



University of
Salford
MANCHESTER

ASSESSMENT REPORT

ON

Big Data Tools & Techniques

Submitted by

TOMISIN SAMUEL ISEDOWO

(@00749815)

In fulfillment of

Big Data Tools & Techniques

For

MASTER OF SCIENCE IN ARTIFICIAL INTELLIGENCE

April 2024.

TASK 1

ANALYSIS OF THE CLINICAL TRIAL AND PHARMA DATASET

INTRODUCTION

In today's digital world, we deal with a lot of data. This brings both opportunities and challenges for businesses. The increasing number of devices like smartphones and computers results in a larger amount of data being generated. Smart devices can provide insights into people's behaviors. To manage this growing volume of data, specialized tools and methods have been developed. They help collect, store, analyze, and show data, often in real-time. These tools let businesses find important insights hidden in big datasets. This era marks a transition from an Informational Society to a Knowledge-Based Society, where the strategic utilization of big data plays a pivotal role. (Youssra et al, 2018).

The analysis will use PySpark methods in RDDs and DataFrames, followed by SQL queries. PySpark is a Python library for large-scale data processing that works well with distributed systems like Apache Spark, allowing parallel computation across multiple nodes. It helps process extensive datasets quickly and efficiently. RDDs are foundational in PySpark, representing distributed object collections processed in parallel, with fault tolerance for seamless recovery from failures. DataFrames organize data in tabular form for structured manipulation. SQL, a common language for database operations, allows analysts to query, filter, and aggregate data efficiently. PySpark enhances SQL operations by offering an SQL interface for RDDs and DataFrames.

Data Cleaning and Preparation

The investigation revealed null values in the Clinical Trials and Pharma CSV files. Removing these nulls might distort analysis results by also eliminating non-null values. These files were imported into the Databricks filesystem, and a notebook was generated, considering the implementation approach (SQL or Python). The listed process was carried out during data cleaning and preparation.

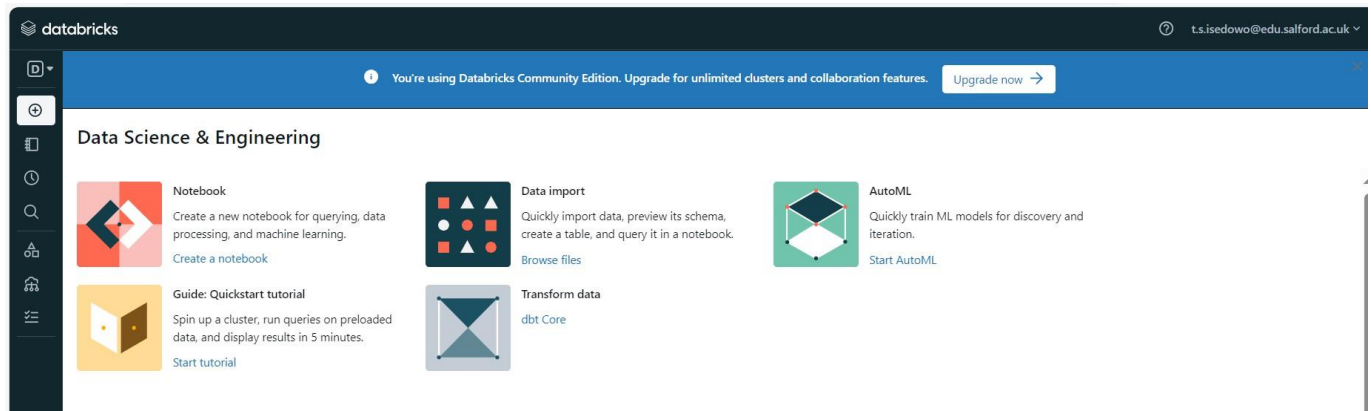


Fig 1: Databricks Environment

The utilization of the 'dbutils.fs.ls' function was pivotal in assessing the presence of the imported CSV files. Upon examination, it was noted that the clinical trials CSV files were successfully uploaded.

```

1 #CHECKING TO SEE THE UPLOADED FILE IN DBFS
2 dbutils.fs.ls('FileStore/tables')

FileInfo(path='dbfs:/FileStore/tables/activations-1.zip', name='activations-1.zip', size=8411369, modificationTime=1706971661000),
FileInfo(path='dbfs:/FileStore/tables/activations-2.zip', name='activations-2.zip', size=8411369, modificationTime=1706972496000),
FileInfo(path='dbfs:/FileStore/tables/activations.csv', name='activations.csv', size=8411369, modificationTime=17069707562000),
FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023.csv', name='clinicaltrial_2023.csv', size=292436366, modificationTime=1712699265000),
FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023.zip', name='clinicaltrial_2023.zip', size=57166668, modificationTime=1709057473000),
FileInfo(path='dbfs:/FileStore/tables/flood.csv', name='flood.csv', size=128984, modificationTime=1707924837000),
FileInfo(path='dbfs:/FileStore/tables/iotstream/', name='iotstream/', size=0, modificationTime=0),
FileInfo(path='dbfs:/FileStore/tables/iotstream.zip', name='iotstream.zip', size=43891, modificationTime=1708611270000),
FileInfo(path='dbfs:/FileStore/tables/logs-1.zip', name='logs-1.zip', size=18168065, modificationTime=1706974312000),
FileInfo(path='dbfs:/FileStore/tables/logs.zip', name='logs.zip', size=18168065, modificationTime=1706712570000),
FileInfo(path='dbfs:/FileStore/tables/movies.csv', name='movies.csv', size=494431, modificationTime=1709731217000),
FileInfo(path='dbfs:/FileStore/tables/myratings.csv', name='myratings.csv', size=10683, modificationTime=1709731677000),
FileInfo(path='dbfs:/FileStore/tables/myratings_1_.csv', name='myratings_1_.csv', size=10683, modificationTime=1709731216000),
FileInfo(path='dbfs:/FileStore/tables/pharma.zip', name='pharma.zip', size=109982, modificationTime=1714039691000),
FileInfo(path='dbfs:/FileStore/tables/ratings.csv', name='ratings.csv', size=2483723, modificationTime=1709731217000),
FileInfo(path='dbfs:/FileStore/tables/steam_200k.csv', name='steam_200k.csv', size=8059447, modificationTime=1712845023000),
FileInfo(path='dbfs:/FileStore/tables/test-1.json', name='test-1.json', size=17958, modificationTime=1706109125000),
FileInfo(path='dbfs:/FileStore/tables/test-2.json', name='test-2.json', size=17958, modificationTime=1706663150000),
FileInfo(path='dbfs:/FileStore/tables/test.json', name='test.json', size=17958, modificationTime=1706109125000),
FileInfo(path='dbfs:/FileStore/tables/webpage/', name='webpage/', size=0, modificationTime=0),
FileInfo(path='dbfs:/FileStore/tables/webpage.zip', name='webpage.zip', size=1582, modificationTime=1707312864000)

```

This line of code sets an environment variable named "fileroot" to the value stored in the variable fileroot. Environment variables are key-value pairs that can be accessed by programs running on a computer. They provide a way to customize the behavior of programs and scripts without modifying their code directly. In this case, the code appears to be assigning a value to an environment variable named "fileroot" using the Python os.environ module. The value assigned to "fileroot" is taken from the variable fileroot, which should be defined elsewhere in the code. This allows other parts of the program or scripts to access the value of "fileroot" using os.environ['fileroot'].

```
Cnd 2
1 #remove clinical trial csv should in case it is available
2 dbutils.fs.rm("/FileStore/tables/clinicaltrial_2023.csv")
3 #declare a variable called fileroot
4 fileroot = "clinicaltrial_2023"
5 # copy the file to the local cluster to unzip the file
6 dbutils.fs.cp("/FileStore/tables/" + fileroot + ".zip", "file:/tmp/")

Out[2]: True
Command took 1.62 seconds -- by t.s.isedow@edu.salford.ac.uk at 5/2/2024, 12:40:39 AM on Task1

Cnd 3
1 import os
2 os.environ['fileroot'] = fileroot

Command took 0.13 seconds -- by t.s.isedow@edu.salford.ac.uk at 5/2/2024, 12:40:47 AM on Task1

Cnd 4
```

Moreso, the existence of the CSV file was checked, and some actions were performed based on its presence.

The existence of the file with the name "\${fileroot}.csv" was confirmed in the "/tmp" directory. "\${fileroot}" is a variable that holds a filename prefix.

If the CSV file exists, the condition [-f "/tmp/\${fileroot}.csv"] evaluates to true, the script proceeds to remove it using the rm command.

Then, the script unzips a file located at "/tmp/\${fileroot}.zip" and extracts its contents into the "/tmp" directory using the unzip command with the -d flag specifying the target directory.

```
Cnd 4
1 $sh
2 # Check if the CSV file exists
3 if [ -f "/tmp/${fileroot}.csv" ]; then
4     # If the CSV file exists, remove it
5     rm "/tmp/${fileroot}.csv"
6 fi
7
8 # Unzip the ZIP file
9 unzip -d "/tmp" "/tmp/${fileroot}.zip"

Archive: /tmp/clinicaltrial_2023.zip
inflating: /tmp/clinicaltrial_2023.csv
Command took 2.44 seconds -- by t.s.isedow@edu.salford.ac.uk at 5/2/2024, 12:41:04 AM on Task1

Cnd 5
1 dbutils.fs.cp("file:/tmp/" + fileroot + ".csv", "/FileStore/tables")

Out[5]: True
Command took 20.35 seconds -- by t.s.isedow@edu.salford.ac.uk at 5/2/2024, 12:41:12 AM on Task1

Cnd 6
```

myRDD1 is a Resilient Distributed Dataset (RDD) created by reading a CSV file located at "FileStore/tables/clinicaltrial_2023.csv" using the sc.textFile() function. Each line in the CSV file becomes an element in myRDD1. It undergoes various Spark operations like map(), filter(), and reduce() to preprocess rows, ensuring consistent formatting and addressing discrepancies in column lengths. Then, a map transformation is applied again to create myRDD2, retaining only the original elements without the indexThe take(5) action is then called on myRDD2 to retrieve the first 5 elements of the RDD,

```
1 #THIS COMMAND HELP TO READ THE CSV FILE AS AN RDD
2 myRDD1= sc.textFile("FileStore/tables/clinicaltrial_2023.csv")

Command took 0.91 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:11:25 AM on sunday

Cnd 8

1 #THIS COMMAND ASSISTS IN DEFINING A FUNCTION DEDICATED TO CLEANING AND TRANSFORMING INDIVIDUAL ROWS WITHIN THE CLINICALTRIAL DATASET
2 def clean_and_transform(row):
3     #THIS AIDS IN ELIMINATING SURPLUS TRAILING COMMAS AND UNNECESSARY DOUBLE QUOTES THROUGH THE UTILIZATION OF THE STRIP METHOD.
4     cleaned_row = row.strip(",").strip('"')
5     #THIS COMMAND FACILITATES THE SPLITTING OF THE ROW BY THE SPECIFIED DELIMITER, WHICH IN THIS CASE IS THE TAB ('\t').
6     cleaned_row = cleaned_row.split('\t')
7
8     return cleaned_row
9
10 # RDD containing rows with varying column lengths
11
12 def fill_empty_columns(row):
13     # Check if the row has less than 14 columns
14     if len(row) < 14:
15         # Fill up the remaining columns with empty strings
16         row += [''] * (14 - len(row))
17     return row
18
19

Command took 0.14 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:11:38 AM on sunday

Cnd 9

1 #THIS COMMAND MAP THE CLEAN_AND_TRANSFORM FUNCTION TO myRDD2.
2 myRDD2 = myRDD1.map(clean_and_transform) \
3     .map(fill_empty_columns) \
4     .zipWithIndex().filter(lambda x: x[1] > 0).map(lambda x: x[0])
5
6 myRDD2.take(5)
```

The function `replace_empty_with_null` takes a list as input. It uses a list comprehension to iterate over each field in the input list. If a field is an empty string (""), it replaces it with the Python keyword "None" to indicate a null value. Otherwise, the field remains unchanged. This function is then applied to each element of an RDD called "myRDD2" using the "map" transformation. Two other functions are also defined. The `clean_date(date)` function appends "-01" to short date strings to complete them, returning the modified or original date. The `fill_empty_date(date, default_value)` function returns the default value if the date string is empty; otherwise, it returns the original date string.

```
1 # Function to replace empty strings with None
2 def replace_empty_with_null(line):
3     return [None if field == "" else field for field in line]
4
5 # Replace empty fields with null values
6 replace_with_null = myRDD2.map(replace_empty_with_null)

Command took 0.14 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:17:38 AM on sunday

Cnd 11

1 def clean_date(date):
2     if len(date) < 8:
3         date = date + "-01"
4     return date
5
6 def fill_empty_date(date, default_value):
7     return default_value if not date else date
8
9 # Select the column indices
10 column_indices = [10, 11]
11
12 # Clean up the start and completion dates
13 default_date = "2024-01-01"
14
15 # First RDD: Fill up empty date columns
16 myRDD3 = replace_with_null.map(lambda row: [
17     fill_empty_date(row[i], default_date) if i in column_indices else row[i] for i in range(len(row))
18 ])
19
20 # Second RDD: Format the date columns
21 myRDD4 = myRDD3.map(lambda row: [
22     clean_date(row[i]) if i in column_indices else row[i] for i in range(len(row))
23 ])
24
25 myRDD4.take(5)
```

▶ (1) Spark Jobs

STATEMENT OF PROBLEM ONE

GETTING THE NUMBER OF STUDIES IN THE CLINICAL TRIAL DATASET

ASSUMPTION MADE

It was assumed that the "Id" serves as the unique identifier in the clinical trial dataset, it is being utilized to ascertain the number of studies contained within the dataset. This approach suggests prioritizing the Id over the study title for its distinctiveness and reliability in accurately representing each study entry.

IMPLEMENTATION IN RESILIENT DISTRIBUTED SYSTEM (RDD)

The implementation counts the number of studies in myRDD4, an RDD with study data. It extracts unique study IDs, removes duplicates, and counts them using the `distinct().count()` function. This is a typical operation in big data processing for distributed datasets. In this case, it reveals that there are 483,422 distinct studies in the clinical trial dataset of 2023.



```
Command took 1.48 seconds -- by t.s.isidow@edu.salford.ac.uk at 5/2/2024, 12:42:10 AM on Task1

Cell 12

1 #Q1 IMPLEMENTATION IN RDD
2
3 # Extract Id
4 Id_rdd = myRDD4.map(lambda x: x[0])
5
6 # Count the total number of Id
7
8 distinct_Id = Id_rdd.distinct().count()
9
10 print("Number of Studies: ", distinct_Id)

▶ (1) Spark Jobs
Number of Studies: 483422

Command took 13.93 seconds -- by t.s.isidow@edu.salford.ac.uk at 5/2/2024, 12:42:12 AM on Task1

Cell 13
```

IMPLEMENTATION IN DATAFRAME

The DataFrame was made with PySpark, a Python tool for Apache Spark, a distributed data processing system. We imported modules like `pyspark.sql.types` for data type definitions. We converted an RDD (myRDD4) to a DataFrame (Clinicaltrial_df) using a preset schema (mySchema). In Spark, it's typical to use DataFrames for organized and efficient distributed data handling.

```
1 #Using DataFrame
2 from pyspark.sql.types import*
3
4 mySchema = StructType([
5     StructField("Id", StringType()),
6     StructField("Study Title", StringType()),
7     StructField("Acronym", StringType()),
8     StructField("Status", StringType()),
9     StructField("Conditions", StringType()),
10    StructField("Interventions", StringType()),
11    StructField("Sponsor", StringType()),
12    StructField("Collaborators", StringType()),
13    StructField("Enrollment", StringType()),
14    StructField("Funder Type", StringType()),
15    StructField("Type", StringType()),
16    StructField("Study Design", StringType()),
17    StructField("Start", StringType()),
18    StructField("Completion", StringType()),
19 ])
20
21
22 Clinicaltrial_df =spark.createDataFrame(myRDD4, mySchema)
23
24 Clinicaltrial_df.show(5)
```

(1) Spark Jobs

Clinicaltrial_df: pyspark.sql.dataframe.DataFrame = [id:string, Study Title:string ... 12 more fields]

Id	Study Title	Acronym	Status	Conditions	Interventions	Sponsor	Collaborators	Enrollment	Funder Type	Type	Study Design	Start	Completion
NCT03630472	Effectiveness of ...	PRIDE	COMPLETED	Mental Health Iss...	BEHAVIORAL: PRIDE...	Sangath	Harvard Medical S...	250.0	OTHER	INTERVENTIONAL	Allocation: RANDO...	2018-08-20	2019-02-28
NCT05992571	Oral ketone Pome...	null	RECRUITING	Cerebrovascular F...	OTHER: Placebo/DI...	McMaster University	Alzheimer's Socie...	30.0	OTHER	INTERVENTIONAL	Allocation: RANDO...	2023-10-25	2024-08-01
NCT08237471	Impact of Tight G...	null	TERMINATED	Myocardial Infarc...	DRUG: Insulin (ti...	Melbourne Health	National Health a...	40.0	OTHER	INTERVENTIONAL	Allocation: RANDO...	2005-10-01	2006-05-01
NCT03020272	New Prognostic Pr...	SUPERHELD	RECRUITING	Decompensated Clr...	OTHER: SuperHELD	Assistance Public...	null	500.0	OTHER	INTERVENTIONAL	Allocation: NA[In...	2020-10-01	2023-10-01
NCT06220171	Intake Care: Deve...	IntakeCare	NOT_YET_RECRUITING	Hypertension	Trea...	adherence ...	Instituto Auxologi...	200.0	OTHER	INTERVENTIONAL	Allocation: RANDO...	2024-10-01	2026-04-01

only showing top 5 rows

Command took 4.93 seconds -- by t.s.isoodow@du.salford.ac.uk at 4/28/2024, 8:49:30 AM on sunday

Then "Id" column from the DataFrame, removes duplicate Ids, and then counts the number of unique Ids, providing a count of distinct studies present in the DataFrame.

```
1 collect_id_rdd = Clinicaltrial_df.select("Id")
2
3 # Count the number of distinct study titles
4 distinct_study_count = collect_id_rdd.distinct().count()
5
6 # Print result
7 print("Number of Studies: ", distinct_study_count)
```

(3) Spark Jobs

collect_id_rdd: pyspark.sql.dataframe.DataFrame = [id:string]

Number of Studies: 483422

Command took 20.34 seconds -- by t.s.isoodow@du.salford.ac.uk at 5/2/2024, 12:36:30 AM on Task1

IMPLEMENTATION IN SQL

The code converts RDD myRDD4 into a DataFrame Clinicaltrial_df using schema mySchema. Using DataFrame is common in Spark for structured and optimized data handling. A temporary view called clinicaltrial_2023 is created from Clinicaltrial_df for Spark SQL operations. The SQL query retrieves the count of distinct values in the "Id" column from the clinicaltrial_2023 table, assigning the result the alias "distinct_studies_count".


```
1 #python
2 from pyspark.sql.types import*
3
4 mySchema = StructType([
5     StructField("id", StringType()),
6     StructField("study_title", StringType()),
7     StructField("acronym", StringType()),
8     StructField("status", StringType()),
9     StructField("conditions", StringType()),
10    StructField("interventions", StringType()),
11    StructField("sponsor", StringType()),
12    StructField("collaborators", StringType()),
13    StructField("enrollment", StringType()),
14    StructField("funder_type", StringType()),
15    StructField("type", StringType()),
16    StructField("study_design", StringType()),
17    StructField("start", StringType()),
18    StructField("completion", StringType()),
19 ])
20 Data = spark.createDataFrame(myRDD4, mySchema)
21
22 Data.createOrReplaceTempView('clinicaltrial_2023')
23
24 # CREATE OR REPLACE TEMPORARY VIEW distinct_study_counts
25
26 Data: pyspark.sql.DataFrame <[(id string, Study Title string ... 12 more fields)]>
27 Command took 7.23 seconds ... by f.s.lindor@edu.salford.ac.uk at 1/2/2024, 12:00:49 AM on Task3

28
29 --Question 1
30 -- This SQL query counts the distinct values of the id column in the clinical_trial table
31 -- and aliases the result as distinct_studies_count.
32
33 SELECT COUNT(DISTINCT id) AS distinct_studies_count
34 FROM clinicaltrial_2023;
35
36 (3) Spark Jobs
37
38 Table
39
40 distinct_studies_count
41
42 1 483422
43
44 1 row | 47.24 seconds runtime
45 Command took 45.28 seconds ... by f.s.lindor@edu.salford.ac.uk at 1/2/2024, 12:10:38 AM on Task3
46 Refreshed 2 hours ago
```

DISCUSSION OF THE RESULT

In our analysis, we used three methods to extract and count distinct studies from a clinical trial dataset represented in RDD, DataFrame, and SQL formats within the Spark framework. Despite employing different approaches, all yielded consistent results, showcasing Spark's reliability across different abstraction levels, from RDD to SQL queries.

STATEMENT OF PROBLEM TWO

LIST OF ALL TYPES OF CLINICAL TRIAL 2023 DATASET AND THEIR FREQUENCIES FROM THE MOST FREQUENT TO THE LEAST FREQUENT

ASSUMPTION MADE

It was assumed that the RDD (myRDD4) contains rows of data, each with a 'Type' field indicating a certain category. The goal of the operation is to determine the frequency of occurrence of each unique type within the RDD. This will show the number of times the unique type occurs in the dataset.

IMPLEMENTATION IN RDD

The process begins by extracting the 'Type' field from each row of RDD (myRDD4) and creating key-value pairs, where the key represents the 'Type', and the value is 1, indicating one occurrence of each type.

Next, a reduction by key operation (`reduceByKey`) is performed to aggregate the counts for each type. The frequencies of each type are then sorted in descending order based on their occurrence frequency. Finally, the sorted frequencies are collected into a list, providing a sorted list of types with their respective frequencies.

```
1 #Q2 IMPLEMENTATION IN RDD
2 # Extract 'Type' field and create key-value pairs
3 type_counts_rdd = myRDD4.map(lambda row: (row[10], 1))
4
5 # get the frequencies of each type
6 type_frequency = type_counts_rdd.reduceByKey(lambda x, y: x + y)
7
8 # Sorting the result by frequency in descending order
9 sorted_frequencies = type_frequency.sortBy(lambda x: x[1], ascending=False)
10
11 sorted_frequencies.collect()
```

▶ (3) Spark Jobs

```
Out[32]: [('INTERVENTIONAL', 371382),
('OBSERVATIONAL', 110221),
('EXPANDED_ACCESS', 928),
(None, 891)]
```

Command took 12.81 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:27:19 AM on sunday

Cmd 14

IMPLEMENTATION IN DATAFRAME

This snippet helps to count the occurrences of each unique value in the "Type" column of a DataFrame (`Clinicaltrial_df`), sort the counts in descending order, and display the results. This from command `pyspark.sql.functions import col`, helps to import the `col` function from `pyspark.sql.functions` which is used to create a Column object representing a column in a DataFrame. The `groupBy(). Count()` function helps to group the DataFrame by the "Type" column and counts the number of occurrences of each unique value in the "Type" column.

```
1 #Question 2
2 from pyspark.sql.functions import col
3
4 type_counts_df = Clinicaltrial_df.select(col("Type").alias("Type")).groupBy("Type").count()
5
6 # Sorting the result by frequency in descending order
7 sorted_frequencies_df = type_counts_df.orderBy("count", ascending=False)
8
9 sorted_frequencies_df.show()
10
```

▶ (2) Spark Jobs


```
type_counts_df: pyspark.sql.dataframe.DataFrame = [Type: string, count: long]
sorted_frequencies_df: pyspark.sql.dataframe.DataFrame = [Type: string, count: long]
```

Type	count
INTERVENTIONAL	371382
OBSERVATIONAL	110221
EXPANDED_ACCESS	928
null	891

Command took 20.18 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:46:28 AM on sunday

IMPLEMENTATION IN SQL

The SQL query effectively analyzes the table to identify the distribution of studies across various types. It groups studies by type and calculates the frequency of each type, providing insights into the most prevalent study categories within the dataset. The snippet retrieves the count of occurrences of each unique value in the Type column from the clinicaltrial_2023 table and presents the results in descending order of count.



The screenshot shows a SQL query execution interface. The query is as follows:

```
-- Question 2
1  SELECT
2    Type,
3    COUNT(*) AS count
4  FROM
5    clinicaltrial_2023
6  GROUP BY
7    Type
8  ORDER BY
9    count DESC;
```

The results are displayed in a table with 4 rows and 2 columns: Type and count. The data is as follows:

Type	count
INTERVENTIONAL	371382
OBSERVATIONAL	110221
EXPANDED_ACCESS	928
null	891

Additional information shown in the interface includes: (2) Spark Jobs, 4 rows, 23.99 seconds runtime, and a refresh button. The command was executed by t.s.lisdown@edu.salford.ac.uk at 5/2/2024, 12:10:56 AM on Task1.

DISCUSSION OF THE RESULT

The analysis showed the same results across three methods: RDD, DataFrame, and SQL. They all found the top 5 types of studies and how often they appear. Types like INTERVENTIONAL, OBSERVATIONAL, and EXPANDED_ACCESS give us a good idea of how common different kinds of clinical research are. Seeing "None" shows us we might need to check the data's quality. Getting consistent results from different methods makes us more confident that our findings are reliable. This helps make better decisions in healthcare and research.

STATEMENT OF PROBLEM THREE

LIST OF THE TOP 5 CONDITIONS AND THEIR FREQUENCY OF OCCURRENCE

ASSUMPTION MADE

Prior to performing the extraction, splitting, and counting operation on the 'Conditions' column of the RDD, it is assumed that the RDD (`myRDD4`) contains rows where each row represents a record, and the 'Conditions' column in each row is a string consisting of multiple conditions separated by the '|' character. Additionally, it is assumed that the RDD is structured such that the first element of each row uniquely identifies the record, and the fifth element contains the 'Conditions' column. This assumption is necessary for the subsequent operations to effectively extract and process the relevant data.

IMPLEMENTATION IN RDD

This script extracts values from the 'Conditions' column of an RDD, splits them by '|' character, counts each condition's occurrences, and sorts them in descending order. It then retrieves the top 5 most frequent conditions. First, a maps function is used to map each row of myRDD4 to a tuple containing the first and fifth elements of the row, where row[0] holds a unique identifier and row[4] holds the 'Conditions' column. Next, the split string function is applied to each element of the RDD from the previous step. This step flattens the resulting list of lists into a single list. Finally, reduceByKey is used to sum up the values (counts) for each key (condition) in the RDD.

```
1  #Q3 IMPLEMENTATION IN RDD
2  # Step 1: Extract the 'Conditions' column
3
4  def split_string(text):
5      if text and text != '':
6          return text.split('|')
7      return ''
8
9  conditions_rdd = myRDD4.map(lambda row: (row[0],row[4])).flatMap(lambda x: split_string(x[1])).map(lambda x: (x,1)).reduceByKey(lambda a,b: a+b)
10 sorted_rdd = conditions_rdd.sortBy(lambda x: x[1], ascending=False)
11 # filter_cond_rdd = conditions_rdd.filter(lambda x: x == '')
12
13 sorted_rdd.take(5)
```

▶ (3) Spark Jobs

```
Out[33]: [('Healthy', 9731),
('Breast Cancer', 7502),
('Obesity', 6549),
('Stroke', 4071),
('Hypertension', 4020)]
```

Command took 15.81 seconds -- by t.s.isedowa@edu.salford.ac.uk at 4/28/2024, 8:27:37 AM on sunday

Cmd 15

IMPLEMENTATION IN DATAFRAME

In PySpark DataFrame analysis of clinical trial data, we start by importing necessary functions from `pyspark.sql.functions`. Then, we select only the 'Id' and 'Conditions' columns from the DataFrame `Clinicaltrial_df`, filtering out rows with null or empty 'Conditions'. Next, we explode the array into separate rows, duplicating the 'Id' column for each condition. After that, we group the DataFrame by 'Conditions' and count the occurrences of each condition, displaying the top 5 most frequent ones in descending order.

```
Cell 12

1 #Question 3
2 from pyspark.sql.functions import split, explode, col, lit
3
4 # Split the string by '|' delimiter if it's not empty
5 # Then, explode the resulting array into separate rows
6 # Finally, count the occurrences of each element
7 conditions_df = Clinicaltrial_df.select(Clinicaltrial_df['Id'], Clinicaltrial_df['Conditions']) \
8     .filter(col('Conditions').isNotNull() & (col('Conditions') != '')) \
9     .withColumn('Conditions', split(col('Conditions'), '|')) \
10    .withColumn('Conditions', explode(col('Conditions'))) \
11    .groupBy('Conditions').count()
12
13 # Sort the DataFrame by count in descending order
14 sorted_df = conditions_df.orderBy(col('count').desc())
15
16 # Take the top 5 rows from the sorted DataFrame
17 top_5_rows = sorted_df.limit(5)
18
19 # Show the top 5 elements
20 top_5_rows.show()

▶ (2) Spark Jobs
▶ conditions_df: pyspark.sql.dataframe.DataFrame = [Conditions: string, count: long]
▶ sorted_df: pyspark.sql.dataframe.DataFrame = [Conditions: string, count: long]
▶ top_5_rows: pyspark.sql.dataframe.DataFrame = [Conditions: string, count: long]

-----
| Conditions|count|
-----
| Healthy| 9751|
| Breast Cancer| 7582|
| Obesity| 6549|
| Stroke| 4871|
| Hypertension| 4828|
-----

Command took 23.39 seconds -- by s.s.lindow@edu.salford.ac.uk at 5/2/2024, 12:18:02 AM on task1
```

IMPLEMENTATION IN SQL

This SQL query identifies the top 5 most common conditions mentioned in a clinical trials dataset along with their counts. It uses a CTE called `condition_counts` to explode the Conditions column, separating multiple conditions with a delimiter '|'. Then, it counts occurrences of each condition using `COUNT(*)` and groups them by condition, ordering the results by count in descending order and limiting output to the top 5 conditions.

```
1 --Question 3
2 -- It first line explodes the 'Conditions' column by '|' delimiter, to separate each condition and filters out non-empty and non-null values.
3 -- Then, it counts the occurrences of each condition and groups the results by condition.
4 -- Finally, it selects the top 5 conditions based on the count of occurrences, sorting them in descending order.
5 WITH condition_counts AS (
6     SELECT
7         condition,
8         COUNT(*) AS count
9     FROM (
10         SELECT
11             id,
12             explode(split(Conditions, '\\|')) AS condition
13         FROM
14             clinicaltrial_2023
15         WHERE
16             Conditions IS NOT NULL AND Conditions != ''
17     ) exploded_conditions
18     GROUP BY
19         condition
20 )
21 SELECT
22     condition,
23     count
24 FROM
25     condition_counts
26 ORDER BY
27     count DESC
28 LIMIT 5;
```

(2) Spark Jobs

condition	count
Healthy	9721
Breast Cancer	7502
Obesity	6549
Stroke	4071
Hypertension	4020

5 rows | 22.84 seconds runtime

Command took 22.84 seconds -- by s.s.lindow@edu.salford.ac.uk at 5/2/2024, 12:13:29 AM on Task1

Refreshed 2 hours ago

DISCUSSION OF THE RESULT

We've identified the top 5 prevalent health conditions in our clinical trial dataset. This information is crucial for healthcare providers, researchers, and policymakers to allocate resources effectively and design interventions. Using RDD, DataFrame, and SQL methods, we ensured consistent results, validating the accuracy of our analysis. This consistency highlights the robustness of our approach in exploring and summarizing the dataset.

STATEMENT OF PROBLEM FOUR

LIST OF THE 10 MOST COMMON SPONSOR THAT ARE NOT PHARMACEUTICAL COMPANIES

ASSUMPTION MADE

To ensure our analysis is reliable, we trust that both the clinical and pharmaceutical datasets are accurately curated, providing a thorough record of clinical studies and pharmaceutical companies. We also assume the pharmaceutical datasets are free from major omissions or biases. Furthermore, we presume the classification of companies as pharmaceutical or non-pharmaceutical is correct according to industry standards and available information. These assumptions are vital for accurately identifying the top ten sponsors that are not pharmaceutical companies.

IMPLEMENTATION IN RDD

The pharmaceutical data file was imported into the Databricks environment. However, the file was in a compressed format which needed to be unzipped. The process of unzipping the file was carried out to ensure its accessibility for further analysis.

```
1 dbutils.fs.cp("/FileStore/tables/pharma.zip", "file:/tmp/")
2

Out[22]: True

Command took 0.34 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:22:47 AM on sunday

Cmd 16

1 %sh
2
3 unzip -d /tmp/ /tmp/pharma.zip
4
5 ls /tmp

inflating: /tmp/pharma.csv
Rserv
RtmpV6V3Ux
chauffeur-daemon-params
```

Then, the pharmaceutical dataset named "pharma.csv" was processed and extraction of the unique pharmaceutical company names was initiated. The `sc.textFile()` function reads the file line by line and creates an RDD where each line is a string and `map()` was used to split each line by commas, effectively converting each line into a list of fields. The lambda function was used to remove and strip any leading or trailing double quotation marks from each field in the row, whitespace from each field in the row before the company name was extracted. The `distinct()` was used to remove duplicate entries from the RDD.

```
Cmd 17

1 dbutils.fs.cp("file:/tmp/pharma.csv", "/FileStore/tables")

Out[17]: True

Command took 0.38 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/2/2024, 12:44:38 AM on Task1

Cmd 18

1 pharma_rdd = sc.textFile("/FileStore/tables/pharma.csv") \
2   .map(lambda line: line.split(",")) \
3   .map(lambda row: [field.strip('"') for field in row]) \
4   .map(lambda row: [field.strip() for field in row]) \
5   .map(lambda x: x[1])
6
7 pharma_companies_rdd = pharma_rdd.distinct()
8

Command took 0.25 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/2/2024, 12:44:13 AM on Task1

Cmd 19
```

For the implementation, two RDDs are used which is `myRDD4` and `pharma_companies_rdd`. `myRDD4` contains data where the 6th element is the identifier for a clinical trial while the `pharma_companies_rdd`

contains data related to pharmaceutical companies, likely with an identifier field. The data allows for joining clinical trials with pharmaceutical companies based on some relationship. The left outer join was performed between `clinicalTrial_pairs` and `pharma_companies_rdd`. In the left outer join, all keys from the left table (clinical trials) are kept, even if there's no match in the right table (pharma companies). The resulting value for unmatched keys in the right table will be `None`.

```
Cmd 19
Python ▶ ▼ - x
1 #Q4 IMPLEMENTATION IN RDD
2 filtered_col = myRDD4.map(lambda x: x[6])
3 clinicalTrial_pairs = filtered_col.map(lambda x: (x, 1))
4 pharma_pairs = pharma_companies_rdd.map(lambda x: (x, 1))
5
6 join_rdd = clinicalTrial_pairs.leftOuterJoin(pharma_pairs)
7
8 filtered_rdd = join_rdd.filter(lambda x: x[1][1] is None)
9 result = filtered_rdd.map(lambda x: (x[0], 1)).reduceByKey(lambda a,b: a+b)
10
11 result = result.map(lambda x: (x[1], x[0]))
12
13 result = result.sortByKey(False).take(10)
14
15 print(result)

▶ (3) Spark Jobs

[(3410, 'National Cancer Institute (NCI)'), (3335, 'Assiut University'), (3023, 'Cairo University'), (2951, 'Assistance Publique - Hôpitaux de Paris'), (2766, 'Mayo Clinic'), (2702, 'M.D. Anderson Cancer Center'), (2393, 'Novartis Pharmaceuticals'), (2340, 'National Institute of Allergy and Infectious Diseases (NIAID)'), (2263, 'Massachusetts General Hospital'), (2181, 'National Taiwan University Hospital')]

Command took 15.44 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/28/2024, 8:32:43 AM on sunday
Cmd 20
```

IMPLEMENTATION IN DATAFRAME

For this implementation, the clinical trial sponsors that are not found in the pharma company dataset were filtered out. Then counts the occurrence of each sponsor, sorts the result by count, selects the top 10 sponsors and displays the result. The pharma dataset was read into DBFS and the column 'c1' was named `Parent_company` before joining with the clinical trial dataset.


```
1 pharma_data = spark.read.csv("/FileStore/tables/pharma.csv", header=False)
2
3 # Rename the column
4 pharma_df = pharma_data.withColumnRenamed("c1", "Parent_company")
5
6 # Left Anti Join to find entries that exist in Clinical Trial Sponsor but not in Pharma Companies
7 filtered_df = clinicaltrial_df.join(pharma_df, clinicaltrial_df.Sponsor != pharma_df.Parent_company, "left_anti")
8
9 # Count occurrences again after filtering
10 result_df = filtered_df.groupBy("Sponsor").count()
11
12 # Sort the result in descending order of counts
13 sorted_result_df = result_df.orderBy("count", ascending=False)
14
15 # Take the top 10 records
16 top_10_sponsors = sorted_result_df.limit(10)
17
18 top_10_sponsors.show()
19
```

```

(4) Spark Jobs:
+-----+
| Sponsor|count|
+-----+
|National Cancer I...| 3419|
|Assut University| 3393|
|Cairo University| 3823|
|Assistance Public...| 2951|
|Mayo Clinic| 2768|
|H.D. Anderson Can...| 2782|
|Novartis Pharmace...| 2393|
|National Institut...| 2348|
|Massachusetts Gen...| 2263|
|National Taiwan U...| 2261|
+-----+

Command took 22.22 seconds -- by t.s.isedow@edu.salford.ac.uk at 4/28/2024, 8:51:17 AM on sunday
End 14
```

IMPLEMENTATION IN SQL

This SQL code is implementing a query to find the top 10 clinical trial sponsors that exist in the "Clinical Trial Sponsor" dataset but not in the "Pharma Companies" dataset. A temporary view called `filtered_df` was created where it selects the sponsor names (`Sponsor`) from the `clinicaltrial_2023` and counts the occurrences of each sponsor.

The LEFT ANTI JOIN operation was used to find entries that exist in the `clinicaltrial_2023` but not in the `pharma` based on the condition that `c.Sponsor = p.Parent_company`. This operation ensures that only sponsors that are present in `clinicaltrial_2023` but not in `pharma` are considered.

After finding the sponsors that meet the criteria, they are sorted in descending order of counts. This means sponsors with the highest number of occurrences will appear first.

```
1 %python
2 #This should run on Python cell
3
4 # Load data from CSV into a temporary view
5 spark.read.csv("/FileStore/tables/pharma.csv", header=True, inferSchema=True).createOrReplaceTempView("pharma")
```

```

(2) Spark Jobs

Command took 2.76 seconds -- by t.s.isedow@edu.salford.ac.uk at 5/2/2024, 12:13:43 AM on Task1
End 15
```

```

1 1. Split
2 --> get unique clinical IDs
3 -- Step 3: Left join table to find entries that exist in Clinical Trial Sponsor but not in Pharma Company
4 CREATE OR REPLACE TMP_VIEW filtered_of_A1
5 SELECT
6   c.sponsor,
7   COUNT(*) AS count
8 FROM
9   clinicaltrial_2021 c
10 LEFT JOIN pharma
11   pharma_p ON c.sponsor = p.pharma_company
12 ORDER BY
13   c.sponsor,
14
15 -- Step 4: Sort the results in descending order of counts
16 CREATE OR REPLACE TMP_VIEW sorted_result_of_A1
17 SELECT
18   sponsor,
19   count
20 FROM
21   filtered_of
22 ORDER BY
23   count DESC;
24
25 -- Step 5: Take the top 10 records
26 CREATE OR REPLACE TMP_VIEW top_10_sponsors_A1
27 SELECT
28   sponsor,
29   count
30 FROM
31   sorted_result_of
32 LIMIT 10;
33
34 -- Show the top 10 sponsors
35 SELECT * FROM top_10_sponsors;

```

• 21 Spark Jobs

Sponsor	count
National Cancer Institute (NCI)	5410
Amgen University	1618
Cisco University	9223
Assistance Publique - Hôpitaux de Paris	2951
Mayo Clinic	2766
NCI Anderson Cancer Center	2762
Novartis Pharmaceuticals	2391
National Institute of Allergy and Infectious Diseases (NIAID)	2340
Massachusetts General Hospital	2263
National Taiwan University Hospital	2181

10 rows | 24.85 seconds runtime

Completed 10 of 10 records -- Top 10 records loaded successfully at 10:00:00 AM on 10/10/2021

DISCUSSION AND COMPARISON OF THE RESULT

The result across the three implementations were consistent. Based on the results, it can be observed that the National Cancer Institute has topped the Non-Pharmaceutical Sponsors list. The analysis sheds light on the prominent role of non-pharmaceutical sponsors in clinical research and underscores the importance of collaboration and resource allocation to drive innovation and address healthcare challenges effectively.

STATEMENT OF PROBLEM FIVE

VISUALS DISPLAYING THE NUMBER OF COMPLETED STUDIES IN THE YEAR 2023

ASSUMPTION MADE

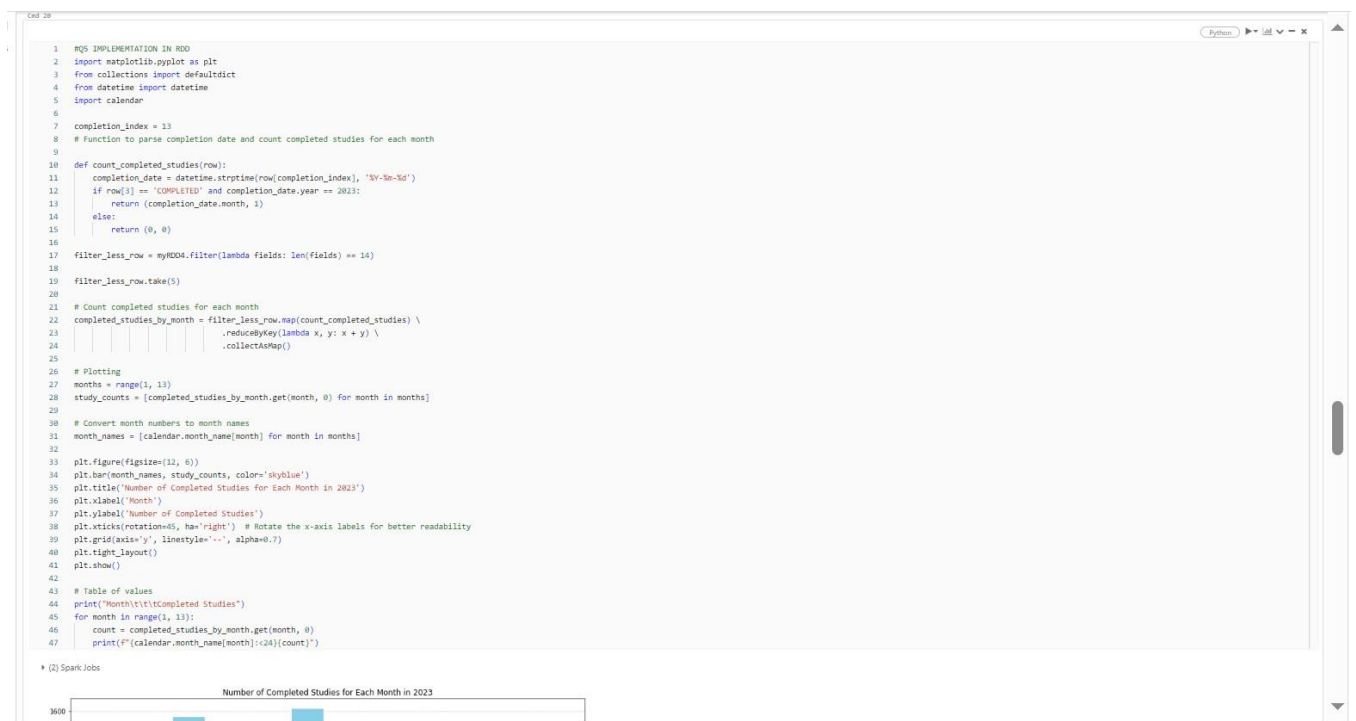
We assume that the completion dates recorded in the dataset are consistent and follow a standard format.

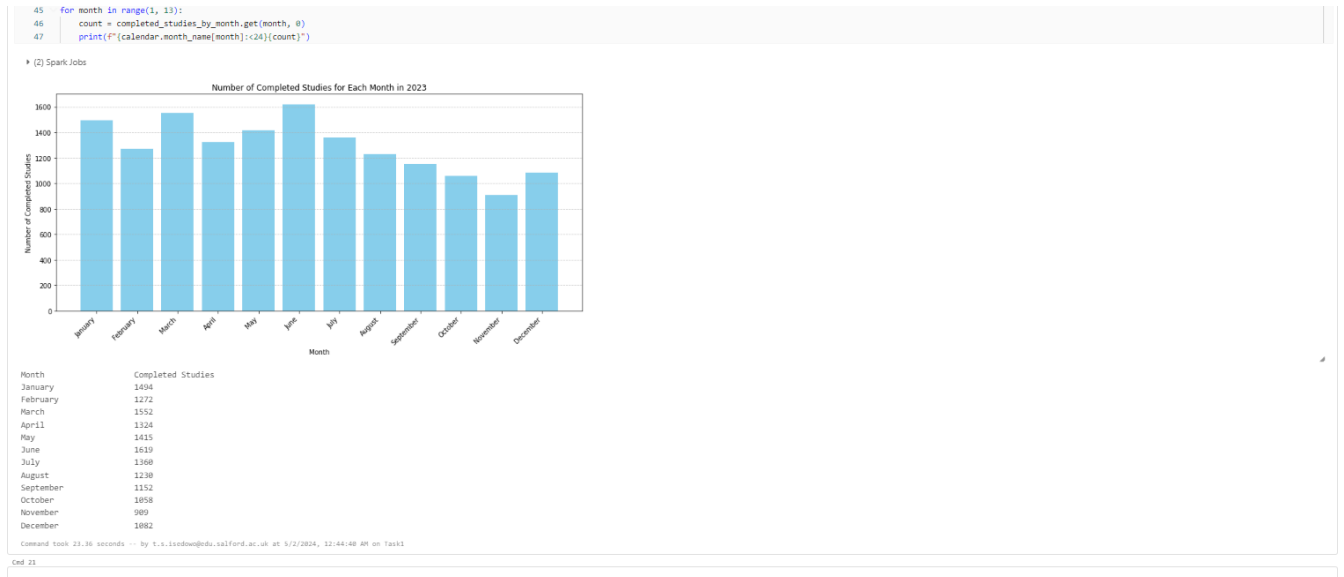
We trust that these dates truly represent when the studies were completed. Also, we believe that the status assigned to each study accurately indicates whether it's completed or not.

Moreover, we're working on the assumption that the dataset includes all studies completed in 2023 without any important missing or incorrect entries. These assumptions are essential to ensure that our analysis accurately shows how many studies were completed each month in 2023.

IMPLEMENTATION IN RDD

This PySpark script analyzes and visualizes the count of completed studies per month in 2023 using RDDs. It employs matplotlib, datetime, and calendar libraries, with `completion_index` set to 13 for identifying completion dates. The `count_completed_studies` function extracts completion dates, checks if studies are completed in 2023, and returns a tuple with the month and 1 if conditions are met, else (0, 0). Data integrity is ensured by filtering out rows with less than 14 fields. The map transformation applies `count_completed_studies` to each row, while `reduceByKey` aggregates counts for each month. Results are collected into a dictionary using `collectAsMap`. The matplotlib creates a bar plot with month names on the x-axis and completed study counts on the y-axis.





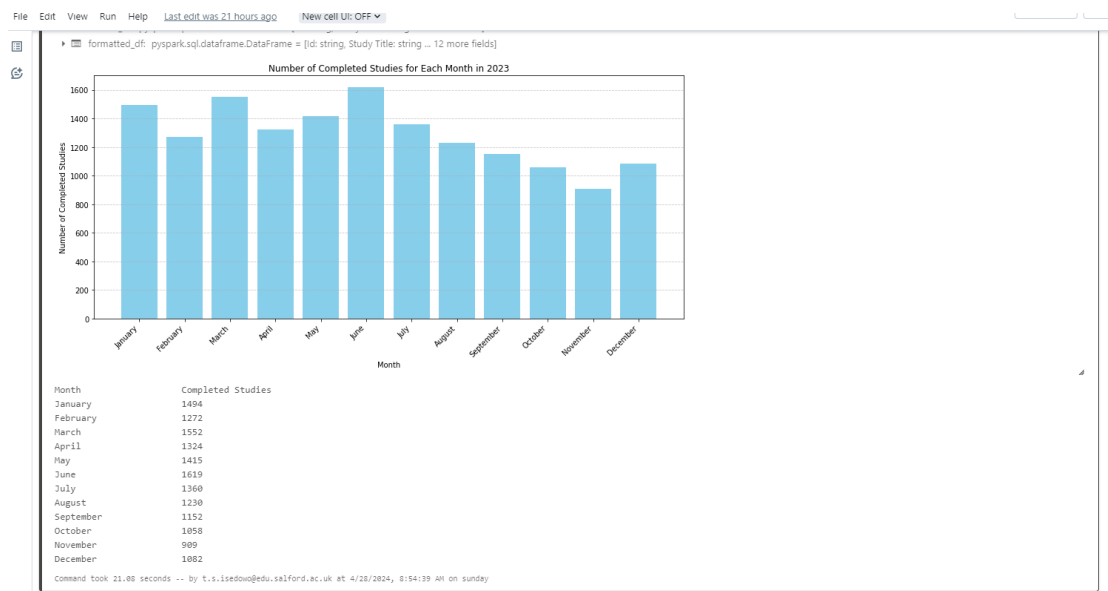
IMPLEMENTATION IN DATAFRAME

The `pyspark.sql.functions` were imported which helps to manipulate park dataframe. The `calendar`, which is a standard python library provides various calendar functions, `completion_year` and `completion_status` are variables used to filter the clinical trial data which filter the data based on completion year. (2023) and status completed.

The `filtered_df` filters the `Clinicaltrial_df` DataFrame to include only rows where the 'Status' column matches the `completion_status` and the year of the 'Completion' column matches the `completion_year`. Then, the `formatted_df` function convert the completion column to a date format.

The `completed_studies_by_month_dict` converts the result of the previous step into a dictionary where the keys are the month numbers, and the values are the corresponding counts of completed studies. The `matplotlib` then visualizes the number of completed studies in a bar plot.

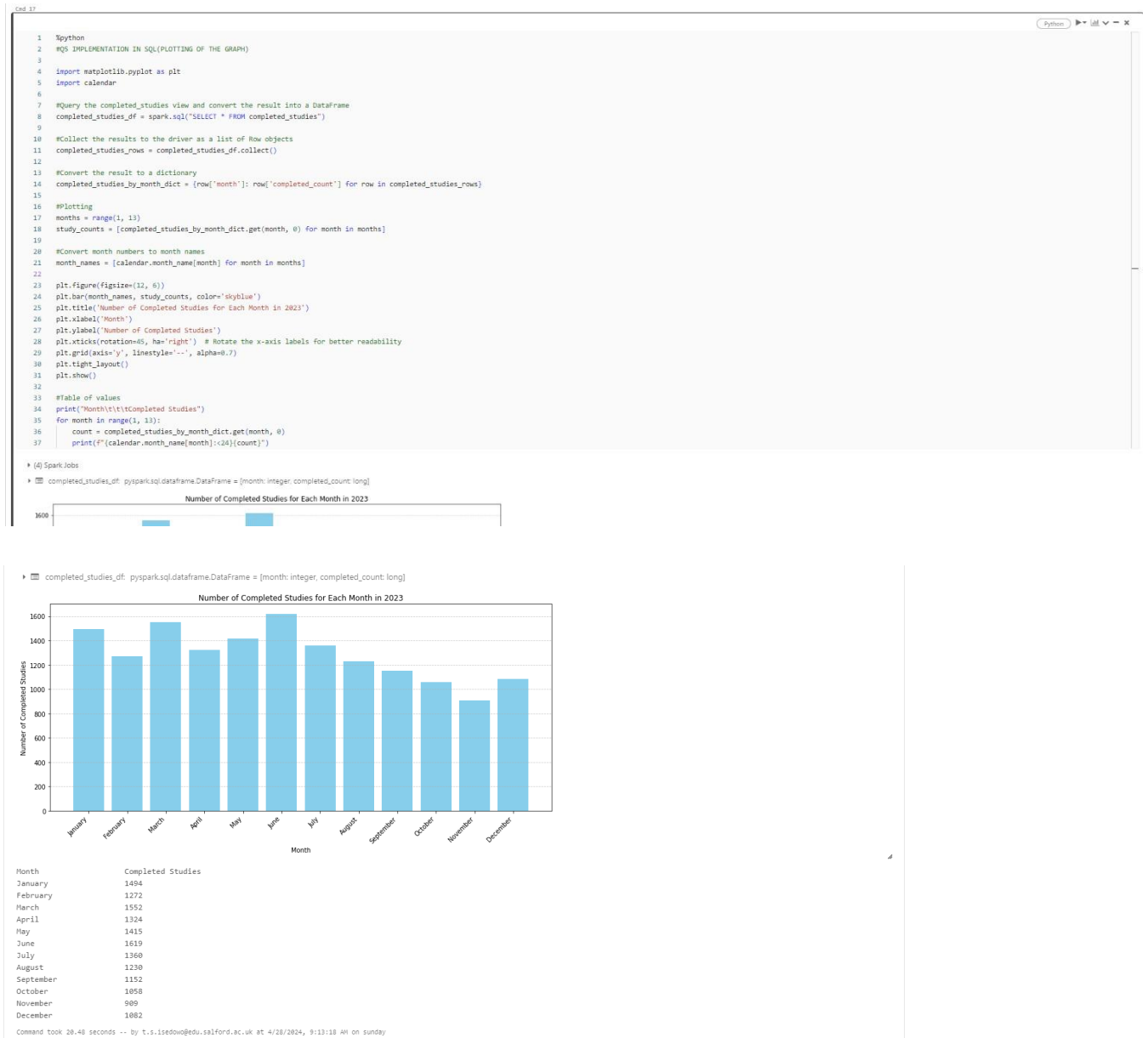
```
1 #QS IMPLEMENTATION IN RDD
2 import matplotlib.pyplot as plt
3 from pyspark.sql.functions import month, year, to_date
4 import calendar
5
6 # Define completion date and status filters
7 completion_year = 2023
8 completion_status = 'COMPLETED'
9
10 # Filter the dataframe to include only rows with the correct number of fields
11 filtered_df = clinicaltrial_df.filter((clinicaltrial_df['Status'] == completion_status) \
12                                     .filter(year('Completion') == completion_year))
13 # Convert the completion column to a date format
14 formatted_df = filtered_df.withColumn('Completion', to_date('Completion', 'yyyy-MM-dd'))
15
16 # Group by month and count completed studies for each month
17 completed_studies_by_month = formatted_df.groupBy(month('Completion').alias('month')) \
18                                     .count() \
19                                     .orderBy('month') \
20                                     .collect()
21
22 # Convert the result to a dictionary
23 completed_studies_by_month_dict = {row['month']: row['count'] for row in completed_studies_by_month}
24
25 # Plotting
26 months = range(1, 13)
27 study_counts = [completed_studies_by_month_dict.get(month, 0) for month in months]
28
29 # Convert month numbers to month names
30 month_names = [calendar.month_name[month] for month in months]
31
32 plt.figure(figsize=(12, 6))
33 plt.bar(month_names, study_counts, color='skyblue')
34 plt.title('Number of Completed Studies for Each Month in 2023')
35 plt.xlabel('Month')
36 plt.ylabel('Number of Completed Studies')
37 plt.xticks(rotation=45, ha='right') # Rotate the x-axis labels for better readability
38 plt.grid(axis='y', linestyle='--', alpha=0.7)
39 plt.tight_layout()
40 plt.show()
41
42 # Table of values
43 print("Month\t\t\tCompleted Studies")
44 for month in range(1, 13):
45     count = completed_studies_by_month_dict.get(month, 0)
46     print(f"{calendar.month_name[month]}:{count}")
```



IMPLEMENTATION IN SQL

A temporary view named `completed_studies` is created to store the aggregated results. This view aggregates the completed studies count for each month in the year 2023.

The `MONTH()` function extracts the month from the Completion date column and the `COUNT(*)` function counts the number of records for each month. The `WHERE` clause filters the data to include only completed studies (`STATUS = 'COMPLETED'`) and studies completed in the year 2023 (`YEAR(Completion) = 2023`). The `GROUPBY` and `ORDER BY` group the data by month and order the result by month.



DISCUSSION OF THE RESULT

The consistent outcomes achieved through all three implementation approaches (RDD, DataFrame, and SQL) highlight the dependability and resilience of the analysis. The tabulated data showcases the quantity of completed studies for each month throughout 2023, revealing fluctuations in completion rates. For instance, June recorded the highest number of studies completed (1619), while November exhibited the lowest (909). These fluctuations might be influenced by factors like funding availability, seasonal variations in research endeavors, or specific project timelines. Moreover, presenting the data visually via a bar chart could offer a clearer depiction of the distribution of completed studies across months, facilitating the identification of any noticeable patterns or trends. In sum, the consistency in findings across diverse

implementation methods bolsters the credibility of the results and underscores the reliability of the analysis conducted.

FURTHER ANALYSIS

REASONS FOR FURTHER ANALYSIS

The initial analysis focused on identifying the top 5 conditions and their frequencies within the clinical trial dataset. However, the broader aim extends beyond merely listing existing conditions. It emphasizes the significance of understanding the current interventions available for preventing, managing, or eliminating these conditions. This understanding can reveal the most common and effective methods currently in use, as well as identify areas where new interventions are needed. By knowing what interventions are already in place, we can work towards contributing to advancing healthcare outcomes and fostering the development of impactful solutions.

EXECUTION STRATEGY

The additional analysis was implemented using the PySpark RDD, the Interventions column data was retrieved from the clinical trial RDD named “myRDD2”. The Intervention column was in the 6th element of myRDD2 which corresponds to index 5. To extract the Intervention column data from myRDD2, I used the map transformation on the myRDD2. to return a new RDD named “interventions_rdd” that only contains the Intervention column data. The interventions_rdd.take(20) retrieves the first 20 elements from the interventions_rdd and the take action retrieves a specified number of elements from the interventions_rdd and returns them as a list

```
1 #Further analysis 1
2
3 #extracting interventions and their frequencies
4 interventions_rdd = myRDD2.map(lambda x: x[5])
5 interventions_rdd.take(10)
6
```

▶ [1] Spark Jobs

```
Out[21]: [{"BEHAVIORAL: PRIDE 'Step 1' problem-solving intervention|BEHAVIORAL: Enhanced usual care",
'OTHER: Placebo|DIETARY_SUPPLEMENT: B-OHB',
'DRUG: Insulin (tight blood glucose control)',
'OTHER: SuperMELD',
'OTHER: adherence support system based on a vocal assistant',
'BEHAVIORAL: Person-centered inhibitory control training|BEHAVIORAL: Active behavioral response training',
'DRUG: Propranolol|DRUG: Placebo',
'BEHAVIORAL: Computerized brief alcohol intervention + IVR booster calls|BEHAVIORAL: Computerized brief alcohol intervention|BEHAVIORAL: Attention Control',
'OTHER: Pharmacist - Smoking Cessation Support',
'PROCEDURE: Biospecimen Collection|OTHER: Genetic Testing|OTHER: Medical Chart Review']
```

Command took 1.27 seconds -- by t.s.isidoro@edu.salford.ac.uk at 5/2/2024, 12:49:13 AM on Task1

After isolating the Intervention column from the clinical trial RDD, some cleaning was performed. A filter operation was performed on `interventions_rdd` to remove any elements (interventions) from the RDD where the value is an empty string ' ' and the `take(10)` collect the first 10 element from the filtered RDD "`interventions_rdd1`" and display them

```
Cnd 22

1 #remove empty column
2 interventions_rdd1 = interventions_rdd.filter(lambda x: x != '')
3 interventions_rdd1.take(10)

Python ▶ ▼ - x

▶ (1) Spark Jobs

Out[111]: ["BEHAVIORAL: PRIDE 'Step 1' problem-solving intervention|BEHAVIORAL: Enhanced usual care",
'OTHER: Placebo|DIETARY_SUPPLEMENT: β-OHB',
'DRUG: Insulin (tight blood glucose control)',
'OTHER: SuperMELD',
'OTHER: adherence support system based on a vocal assistant',
'BEHAVIORAL: Person-centered inhibitory control training|BEHAVIORAL: Active behavioral response training',
'DRUG: Propranolol|DRUG: Placebo',
'BEHAVIORAL: Computerized brief alcohol intervention + IVR booster calls|BEHAVIORAL: Computerized brief alcohol intervention|BEHAVIORAL: Attention Control',
'OTHER: Pharmacist - Smoking Cessation Support',
'PROCEDURE: Biospecimen Collection|OTHER: Genetic Testing|OTHER: Medical Chart Review']

Command took 1.42 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 12:35:32 PM on B00TT
```

```
Cnd 23
```

The elements of `interventions_rdd1` were split by comma, the `flatMap` transformation applies a function to each element of the `interventions_rdd1` and flattens the result. `take(10)` action is used to retrieve the first 10 elements of the resulting `interventions_rdd2`.

```
LOADING TOOL: 4.19 seconds -- BY t.s.isedowo@edu.salford.ac.uk AT 5/1/2024, 12:35:39 PM ON TASK4

Cnd 23

1 #split by comma
2 interventions_rdd2 = interventions_rdd1.flatMap(lambda x: x.split(","))
3 interventions_rdd2.take(10)

Python ▶ ▼ - x

▶ (1) Spark Jobs

Out[10]: ["BEHAVIORAL: PRIDE 'Step 1' problem-solving intervention|BEHAVIORAL: Enhanced usual care",
'OTHER: Placebo|DIETARY_SUPPLEMENT: β-OHB',
'DRUG: Insulin (tight blood glucose control)',
'OTHER: SuperMELD',
'OTHER: adherence support system based on a vocal assistant',
'BEHAVIORAL: Person-centered inhibitory control training|BEHAVIORAL: Active behavioral response training',
'DRUG: Propranolol|DRUG: Placebo',
'BEHAVIORAL: Computerized brief alcohol intervention + IVR booster calls|BEHAVIORAL: Computerized brief alcohol intervention|BEHAVIORAL: Attention Control',
'OTHER: Pharmacist - Smoking Cessation Support',
'PROCEDURE: Biospecimen Collection|OTHER: Genetic Testing|OTHER: Medical Chart Review']

Command took 1.14 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 12:36:12 AM on Task3
```

```
Cnd 24
```

Another `flatMap` transformation is applied to split each intervention record by the '|' delimiter. This transformation flattens the resulting nested lists into a single list. For each record `x`, the `split("|")` function is used to split the interventions separated by '|

```
Cnd 24

1 #split by pipe
2 interventions_rdd3 = interventions_rdd2.flatMap(lambda x: x.split("|"))
3 interventions_rdd3.take(10)

Python ▶ ▼ - x

▶ (1) Spark Jobs

Out[48]: ["BEHAVIORAL: PRIDE 'Step 1' problem-solving intervention",
'BEHAVIORAL: Enhanced usual care',
'OTHER: Placebo',
'DIETARY_SUPPLEMENT: β-OHB',
'DRUG: Insulin (tight blood glucose control)',
'OTHER: SuperMELD',
'OTHER: adherence support system based on a vocal assistant',
'BEHAVIORAL: Person-centered inhibitory control training',
'BEHAVIORAL: Active behavioral response training',
'DRUG: Propranolol']

Command took 1.41 seconds -- by t.s.isedowo@edu.salford.ac.uk AT 5/1/2024, 2:21:13 PM ON TASK2
```

```
Cnd 25
```

A new RDD (interventions_rdd3) was created by transforming each element of the interventions_rdd2 RDD into a tuple where the intervention is the key, and the value is 1. Essentially, it's mapping each intervention to the number 1, indicating that each intervention has occurred once.



```
1 #pair rdd to count the frequency
2 interventions_rdd4 = interventions_rdd3.map(lambda c: (c, 1))
3 interventions_rdd4.take(10)
```

▶ (1) Spark Jobs

```
Out[50]: [('BEHAVIORAL: PRIDE 'Step 1' problem-solving intervention', 1),
('BEHAVIORAL: Enhanced usual care', 1),
('OTHER: Placebo', 1),
('DIETARY_SUPPLEMENT: β-OHB', 1),
('DRUG: Insulin (tight blood glucose control)', 1),
('OTHER: SuperMELD', 1),
('OTHER: adherence support system based on a vocal assistant', 1),
('BEHAVIORAL: Person-centered inhibitory control training', 1),
('BEHAVIORAL: Active behavioral response training', 1),
('DRUG: Propranolol', 1)]
```

Command took 1.24 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 2:22:22 PM on TASK2

The reduceByKey transformation in PySpark was used to sum up the values associated with each key in the RDD interventions_rdd4.



```
1 #reduce by key to sum up the values
2 interventions_rdd5 = interventions_rdd4.reduceByKey(lambda x,y: x+y)
3 interventions_rdd5.take(10)
```

▶ (1) Spark Jobs

```
Out[44]: [('BEHAVIORAL: Enhanced usual care', 25),
('BEHAVIORAL: Attention Control', 66),
('OTHER: Pharmacist - Smoking Cessation Support', 1),
('OTHER: Blood draw', 86),
('BEHAVIORAL: SwaysA Inactive wait-list control', 1),
('DRUG: human fetal-derived Neural Stem Cells (hNSCs)', 1),
('DRUG: PREDNISOLONE', 1),
('PROCEDURE: Doing Different ICU techniques and skills', 1),
('DEVICE: rTMS', 107),
('OTHER: Intra-abdominal pressure during laparoscopy', 1)]
```

Command took 10.73 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 2:17:31 PM on TASK2

In addition, each intervention is split by ":" using the split(":") method, and only the first part (intervention type) is extracted. The strip() method is used to remove any leading or trailing whitespace that is still present. The reduceByKey() transformation is applied to aggregate the counts for each intervention type and sortBy() was used to sort the interventions by their count in descending order. The take() action is used to extract the top 10 interventions after sorting.

```
Cmd 27

1 # Split each intervention by ":" and extract the intervention type
2 interventions_rdd5_cleaned = interventions_rdd5.map(lambda x: (x[0].split(":")[0].strip(), x[1]))
3
4 # Aggregate the counts for each intervention type
5 interventions_rdd5_aggregated = interventions_rdd5_cleaned.reduceByKey(lambda x, y: x + y)
6
7 # Sort the aggregated interventions by count in descending order
8 interventions_rdd5_sorted = interventions_rdd5_aggregated.sortBy(lambda x: x[1], ascending=False)
9
10 # Take the first 20 interventions after aggregation and sorting
11 interventions_rdd6 = interventions_rdd5_sorted.take(20)
12
13 # Display the top 20 interventions
14 interventions_rdd6
15

▶ (3) Spark Jobs

Out[45]: [('DRUG', 353617),
('OTHER', 135985),
('DEVICE', 81400),
('BEHAVIORAL', 74219),
('PROCEDURE', 71472),
('BIOLOGICAL', 40817),
('DIETARY_SUPPLEMENT', 24422),
('DIAGNOSTIC_TEST', 19660),
('RADIATION', 10937),
('GENETIC', 3557)]

Command took 2.89 seconds -- By t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 2:17:46 PM on TASK2

Cmd 28
```

Finally, each pair has its elements swapped, so that the count comes first, and the intervention type comes second.

```
Cmd 28

1 # Swap the positions of count and intervention type
2 interventions_top_10 = [(count, intervention) for intervention, count in interventions_rdd6]
3
4 # Display the top 20 interventions with count first
5 interventions_top_10
6

Out[47]: [(353617, 'DRUG'),
(135985, 'OTHER'),
(81400, 'DEVICE'),
(74219, 'BEHAVIORAL'),
(71472, 'PROCEDURE'),
(40817, 'BIOLOGICAL'),
(24422, 'DIETARY_SUPPLEMENT'),
(19660, 'DIAGNOSTIC_TEST'),
(10937, 'RADIATION'),
(3557, 'GENETIC')]

Command took 0.04 seconds -- By t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 2:18:31 PM on TASK2

Cmd 29
```

After further analysis was carried out to understand the current interventions available for preventing, managing, or eliminating the conditions, it was discovered that the highest intervention to the clinical trial conditions is Drug based on the analysis of the Top 10 interventions.

TASK 2

COLLABORATIVE FILTERING RECOMMENDER SYSTEM

BUILDING A COLLABORATIVE FILTERING RECOMMENDER SYSTEM FOR GAME RECOMMENDATIONS USING MLLIB AND MLFLOW

In the realm of recommendation systems, collaborative filtering stands out as a powerful technique for suggesting items based on user preferences and behaviors. Leveraging Apache Spark's MLlib, we embark on a journey to construct a collaborative filtering recommender system. This guide navigates through the process of utilizing the steam-200k.csv dataset, a rich collection of user interactions with video games on the Steam platform.

DESCRIPTION OF THE SET UP USED FOR THE TASK.

1. Databricks Account Creation

The initial phase commenced with registering for a Databricks account, granting access to the platform's workspace for managing project-related tasks. In this experiment, we opted for the community edition to fulfill our requirements.

2. Cluster Configuration

We configured a Databricks cluster optimized to meet the computational requirements of training our recommender system. For this experiment, we selected Databricks runtime version 15.0 ML (Scala 2.12, Spark 3.5.0), ensuring compatibility and performance for our machine learning tasks.

3. Notebook Creation:

Within our Databricks workspace, we established a notebook named "Tomisin_Isedowo_ml" to facilitate the development of our collaborative filtering recommender system. This notebook served as an interactive environment for writing and executing code, allowing for seamless development and experimentation.

4 Data Preparation:

We imported the "steam-200k.csv" dataset, which captures user interactions with video games on the Steam platform, into the Databricks workspace. The dataset was stored in the DBFS file layer with the file path "/FileStore/tables".

5 Library Installation:

We ensured the installation of essential libraries and dependencies within the Databricks cluster to support the development process. This included installing Apache Spark's MLflow, which plays a crucial role in managing the machine learning lifecycle. Additionally, MLLib, a component of Apache Spark, was readily available within the Databricks platform to support our recommender system development.



The image shows two screenshots of a Databricks code editor. The first screenshot, labeled 'Cell 1', contains the following code:

```
1 import mlflow
2 mlflow.pyspark.ml.autolog()
```

Below the code, it says 'Command took 1.44 seconds -- By t.s.isidw@edu.salford.ac.uk at 5/2/2024, 1:02:39 AM on new'. The second screenshot, labeled 'Cell 2', contains the following code:

```
1 dbutils.fs.ls("/FileStore/tables")
```

Below the code, it shows the output of the command: `["steam_200k.csv"]`. The third screenshot, labeled 'Cell 8', contains the following code:

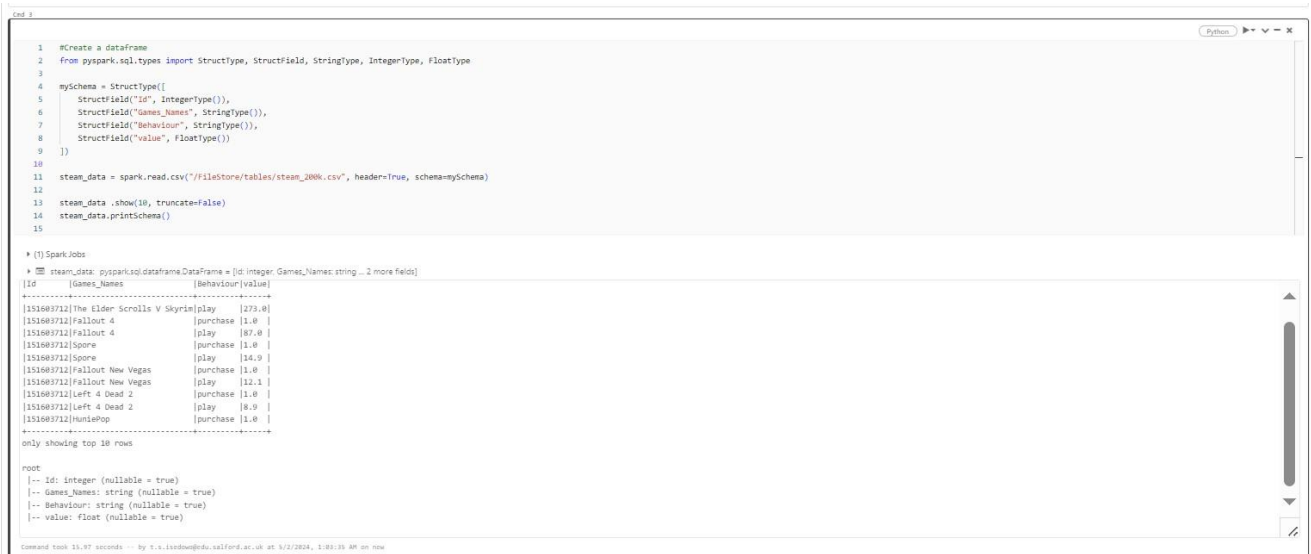
```
1 # Import necessary libraries
2 from pyspark.sql import SparkSession
3 from pyspark.ml.evaluation import RegressionEvaluator
4 from pyspark.ml.recommendation import ALS
5 from pyspark.sql.functions import col
6 from pyspark.sql.window import Window
7 from pyspark.sql.functions import row_number
```

Below the code, it says 'Command took 0.11 seconds -- By t.s.isidw@edu.salford.ac.uk at 5/2/2024, 1:04:37 AM on new'.

LOADING DATA INTO SPARK DATAFRAME AND ANY EXPLORATORY ANALYSIS OR VISUALISATION CARRIED OUT PRIOR TO TRAINING.

- 1) **Checking Dataset Availability:** We use a code `dbutils.fs.ls("/FileStore/tables")` in Databricks to see what files are in a specific directory (`/FileStore/tables`). This helps us confirm if our dataset is there before we start working with it.
- 2) **Preparing for Data Analysis:** Before diving into analyzing the data, we set up a plan for how the data should look. This plan, called a schema, tells us what kind of information each column in our dataset should contain. For example, we decided that our dataset should have columns named "Id", "Games_Names", "Behaviour", and "value", each with specific types of data like numbers or text.
- 3) **Loading the Dataset:** Once we've made our plan, we read the actual dataset into our analysis environment. The dataset is a CSV file located at `/FileStore/tables/steam_200k.csv`. We use the plan we made (our schema) to make sure the data is organized correctly. By setting `header=True`, we tell the system that the first row of the CSV file has the names of the columns. This helps the system understand how the data is organized.

- 4) **Everythingthe Data:** After we load the data, we take a quick look at the first 10 rows to see if everything looks right. This helps us make sure that the data was loaded correctly and matches our schema.



```
1 #Create a dataframe
2 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, FloatType
3
4 mySchema = StructType([
5     StructField("Id", IntegerType()),
6     StructField("Games_Names", StringType()),
7     StructField("Behaviour", StringType()),
8     StructField("value", FloatType())
9 ])
10
11 steam_data = spark.read.csv("/FileStore/tables/steam_200k.csv", header=True, schema=mySchema)
12
13 steam_data.show(10, truncate=False)
14 steam_data.printSchema()
15
```

*(1) Spark Jobs

Id	Games_Names	Behaviour	value
151089712	The Elder Scrolls V Skyrim	play	129.0
151089712	Fallout 4	purchase	1.0
151089712	Fallout 4	play	187.0
151089712	Spore	purchase	1.0
151089712	Spore	play	14.0
151089712	Fallout New Vegas	purchase	1.0
151089712	Fallout New Vegas	play	112.1
151089712	Left 4 Dead 2	purchase	1.0
151089712	Left 4 Dead 2	play	18.9
151089712	Huntahop	purchase	1.0

only showing top 10 rows

```
root
|-- Id: integer (nullable = true)
|-- Games_Names: string (nullable = true)
|-- Behaviour: string (nullable = true)
|-- value: float (nullable = true)
```

Command took 15.97 seconds -- by t.v.lindner@ku.salford.ac.uk at 3/2/2024, 1:00:19 AM on row

DATA PREPARATION AND PRE-PROCESSING CARRIED OUT PRIOR TO TRAINING THE MODEL.

Prior to training the model, data preparation and preprocessing were undertaken to ensure the dataset's suitability for training a machine learning model. The following grouping was done to understand the structure of the dataset.

1) Grouping by "Behaviour" Column

The `groupBy("Behaviour")` function groups the data in the `steam_data` based on the unique values in the "Behaviour" column. This operation creates groups where all rows with the same value in the "Behaviour" column are grouped together. The `count()` function is applied to each group, which calculates the number of rows in each group. This results in `steam_data` with two columns: "Behaviour" and "count". The "Behaviour" column contains unique values from the original "Behaviour" column, and the "count" column contains the number of rows corresponding to each unique "Behaviour" value. Then the `display()` function is used to display the result of `steam_data` in a visually appealing format.

Cell 4

```
1 steam_data.groupBy("Behaviour").count().display()
```

(2) Spark Jobs

Table

	Behaviour	count
1	purchase	129510
2	play	70489

2 rows | 5.79 seconds runtime

Command took 5.79 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/30/2024, 10:57:17 PM on anoo

Refreshed 24 minutes ago

Cell 5

2) Grouping by "Id" Column

The `groupBy("Id")` function collects all rows with the same value in the "Id" column into groups. After grouping by "Id", the `count()` function is applied to each group. This function calculates the number of rows in each group, effectively counting the occurrences of each unique "Id" value in the dataset. Then, the `display()` function is used to visually display the result.

Cell 5

```
1 steam_data.groupBy("Id").count().display()
```

(2) Spark Jobs

Table

	Id	count
1	16167221	57
2	166705920	11
3	244878937	2
4	99992274	2
5	174415183	10
6	156156544	2
7	152861732	34

10,000 rows | Truncated data | 4.45 seconds runtime

Command took 4.45 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/2/2024, 1:04:13 AM on now

Refreshed 1 hour ago

Cell 6

3) Grouping by "Games_Names" Column

The `groupBy("Games_Names")` function collects all rows with the same value in the "Games_Names" column into groups. Then, the `count()` function is applied to each group. This function calculates the number of rows in each group, effectively counting the occurrences of each unique game name in the dataset and the `display()` function is used to visually display the result

cmd 6

```
1 steam_data.groupby("Games_Names").count().display()
```

(2) Spark Jobs

Table

	Games_Names	count
1	Dota 2	9692
2	METAL GEAR SOLID V THE PHANTOM PAIN	134
3	LEGO Batman The Videogame	17
4	RIFT	236
5	Anodyne	14
6	Legend of Grimrock	93
7	Divinity Original Sin	138
8	Meltdown	4
9	SanctuaryRPG Black Edition	6
10	Snuggle Truck	14
11	Lunar Flight	16
12	Dungeons 2	12
13	Zuma's Revenge	7
14	HustleHeart	1
15	Ihf Handball Challenge 12	1

5,155 rows | 3.21 seconds runtime

Command took 3.21 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/2/2024, 1:04:21 AM on new

cmd 7

The `dropna()` function was applied to `steam_data`, which removes all rows containing any missing values.

The resulting DataFrame with missing values removed is assigned to a new variable `df_cleaned`. The `df_cleaned`, contains only the rows from the original DataFrame `steam_data` where all values are present (i.e., no missing values).

cmd 7

```
1 df_cleaned = steam_data.dropna()
```

df_cleaned: pyspark.sql.dataframe.DataFrame = [Id: integer, Games_Names: string ... 2 more fields]

Command took 0.15 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/30/2024, 10:57:46 PM on anoo

cmd 8

Moreso, The DataFrame `df_cleaned` contains rows representing user interactions with video games. The `filter()` function is applied to this DataFrame to select only those rows where the value in the "Behaviour" column is equal to "purchase". After filtering, the `select()` function is used to retain only specific columns from the filtered DataFrame.

In this case, the "Id", "Games_Names", and "value" columns are selected. These columns represent the user ID, the name of the game purchased, and any associated value respectively which is the focus for the task. The filtered and selected DataFrame is assigned to a new variable `data`. This new DataFrame contains only the rows where users have made purchases and includes the selected columns: "Id", "Games_Names", and "value".

```
Cmd 9
1 #Training the recommender system using Purchase
2 # This checked the dataset and find the behaviour column. Then it picks all parameter with purchase
3 data = df_cleaned.filter(col("Behaviour") == "purchase").select("Id", "Games_Names", "value")

data: pyspark.sql.dataframe.DataFrame = [Id: integer, Games_Names: string ... 1 more field]
Command took 0.22 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 12:48:07 AM on ooo

Cmd 10
```

The line `window = Window.orderBy("Games_Names")` defines a window specification for ordering rows by the "Games_Names" column. This window specification is used in subsequent operations involving window functions. The `unique_gameID = data.select("Games_Names").distinct()` selects the distinct game names from the DataFrame data. Then, `withColumn("GameID", row_number().over(window))` adds a new column "GameID" to the DataFrame `unique_gameID`. This column is populated with unique row numbers assigned to each distinct game name based on the order specified by the window. This effectively assigns a unique ID to each game name. After that, `joined_data = data.join(unique_gameID, "Games_Names", "left")` performs a left join between the original DataFrame data and the DataFrame `unique_gameID` on the "Games_Names" column. This helps to join the unique game IDs with the corresponding game names in the original DataFrame and `.drop("Games_Names").withColumnRenamed("GameID", "GameID")` drops the original "Games_Names" column from the joined DataFrame and renames the newly generated "GameID" column to "GameID". The `joined_data.display()` function is used to visually display the joined DataFrame `joined_data` with the game names replaced by their corresponding IDs.

```
1 # Window function help to assign IDs to the game names
2 window = Window.orderBy("Games_Names")
3 #We get the distinct game names so that the game name is not repeated.
4 unique_gameID = data.select("Games_Names").distinct().withColumn("GameID", row_number().over(window))
5 #The new dataframe is joined with existing dataframe, then replace the Game_names with the GameID generated in the new dataframe
6 joined_data = data.join(unique_gameID, "Games_Names", "left").drop("Games_Names").withColumnRenamed("GameID", "GameID")
7
8 joined_data.display()
```

▶ (4) Spark Jobs

unique_gameID: pyspark.sql.dataframe.DataFrame = [Games_Names: string, GameID: integer]

joined_data: pyspark.sql.dataframe.DataFrame = [Id: integer, value: float ... 1 more field]

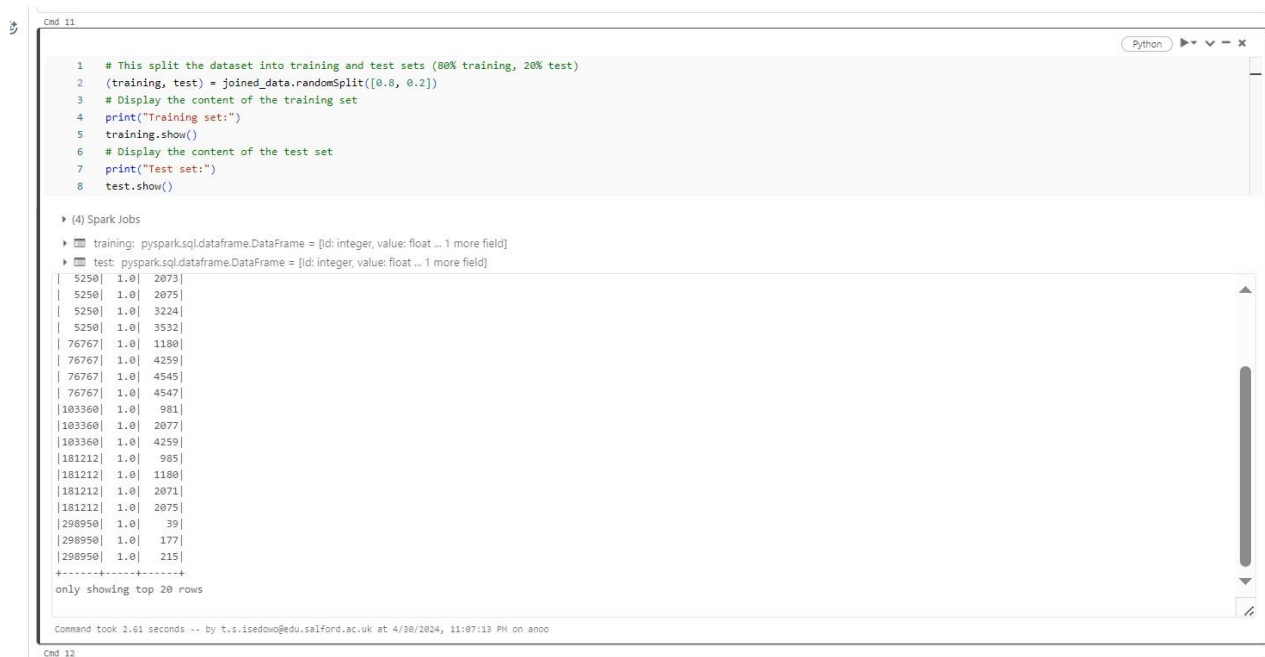
	Id	value	GameID
1	151603712	1	1679
2	151603712	1	3998
3	151603712	1	1680
4	151603712	1	2476
5	151603712	1	2194
6	151603712	1	3115
7	151603712	1	3214

10,000 rows | Truncated data | 5.77 seconds runtime

Command took 5.77 seconds -- by t.s.isedowo@edu.salford.ac.uk at 4/30/2024, 11:02:39 PM on anoo

Data Splitting

The dataset was split into training and test sets to evaluate the model's performance on the dataset. The `joined_data.randomSplit([0.8, 0.2])`, splits the DataFrame into two sets, training set and test set. The first argument `[0.8, 0.2]` specifies the proportions for splitting the data, with 80% allocated to the training set and 20% to the test set. Then the `test.show()` displays the content of the test set.



The screenshot shows a Jupyter Notebook interface with a code cell and its output. The code cell contains the following Python code:

```
1 # This split the dataset into training and test sets (80% training, 20% test)
2 (training, test) = joined_data.randomSplit([0.8, 0.2])
3 # Display the content of the training set
4 print("Training set:")
5 training.show()
6 # Display the content of the test set
7 print("Test set:")
8 test.show()
```

The output of the code cell shows the Spark Jobs and the content of the test set. The test set is a DataFrame with 20 rows, each containing three columns: an integer ID, a float value of 1.0, and another integer value. The output is as follows:

```
> (4) Spark Jobs
> training: pyspark.sql.dataframe.DataFrame = [id: integer, value: float ... 1 more field]
> test: pyspark.sql.dataframe.DataFrame = [id: integer, value: float ... 1 more field]

+-----+-----+
| 5250 | 1.0 | 2073 |
| 5250 | 1.0 | 2075 |
| 5250 | 1.0 | 3224 |
| 5250 | 1.0 | 3532 |
| 76767 | 1.0 | 1180 |
| 76767 | 1.0 | 4259 |
| 76767 | 1.0 | 4545 |
| 76767 | 1.0 | 4547 |
| 103360 | 1.0 | 981 |
| 103360 | 1.0 | 2077 |
| 103360 | 1.0 | 4259 |
| 181212 | 1.0 | 985 |
| 181212 | 1.0 | 1180 |
| 181212 | 1.0 | 2071 |
| 181212 | 1.0 | 2075 |
| 298950 | 1.0 | 39 |
| 298950 | 1.0 | 177 |
| 298950 | 1.0 | 215 |
+-----+-----+
only showing top 20 rows
```

The command took 2.61 seconds to execute. The output is displayed in a scrollable area with a vertical scrollbar on the right.

SELECTION OF HYPERPARAMETERS AND MODEL TRAINING AND EVALUATION AND MLFLOW EXPERIMENT TRACKING.

TRAINING A RECOMMENDATION MODEL USING THE ALTERNATING LEAST SQUARES (ALS) ALGORITHM

This ALS (maxIter=5, regParam=0.01, userCol="Id", itemCol="GameID", ratingCol="value", coldStartStrategy="drop"), initializes an ALS (Alternating Least Squares) model object. ALS is a collaborative filtering algorithm commonly used for recommendation systems. The maxIter=5 specifies the maximum number of iterations to run during training. In this case, the model will run for a maximum of 5 iterations and regParam=0.01 which is a regularization parameter was used to prevent overfitting by penalizing large parameter values. Also, the userCol="Id", itemCol="GameID", ratingCol="value" specify the names of the columns in the DataFrame training that represent the user ID, item ID (game ID), and rating (value) respectively. The coldStartStrategy="drop" specifies the strategy to handle cold-start scenarios, where new users or items are encountered during prediction. Setting it to "drop" ensures that rows with missing values in the user or item columns are dropped during prediction.

The model = als.fit(training), fits the ALS model to the training data. The fit() method takes the training data (training DataFrame) as input and learns the model parameters (user and item embeddings) based on the specified configuration and optimization objective resulting in a trained recommendation model (model) that can be used to make predictions on unseen data.



```
1 #training a recommendation model using the Alternating Least Squares (ALS) algorithm
2 als = ALS(maxIter=5, regParam=0.01, userCol="Id", itemCol="GameID", ratingCol="value",
3           coldStartStrategy="drop")
4 model = als.fit(training)
```

▶ (5) Spark Jobs

▼ (1) MLflow run

Logged 1 run to an experiment in MLflow. [Learn more](#)

2024/04/30 23:49:06 INFO mlflow.utils.autologging_utils: Created MLflow autologging run with ID '5b63eb9f10004a73b6f6ff4cf78bc2a2', which will track hyperparameters, performance metrics, model artifacts, and lineage information for the current pyspark.ml workflow

2024/04/30 23:49:52 WARNING mlflow.pyspark.ml: Model ALS_e2ecd51e282d will not be autologged because it is not allowlisted or because one or more of its nested models are not allowlisted. Call mlflow.spark.log_model() to explicitly log the model, or specify a custom allowlist via the spark.mlflow.pysparkml.autolog.logModelAllowlistFile Spark conf (see mlflow.pyspark.ml.autolog docs for more info).

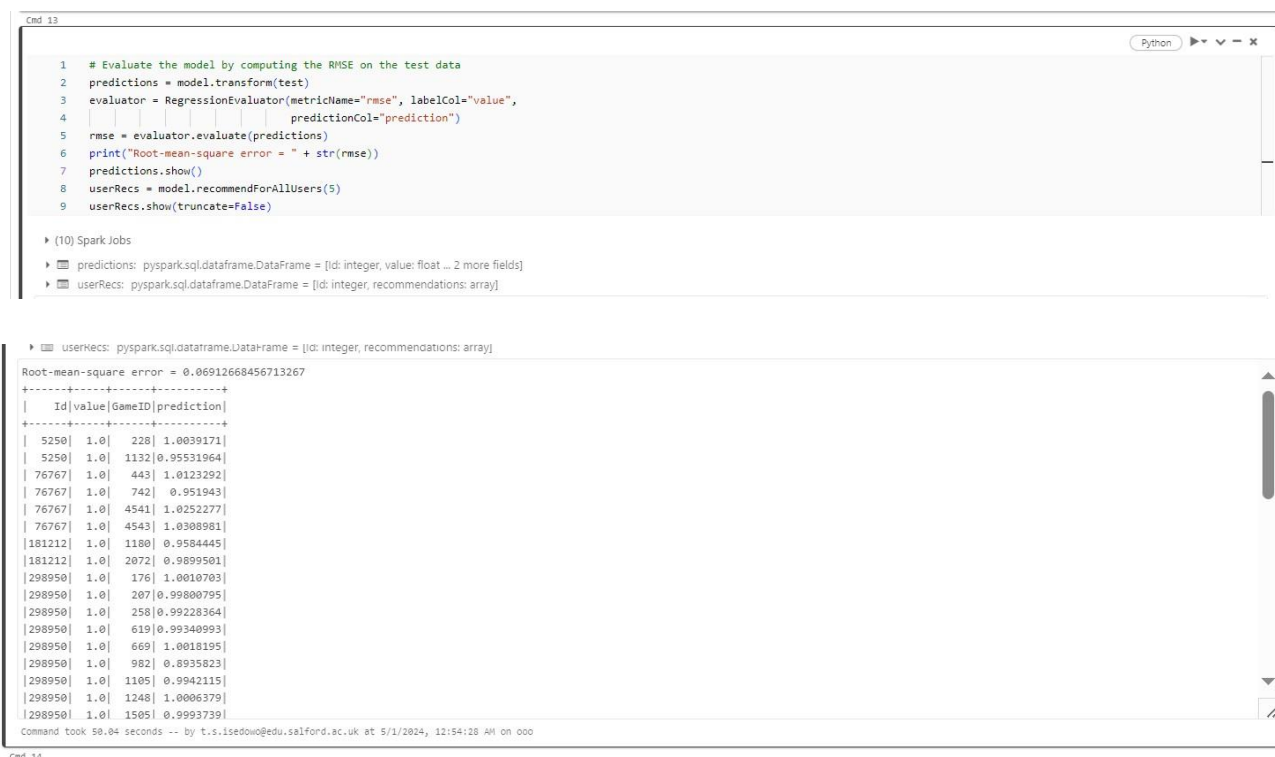
Command took 48.22 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 12:49:04 AM on ooo

MODEL EVALUATION

Before evaluation was initiated, `predictions = model.transform(test)` generates predictions on the test dataset using the trained ALS model (`model`). The `transform()` method applies the model to the test dataset and adds a new column "prediction" to the DataFrame predictions containing the predicted ratings for each user-item pair.

Evaluation was done using RMSE, the `evaluator = RegressionEvaluator(metricName="rmse", labelCol="value", predictionCol="prediction")`, initializes a regression evaluator object 'evaluator' to compute the Root Mean Square Error (RMSE) metric. The `rmse = evaluator.evaluate(predictions)`, computes the RMSE by comparing the actual ratings (`labelCol="value"`) in the test dataset with the predicted ratings (`predictionCol="prediction"`) generated by the model. The RMSE is a measure of the differences between predicted and observed values, providing an indication of the model's accuracy.

The `print("Root-mean-square error = " + str(rmse))`, helps to print out the computed RMSE, indicating the level of error in the model's predictions on the test dataset and `predictions.show()` displays the predictions DataFrame, showing the actual ratings, predicted ratings, and other relevant columns.



```
1 # Evaluate the model by computing the RMSE on the test data
2 predictions = model.transform(test)
3 evaluator = RegressionEvaluator(metricName="rmse", labelCol="value",
4                               predictionCol="prediction")
5 rmse = evaluator.evaluate(predictions)
6 print("Root-mean-square error = " + str(rmse))
7 predictions.show()
8 userRecs = model.recommendForAllUsers(5)
9 userRecs.show(truncate=False)
```

▶ (10) Spark Jobs

▶ predictions: pyspark.sql.dataframe.DataFrame = [id: integer, value: float ... 2 more fields]

▶ userRecs: pyspark.sql.dataframe.DataFrame = [id: integer, recommendations: array]

▶ userRecs: pyspark.sql.dataframe.DataFrame = [id: integer, recommendations: array]

Root-mean-square error = 0.06912668456713267

	Id	value	GameID	prediction
	5250	1.0	228	1.0039171
	5250	1.0	1132	0.95531964
	76767	1.0	443	1.0123292
	76767	1.0	742	0.951943
	76767	1.0	4541	1.0252277
	76767	1.0	4543	1.0308981
	181212	1.0	1180	0.9584445
	181212	1.0	2072	0.9899501
	298950	1.0	176	1.0010703
	298950	1.0	207	0.99800795
	298950	1.0	258	0.99228364
	298950	1.0	619	0.99340993
	298950	1.0	669	1.0018195
	298950	1.0	982	0.8935823
	298950	1.0	1105	0.9942115
	298950	1.0	1248	1.0006379
	298950	1.0	1505	0.9993739

Command took 50.04 seconds -- by t.s.isedowo@edu.salford.ac.uk at 5/1/2024, 12:54:28 AM on ooo

Recommendations for the test dataset

was beneficial for improving model performance.

Regularization Parameter (RegParam)

Regularization is a technique used to prevent overfitting by penalizing large parameter values. The regularization parameter controls the strength of this penalty. A higher regularization parameter imposes a stronger penalty, which can help prevent the model from fitting the noise in the data but may also lead to underfitting if set too high. Alternatively, a lower regularization parameter allows the model to fit the training data more closely but risks overfitting. In this task, I tested regularization parameters of 0.01 and 0.1. The best model selected had a regularization parameter of 0.01, indicating that a relatively low regularization strength was sufficient to prevent overfitting while still allowing the model to capture the underlying patterns in the data effectively.

```
1 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
2
3 # Define the parameter grid
4 param_grid = ParamGridBuilder() \
5     .addGrid(als.rank, [5,7]) \
6     .addGrid(als.maxIter, [5,7]) \
7     .addGrid(als.regParam, [0.01, 0.1]) \
8     .build()
9
10 # Define evaluator
11 evaluator = RegressionEvaluator(metricName="rmse", labelCol="value", predictionCol="prediction")
12
13 # Define cross-validator
14 cross_validator = CrossValidator(estimator=als,
15                                 estimatorParamMaps=param_grid,
16                                 evaluator=evaluator,
17                                 numFolds=5) # 5-fold cross-validation
18
19 # Fit cross-validator to the training data
20 cv_model = cross_validator.fit(training)
21
22 # Get the best model from cross-validation
23 best_model = cv_model.bestModel
24
25 # Evaluate the best model on the test set
26 predictions = best_model.transform(test)
27 rmse = evaluator.evaluate(predictions)
28 print("Root-mean-square error (Best Model) = " + str(rmse))
29 predictions.show()
30
31 # Print the best hyperparameters
32 print("Best Rank: ", best_model.rank)
33 print("Best Max Iterations: ", best_model._java_obj.parent().getMaxIter())
34 print("Best Regularization Parameter: ", best_model._java_obj.parent().getRegParam())
35
```

End 14

Python

21 Spark Jobs

9 MLflow runs

Logged 9 runs to an experiment in MLflow. Learn more

predictions: pyspark.sql.DataFrame = [id:integer,value:float... 2 more fields]

2024/05/01 21:06:59 WARNING mlflow.pyspark.ml: Model crossValidatorModel_4d51a5fe0b1d will not be autologged because it is not allowed or because one or more of its nested models are not autologged. Call mlflow.spark.log_model() to explicitly log the model, or specify a custom allowlist via the spark.mlflow.pyspark.ml.autolog.logModelAllowlistFile Spark conf (see mlflow.pyspark.ml.autolog docs for more info).

Root-mean-square error (Best Model) = 0.048343638428558915

id	value	rmse	prediction
5250	1.0	1152	6.5635958
76767	1.0	323	0.9885943
76767	1.0	727	0.9677381
76767	1.0	764	0.9711115
76767	1.0	4547	0.0626854
183360	1.0	979	0.9759166
183360	1.0	1152	0.9758787
183360	1.0	2678	0.9921936
183360	1.0	2879	0.9932121
181212	1.0	2678	0.98570784
208950	1.0	175	0.98243535
208950	1.0	176	0.9842177
208950	1.0	215	0.9873877
208950	1.0	388	0.9893474
208950	1.0	629	0.9833585
208950	1.0	659	0.9905615
208950	1.0	9821	0.9511981

Command took 36.09 minutes By f-a.lindood@brs.salford.ac.uk at 5/1/2024, 9:58:55 PM on Task555

End 15

MLFLOW EXPERIMENT TRACKING

MLflow plays a crucial role in machine learning development by ensuring transparency, reproducibility, and accountability. It allows us to track experiment details comprehensively, including parameters and metrics. This capability enables users to recreate experiments and validate results effectively.

Experiments

Tomisin_Isedowo_ml

Provide Feedback

Experiment ID: 3933729394456198 Artifact Location: dbfs:/databricks/mlflow-tracking/3933729394456198

Description Edit

Q metrics.rmse < 1 and params.model = "tree"

Time created

State: Active

Datasets

Sort: Created

Columns

Expand rows

Table

Chart

Evaluation

Preview

	Run Name	Created	Dataset	Duration	Source	Models
	nebulous-ram-805	20 hours ago		25.3s	Tomisin_...	-
	nervous-snipe-351	20 hours ago	-	46.8s	Tomisin_...	-
	gregarious-snail-900	22 hours ago	-	34.5s	Tomisin_...	-
	hilarious-horse-979	3 days ago	dataset (db44f644) Eval dataset (75a...	20.1s	Recomm...	-
	stately-mink-383	3 days ago	dataset (81b0b01a) Eval dataset (b3L...	24.7s	Recomm...	-
	worried-stork-10	3 days ago	dataset (66d5ce5d) Eval dataset (848...	31.6s	Recomm...	-
	incongruous-stag-121	20 days ago	dataset (2b03434b) Eval dataset (495...	18.3s	Recomm...	-
	amazing-asp-413	20 days ago	dataset (5fae253d) Eval dataset (905...	28.4s	Recomm...	-

Show more columns
(24 total)

CONCLUSION

The RMSE value of approximately 0.069 indicates the average difference between the actual ratings and the predicted ratings on the test dataset. The low RMSE value indicates that the ALS model performs well in predicting user ratings for video games, as the model's predictions are generally close to the actual ratings. The model's ability to generate accurate predictions and relevant recommendations demonstrates its effectiveness in capturing user preferences and delivering personalized content in the gaming domain. the results suggest that the ALS recommendation model trained on the dataset has the potential to improve user satisfaction and retention by providing tailored recommendations for video games.

The hyperparameter tuning process identified a combination of hyperparameters (rank=5, max iterations=7, regularization parameter=0.01) that resulted in a lower root-mean-square error (RMSE) on the test data compared to the initial model. This suggests that the tuned model is better at predicting user-item interactions. Additionally, the impact of each hyperparameter can be understood in terms of its effect on model complexity, training time, and regularization strength.

