercedatatable.description |
+-----+-----+-----+-----+-----+
| u_id | name | offer_price | original_price | rating |
| off_now | total_reviews | total_ratings | |
| description | |
+-----+-----+-----+-----+
| BRULJIVR | Adsun 80 cm (32 inch) HD Ready LED Smart Android Based TV | 7849 | 21999 | 3.8
| 64% off | 11038 | 1587 | |
| "[Operating System: Android Based'" | |
| HOHZ3IE | Adsun 80 cm (32 inch) HD Ready LED TV | 6653 | 16999 | 3.8
| 60% off | 11038 | 1587 | |
| "[HD Ready 1366 x 768 Pixels'" | |
| 1SSBTPAH | Mi 5A 80 cm (32 inch) HD Ready LED Smart Android TV with Dolby Audio (2022 Model) | 12499 | 28889 |
| 24999 | 303019 | 50% off | |
| 4.4 | "[Operating System: Android'" | |
| 16405M10 | realme 80 cm (32 inch) HD Ready LED Smart Android TV | 15499 | 17999 | 4.3
| 13% off | 209810 | |
| "[Operating System: Android'" | |
| UC5TBKEA | LG 80 cm (32 inch) HD Ready LED Smart WebOS TV | 12980 | 21990 | 4.4
| 40% off | 42111 | 3781 | |
```

We have created the table by fetching the data from hadoop, now it's time to run the Hive Queries.

Check if queries are running fine

Query: Get count and data of all records where rating is greater than 4.5

```
docker-hive - docker-compose - docker compose exec hive-server bash - 138x30
0: jdbc:hive2://localhost:10000> select count(*) from EcommerceDataTable;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
+-----+
| _c0 |
+-----+
| 937 |
+-----+
1 row selected (1.667 seconds)
0: jdbc:hive2://localhost:10000> select count(*) from EcommerceDataTable where ecommercedatatable.rating > 4.5;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
+-----+
| _c0 |
+-----+
| 59 |
+-----+
1 row selected (1.517 seconds)
0: jdbc:hive2://localhost:10000>
```

```
docker-hive - docker-compose - docker compose exec hive-server bash - 138x30
+-----+
| 30 | Blaupunkt 189 cm (75 inch) Ultra HD (4K) LED Smart Android TV with Dolby Atmos & Dolby Vision | 149999
| 149999 | 0 | 3904 | 1026
| 4.6 | "[ 'Operating System: Android'" | |
| HO9RTMSZ | LG 139 cm (55 inch) Ultra HD (4K) LED Smart WebOS TV | 104990 | 104990
| 0 | 27 | 4 | 4.8
| "[ 'Operating System: WebOS'" | |
| T85JE998 | "Dyanora 80 cm (32 inch) HD Ready LED TV with Noise Reduction | Cinema Zoom | Powerful A
udio Box Speakers" | 6699 | 15999 | 58% off | 795
| 105 | |
| HQSLDN8J | "Dyanora 127 cm (50 inch) Ultra HD (4K) LED Smart Android TV with Noise Reduction | Android 9.0
| Google ..." | 26499 | 43999 | 39% off
| 881 | 133 | |
| SGSAHWZK | Lloyd Clara 163 cm (65 inch) Ultra HD (4K) LED Smart Linux based TV | 74999 | 15399
| 51% off | 6 | 1 | 4.8
| "[ 'Operating System: Linux based'" | |
| 1C8HG07T | SAMSUNG Q60RAK 123 cm (49 inch) QLED Ultra HD (4K) Smart Tizen TV | 93999 | 101900
| 7% off | 44 | 5 | 4.6
| "[ 'Operating System: Tizen'" | |
| 60E1RFQK | LG 164 cm (65 inch) Ultra HD (4K) LED Smart WebOS TV | 102060 | 219990
| 53% off | 3 | 1 | 5
| "[ 'Operating System: WebOS'" | |
| 09CRJAKF | SONY Bravia 108 cm (43 inch) Ultra HD (4K) LED Smart TV | 40999 | 61900
| 33% off | 286 | 54 | 4.6
| "[ 'Ultra HD (4K) 3840 x 2160 pixels Pixels'" | |
+-----+
59 rows selected (0.242 seconds)
0: jdbc:hive2://localhost:10000>
```

Query: get 20 records where offer_price >= 5000 and rating > 4.5 and order by rating

```
docker-hive — docker-compose - docker compose exec hive-server bash — 134x45
0: jdbc:hive2://localhost:10000> select ecommercedatatable.name, ecommercedatatable.rating from EcommerceDataTable where
atable.offer_price >= 5000 and ecommercedatatable.rating > 4.5 order by rating limit 20;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different e
ne (i.e. spark, tez) or using Hive 1.X releases.
+-----+-----+
| ecommercedatatable.name | ecommercedatatable.rating |
+-----+-----+
| Blaupunkt Cybersound 108 cm (43 inch) Ultra HD (4K) LED Smart Android TV with Dolby MS12 & 50W Speaker... | 4.6
|
| SONY BRAVIA 80 cm (32 inch) HD Ready LED Smart Linux based TV | 4.6
|
| SAMSUNG 138 cm (55 inch) QLED Ultra HD (4K) Smart TV | 4.6
|
| Xiaomi OLED Vision 138.8 cm (55 inches) 4K Ultra HD Smart Android TV with Dolby Vision IQ and Dolby At... | 4.6
|
| LG 109.22 cm (43 inch) Ultra HD (4K) LED Smart TV | 4.6
|
| SONY Bravia 108 cm (43 inch) Ultra HD (4K) LED Smart Android TV | 4.6
|
| SONY Bravia 123 cm (49 inch) Ultra HD (4K) LED Smart Android TV | 4.6
|
| SONY Bravia 108 cm (43 inch) Ultra HD (4K) LED Smart TV | 4.6
|
| SAMSUNG Q60RAK 123 cm (49 inch) QLED Ultra HD (4K) Smart Tizen TV | 4.6
|
| Blaupunkt 189 cm (75 inch) Ultra HD (4K) LED Smart Android TV with Dolby Atmos & Dolby Vision | 4.6
|
| LG 139 cm (55 inch) OLED Ultra HD (4K) Smart WebOS TV | 4.6
|
| Blaupunkt Cyber Sound 164 cm (65 inch) Ultra HD (4K) LED Smart Android TV with Dolby MS12 & 60W Speake... | 4.6
|
| SAMSUNG 139 cm (55 inch) QLED Ultra HD (4K) Smart Tizen TV | 4.6
|
| Blaupunkt Cybersound 139 cm (55 inch) Ultra HD (4K) LED Smart Android TV with Dolby MS12 & 60W Speaker... | 4.6
|
| Blaupunkt Cybersound 126 cm (50 inch) Ultra HD (4K) LED Smart Android TV with Dolby MS12 & 60W Speaker... | 4.6
|
| Thomson 139 cm (55 inch) QLED Ultra HD (4K) Smart Google TV With Dolby Vision & Dolby Atmos | 4.7
|
| SONY Bravia 125.7 cm (50 inch) Ultra HD (4K) LED Smart Google TV | 4.7
|
| SONY Bravia 138.8 cm (55 inch) Ultra HD (4K) LED Smart Google TV | 4.7
|
| Blaupunkt 139 cm (55 inch) QLED Ultra HD (4K) Smart Google TV With Dolby Atmos & Far-Field Mic | 4.7
|
| SONY Bravia 189 cm (75 inch) Ultra HD (4K) LED Smart Google TV | 4.7
+
20 rows selected (1.561 seconds)
0: jdbc:hive2://localhost:10000>
```

Query: get 20 records where total reviews are less than 100 and also total ratings is less than 10000.

```
0: jdbc:hive2://localhost:10000> select ecommercedatatable.name, ecommercedatatable.rating from EcommerceDataTable where
atable.total_review < 100 and ecommercedatatable.total_ratings < 10000 limit 20;
+-----+-----+
| ecommercedatatable.name | ecommercedatatable.rating |
+-----+-----+
| Mi X Series 125 cm (50 inch) Ultra HD (4K) LED Smart Android TV with Dolby Vision and Dolby Audio (202... | 4.5
|
| Nokia 109 cm (43 inch) Ultra HD (4K) LED Smart Android TV with Dolby Atmos and Dolby Vision | 4.3
|
| iFFALCON by TCL F53 100 cm (40 inch) Full HD LED Smart Android TV with Android 11 | 4.3
|
| Mi X Series 138 cm (55 inch) Ultra HD (4K) LED Smart Android TV with Dolby Vision and Dolby Audio (202... | 4.5
|
| MarQ by Flipkart 60 cm (24 inch) HD Ready LED TV | 4.2
|
| Nokia 139 cm (55 inch) Ultra HD (4K) LED Smart Android TV with Dolby Atmos and Dolby Vision | 4.3
|
| MOTOROLA Revou 2 80 cm (32 inch) HD Ready LED Smart Android TV with Sound by boAt | 4.2
|
| MOTOROLA Revou 2 109 cm (43 inch) Ultra HD (4K) LED Smart Android TV with Sound by boAt and Dolby Visi... | 4.2
|
| TCL C715 Series 139 cm (55 inch) QLED Ultra HD (4K) Smart Android TV with Handsfree Voice Control & Do... | 4.4
|
| SONY Bravia 80 cm (32 inch) HD Ready LED Smart Google TV with Google TV | 4.7
|
| TOSHIBA C350LP 108 cm (43 inch) Ultra HD (4K) LED Smart Google TV with Dolby Vision Atmos and REGZA En... | 4.3
|
| Mi 5A Pro 80 cm (32 inch) HD Ready LED Smart Android TV with 24W Dolby Audio & 1.5GB RAM (2022 Model) | 4.5
|
| Vu 139 cm (55 inch) QLED Ultra HD (4K) Smart Android TV | 4.3
|
| TCL P615 108 cm (43 inch) Ultra HD (4K) LED Smart TV with Dolby Audio | 4.4
|
| Thomson 126 cm (50 inch) QLED Ultra HD (4K) Smart Google TV With Dolby Vision & Dolby Atmos | 4.7
|
| iFFALCON by TCL 126 cm (50 inch) Ultra HD (4K) LED Smart Google TV | 0
|
| Thomson 139 cm (55 inch) QLED Ultra HD (4K) Smart Google TV With Dolby Vision & Dolby Atmos | 4.7
|
| Blaupunkt 139 cm (55 inch) QLED Ultra HD (4K) Smart Google TV With Dolby Atmos & Far-Field Mic | 4.7
|
| Compaq HUEO G50B 127 cm (50 inch) Ultra HD (4K) LED Smart Android TV | 4.4
|
| iFFALCON by TCL 139 cm (55 inch) Ultra HD (4K) LED Smart Google TV | 0
+
20 rows selected (0.232 seconds)
0: jdbc:hive2://localhost:10000>
```