

## **Big Data Tools and Technique Assignment**

By

Syed Ali Murtaza

## 1. Part-01

### STEPS TO IMPLEMENT THE PROJECT

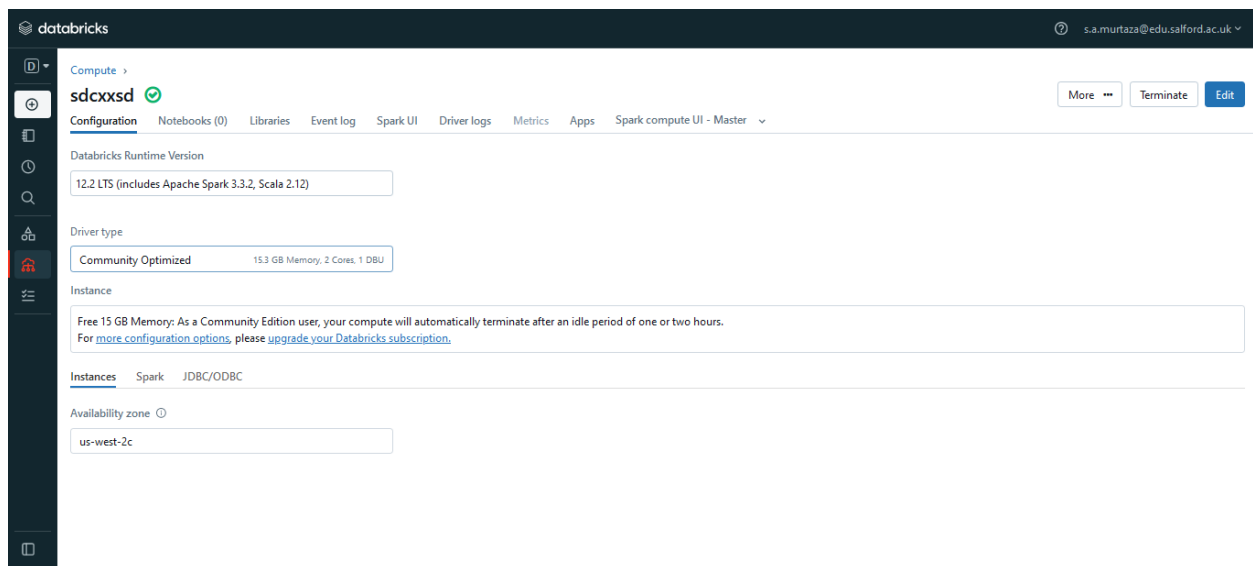
#### Description of the necessary configuration or prerequisites needed to successfully accomplish this task:-

In this part we have given different files with the name of clinicaltrial-2023, clinicaltrial-2021 and clinicaltrial-2020. So we need databricks platform to perform this task and we have to clean the data and answer the questions given in this task. So I follow several steps to answer the question.

#### 1. Fire up the Databricks workspace:-

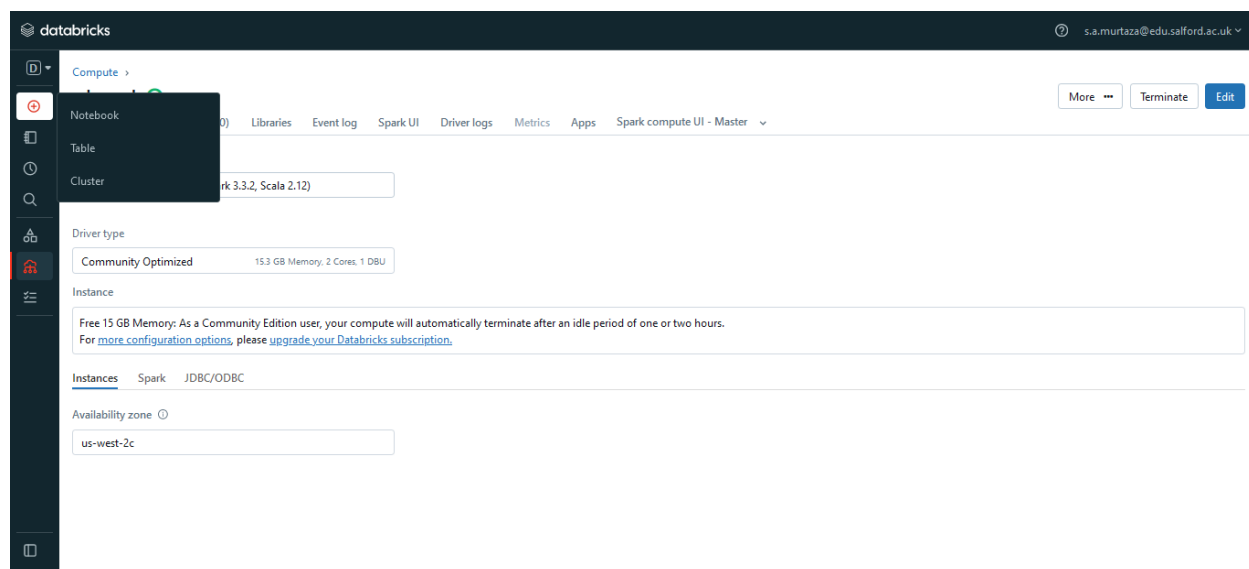
I logged in to the data bricks community edition account.

Then I created a cluster to start my analyses, because a cluster will give me access to use a machine. So for this purpose I clicked on compute and type the name of the cluster and select the most recent runtime as in this screen shot.

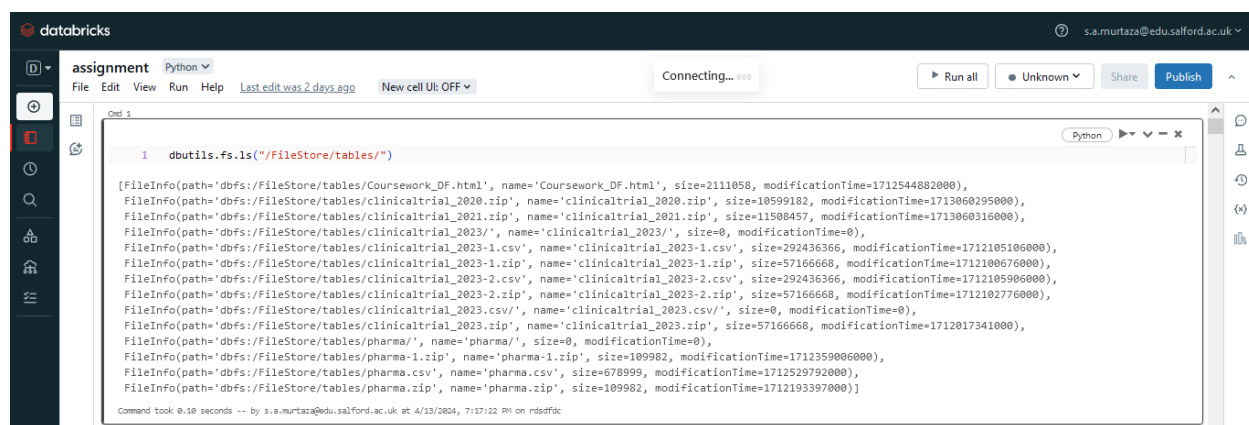


#### 2. Creating a Notebook:-

Then I click on “PLUS” icon and select notebook.



Then I name the notebook and attach the cluster (that I started before) to this notebook.



### 3. Data cleaning and preparation:-

Before answering the questions; we need to clean and prepare the data because the data is not in a correct format it contains all the values in a single row and single column, So I follow various steps to unzip the file and handle missing values, transform the data types and completely ensure data integrity.

1. **Load the dataset:** Load the files into the databricks.

## Create New Table

Data source ?

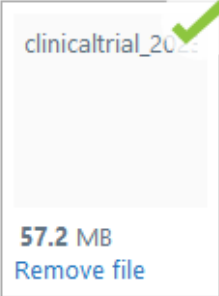
Upload File S3 Other Data Sources

DBFS Target Directory ?

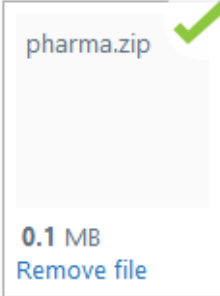
/FileStore/tables/ (optional) Select

Files uploaded to DBFS are accessible by everyone who has access to this workspace. [Learn more](#)

Files ?




57.2 MB  
[Remove file](#)



0.1 MB  
[Remove file](#)

✓File uploaded to /FileStore/tables/pharma-2.zip  
✓File uploaded to /FileStore/tables/clinicaltrial\_2023-4.zip

Create Table with UI
 Create Table in Notebook ?

All the files that I uploaded is present in this path /FileStore/tables/ by default. Now to ensure that if the file is uploaded successfully or not I use the `dbutils.fs.ls` command.

```

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

[FileInfo(path='dbfs:/FileStore/tables/Coursework_DF.html', name='Coursework_DF.html', size=2111058, modificationTime=1712544882000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2020.zip', name='clinicaltrial_2020.zip', size=10599182, modificationTime=1713060295000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2021.zip', name='clinicaltrial_2021.zip', size=11508457, modificationTime=1713060316000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023/', name='clinicaltrial_2023/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023-1.csv', name='clinicaltrial_2023-1.csv', size=292436366, modificationTime=1712105106000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023-1.zip', name='clinicaltrial_2023-1.zip', size=57166668, modificationTime=1712100676000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023-2.csv', name='clinicaltrial_2023-2.csv', size=292436366, modificationTime=1712105906000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023-2.zip', name='clinicaltrial_2023-2.zip', size=57166668, modificationTime=1712102776000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023.csv/', name='clinicaltrial_2023.csv/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023.zip', name='clinicaltrial_2023.zip', size=57166668, modificationTime=1712017341000),
 FileInfo(path='dbfs:/FileStore/tables/pharma/', name='pharma/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/pharma-1.zip', name='pharma-1.zip', size=109982, modificationTime=1712359006000),
 FileInfo(path='dbfs:/FileStore/tables/pharma.csv', name='pharma.csv', size=678999, modificationTime=1712529792000),
 FileInfo(path='dbfs:/FileStore/tables/pharma.zip', name='pharma.zip', size=109982, modificationTime=1712193397000)]

Command took 0.10 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/13/2024, 7:17:22 PM on rdsdfdc

```

Now from above we can see that both the files are present.

## 2. Unzipping the file:-

As the uploaded files are in zip form and dbutils toolkit has no tool to unzip So I copy the files to local file system and extract here using the shell command. After unzipping I put the extracted content back into DBFS. So first I am going to copy both file from DBFS to local file system using dbutils command.

```
Cmd 3
1 dbutils.fs.cp("/FileStore/tables/clinicaltrial_2023.zip", "file:/tmp/")
Out[2]: True
Command took 4.17 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/22/2024, 4:14:32 PM on dfvcds
```

```
Cmd 6
1 dbutils.fs.cp("/FileStore/tables/pharma.zip", "file:/tmp/")
Out[5]: True
Command took 0.69 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/22/2024, 4:15:07 PM on dfvcds
```

Now using ls command I checked that whether it is in shell command or not.

```
Cmd 7
1 %sh
2 ls /tmp/

Rserv
RtmpxJlAIn
chauffeur-daemon-params
chauffeur-daemon.pid
chauffeur-env.sh
clinicaltrial_2021.zip
clinicaltrial_2023.zip
custom-spark.conf
driver-daemon-params
driver-daemon.pid
driver-env.sh
hsperepdata_root
pharma.zip
python_lsp_logs
systemd-private-6271a62ed2b044e18cf155510f6bfdd-apache2.service-Yz0XGi
systemd-private-6271a62ed2b044e18cf155510f6bfdd-ntp.service-5e1kfj
systemd-private-6271a62ed2b044e18cf155510f6bfdd-systemd-logind.service-Squ12g
systemd-private-6271a62ed2b044e18cf155510f6bfdd-systemd-resolved.service-3fdDWi
tmp.hQysW0PF6
Command took 0.25 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/22/2024, 4:15:11 PM on dfvcds
```

From above SS we can see that both files are in tmp directory.

Then I have used UNIX command (Unzip) to unzip the files.

```
Cmd 8
1 %sh
2 unzip -d /tmp/ /tmp/clinicaltrial_2023.zip

unzip: cannot find or open /tmp/clinicaltrial_2023.zip, /tmp/clinicaltrial_2023.zip.zip or /tmp/clinicaltrial_2023.zip.ZIP.
Command took 0.12 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/13/2024, 7:19:43 PM on rdsdfdc

Cmd 11
1 %sh
2 unzip -d /tmp/ /tmp/pharma.zip

Archive: /tmp/pharma.zip
inflating: /tmp/pharma.csv
Command took 0.19 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk
```

Now, I aim to transfer the files from the local file system to DFBS. So for this purpose I make the directories.

```
Cmd 17
1 dbutils.fs.mkdirs("FileStore/tables/clinicaltrial_2023")

Out[36]: True
Command took 0.10 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk

Cmd 20
1 dbutils.fs.mkdirs("FileStore/tables/pharma")

Out[37]: True
Command took 0.10 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk
```

Now I'm transferring the files from the local file system to DBFS using the command `fs.mv`.

```
Cmd 21
1 dbutils.fs.mv("file:/tmp/clinicaltrial_2023.csv", "FileStore/tables/clinicaltrial_2023", True)

Out[38]: True
Command took 10.53 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk

Cmd 24
1 dbutils.fs.mv("file:/tmp/pharma.csv", "FileStore/tables/pharma", True)

Out[39]: True
Command took 0.39 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk
```

I verify the contents of the file located at `'FileStore/tables/FileName'` to ensure it contains the expected files

```
Cmd 25
1 dbutils.fs.ls("FileStore/tables/clinicaltrial_2023/")

Out[40]: [FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023/cleaned_clinicaltrial_2023.csv/', name='cleaned_clinicaltrial_2023.csv/', size=0, modificationTime=0),
  FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2023/clinicaltrial_2023.csv', name='clinicaltrial_2023.csv', size=292436366, modificationTime=1712880714000)]
Command took 0.19 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk

Cmd 28
1 dbutils.fs.ls("FileStore/tables/pharma/")

Out[41]: [FileInfo(path='dbfs:/FileStore/tables/pharma/pharma.csv', name='pharma.csv', size=678999, modificationTime=1712880717000)]
Command took 0.19 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk
```

Then I use the `dbutils.fs.head` command to see the format of file

clinicaltrial-2023. Although it is not well formatted and looks alike as in this screen shot.

```

1 dbutils.fs.head("FileStore/tables/clinicaltrial_2023/clinicaltrial_2023.csv")

[Truncated to first 65536 bytes]
Out[44]: "ID\tStudy Title\tAcronym\tStatus\tConditions\tInterventions\tSponsor\tCollaborators\tEnrollment\tFunder Type\tType\tStudy Design\tStart\tCompleto
n"
.....
.....\r\n"NCT03638471\tEffectiveness of a Problem-solving Intervention for Common Adolescent Mental Health Problems in Indi
a\tPRIDE\tCOMPLETED\tMental Health Issue (E.G., Depression, Psychosis, Personality Disorder," Substance Abuse)\tBEHAVIORAL: PRIDE 'Step 1'\tproblem-solving interven
tion\tBEHAVIORAL: Enhanced usual care\tSangath\tHarvard Medical School (HMS and HSDM)\tLondon School of Hygiene and Tropical Medicine\t250.0\t0\tOTHER\tINTERVENTIONAL\tAl
location: RANDOMIZED\tIntervention Model: PARALLEL\tMasking: DOUBLE (INVESTIGATOR", " OUTCOMES_ASSESSOR)|Primary Purpose: TREATMENT\t2018-08-20\t2019-02-2
8"
.....
.....\r\n"NCT05992571\tOral Ketone Monoester Supplementation and Resting-state Brain Connectivity\t\tRECRUITING\tCerebrovascular
Function\tCOGNITION\tOTHER: Placebo\tDIETARY SUPPLEMENT:  $\beta$ -OHB\tMcMaster University\tAlzheimer's Society of Brant", Haldimand Norfolk", Hamilton Halton\t30.0\t0\tOTHER\t
INTERVENTIONAL\tAllocation: RANDOMIZED\tIntervention Model: CROSSOVER\tMasking: TRIPLE (PARTICIPANT", INVESTIGATOR", OUTCOMES_ASSESSOR)|Primary Purpose: BASIC_SCIENCE
\t2023-10-25\t2024-0
8"
.....
.....\r\n"NCT00237471\tImpact of Tight Glycemic Control in Acute Myocardial Infarction\t\tTERMINATED\tMyocardial Infarct\tHyper
glycemia\tDRUG: Insulin (tight blood glucose control)\tMelbourne Health\tNational Health and Medical Research Council", Australian\tBristol-Myers Squibb\t40.0\t0\tOTHER
\tINTERVENTIONAL\tAllocation: RANDOMIZED\tIntervention Model: PARALLEL\tMasking: NONE\tPrimary Purpose: TREATMENT\t2005-10\t2006-0
5"
.....
.....\r\n"NCT03820271\tNew Prognostic Predictive Models of Mortality of Decompensated Cirrhotic Patients Waiting for Liver T
ransplantation\tSUPERHELD\tRECRUITING\tDecompensated Cirrhosis\tLiver Transplantation\tOTHER: SuperMELD\tAssistance Publique - Hôpitaux de Paris\t500.0\t0\tOTHER\tINTE
RVENTIONAL\tAllocation: NA\tIntervention Model: SINGLE_COHORT\tMasking: NONE\tPrimary Purpose: OTHER\t2020-10-01\t2023-10-0
1"
.....
Command took 0.39 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/11/2024, 5:11:29 PM on hjnk

```

### 3. Cleaning and Preparation of Data:-

## 1. Initializing Spark:-

First of all I initialize the spark session.

```
Cmd 2

1
2 from pyspark.sql import SparkSession
3 from pyspark.sql.types import StructType, StructField, StringType
4
5 # Initialize Spark session
6 spark = SparkSession.builder.appName("Data Cleaning").getOrCreate()
```

Command took 0.13 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 6:45:58 PM on fdcscx

## 2. Defining Schema:-

I defined schema using StructType and StructFeild. StructFeild is used to define the column names of the data frame while StructType is used to define the data type of the column. I define all the entries as string. This explicit schema definition ensures data consistency and integrity.

```
Cmd 3

1
2 schema = StructType([
3     StructField("Id", StringType(), True),
4     StructField("Study Title", StringType(), True),
5     StructField("Acronym", StringType(), True),
6     StructField("Status", StringType(), True),
7     StructField("Conditions", StringType(), True),
8     StructField("Interventions", StringType(), True),
9     StructField("Sponsor", StringType(), True),
10    StructField("Collaborators", StringType(), True),
11    StructField("Enrollment", StringType(), True),
12    StructField("Funder Type", StringType(), True),
13    StructField("Type", StringType(), True),
14    StructField("Study Design", StringType(), True),
15    StructField("Start", StringType(), True),
16    StructField("Completion", StringType(), True),
17 ])

Command took 0.07 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 4:41:51 PM on fdscsx
```

### 3. Load the data:-

Loading data from a file is necessary to perform data cleaning. So I loaded the data into an rdd using `spark.sparkContext.textFile(file_path)`.

```
Cmd 4

1
2 file_path = "/FileStore/tables/clinicaltrial_2023/clinicaltrial_2023.csv"
3 raw_rdd = spark.sparkContext.textFile(file_path)

Command took 0.55 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 4:42:52 PM on fdscsx
```

### 4. Deleting a directory and its content:-

In this code I makes a function name `delete_dir` to recursively delete directories and their contents in Databricks FileStore. Then I specifies directories or files to be deleted and deletes each path using the defined function. Finally, I print a message indicating the completion of cleanup.

```
1 # No need to mount DBFS root (it's already mounted at `/FileStore`)
2 # Function to recursively delete a directory and its contents
3 def delete_dir(path):
4     try:
5         for file in dbutils.fs.ls(path):
6             if file.isDir():
7                 delete_dir(file.path) # Recursive call for subdirectories
8             else:
9                 dbutils.fs.rm(path, recurse=True) # Use fs.rm for files in FileStore
10        dbutils.fs.rm(path, recurse=True) # Delete the empty directory itself
11    except Exception as e:
12        print(f"Error deleting directory {path}: {e}")
13
14 # Specify directories or files to be deleted (replace with your specific paths)
15 paths_to_delete = [
16     "/FileStore/tables/clinicaltrial_2023-2.csv"
17 ]
18
```



```

18
19 # Delete each path
20 for path in paths_to_delete:
21     delete_dir(path)
22
23 # Unmount DBFS (optional, good practice for security)
24 # Since we didn't mount anything, this line is not necessary
25
26 print("Cleanup completed!")
27

```

Cleanup completed!

Command took 0.42 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 5:50:41 AM on dscdsxc

## 4. Processing and Cleaning Data:-

The file contains the delimiter /t between each line and various unnecessary commas, strips leading and trailing quotes, and additional whitespaces that I realize after seeing the file.

So to remove these I defined a UDF which will remove commas, strips leading and trailing quotes, and remove any additional whitespace. This UDF also ensure that each row has same length as the header by padding missing values with empty string.

Than I use a map function on RDD to split each line by a delimiter /t.

```

1
2 def clean_and_pad(parts):
3     # Remove commas, strip quotes and whitespace
4     cleaned_parts = [part.replace(",", "").strip().strip('\"') for part in parts]
5     # Pad the row if it has fewer elements than expected
6     if len(cleaned_parts) < 14:
7         cleaned_parts += [""] * (14 - len(cleaned_parts))
8     return cleaned_parts
9
10 processed_rdd = raw_rdd.map(lambda line: line.split("\t")).map(clean_and_pad)
11
12
13 # Filter out the header if it's the first row and matches expected headers
14 header = processed_rdd.first() # Assuming the first row is the header
15 data_rdd = processed_rdd.filter(lambda row: row != header and len(row) == 14) # Ensure all rows have exactly 14 elements
16
17 # Create DataFrame
18 df = spark.createDataFrame(data_rdd, schema=schema)
19 df.show()
20
21

```

▶ (2) Spark Jobs

**Output:-**

Command took 14.71 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 4:43:58 PM on fdscsx

Type	Id	Study Title	Acronym	Status	Conditions	Interventions	Sponsor	Collaborators	Enrollment	Funder
Type	Type	Study Design	Start	Completion						
	NCT03630471	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 Monoe...		RECRUITING	Cerebrovascular F...	OTHER: Placebo DI...	McMaster University	Alzheimer's Socie...	30.0	
OTHER	INTERVENTIONAL	Allocation: RANDO...	2023-10-25	2024-08						
	NCT00237471	Impact of Tight G...		TERMINATED	Myocardial Infarc...	DRUG: Insulin (ti...	Melbourne Health	National Health a...	40.0	
OTHER	INTERVENTIONAL	Allocation: RANDO...	2005-10	2006-05						
	NCT03820271	New Prognostic Pr...	SUPERMELD	RECRUITING	Decompensated Cir...	OTHER: SuperMELD	Assistance Publi...		500.0	
OTHER	INTERVENTIONAL	Allocation: NA In...	2020-10-01	2023-10-01						
	NCT06229171	InTake Care: Deve...	InTakeCare	NOT_YET_RECRUITING	Hypertension Trea...	OTHER: adherence ...	Istituto Auxologi...	Istituti Clinici ...	206.0	
OTHER	INTERVENTIONAL	Allocation: RANDO...	2024-10-01	2026-04-01						
	NCT02945371	Tailored Inhibito...	REV	COMPLETED	Smoking Alcohol D...	BEHAVIORAL: Perso...	University of Oregon		103.0	
OTHER	INTERVENTIONAL	Allocation: RANDO...	2014-09	2016-05						
	NCT01055171	Neuromodulation o...		COMPLETED	Alcohol Dependenc...	DRUG: Propranolol...	Medical Universit...	National Institut...	44.0	
OTHER	INTERVENTIONAL	Allocation: RANDO...	2010-01	2012-08						
	NCT01125371	Computerized Brie...		COMPLETED	Alcohol: Harmful ...	BEHAVIORAL: Compu...	Johns Hopkins Uni...	National Institut...	439.0	

Now from the output we can see that the data is now cleaned and is prepare to analyze.

## 5. Date format:-

As the file contain various date format in the Start and Completion column of the file. So this correction is achieved by using 2 UDF.

### First UDF (delete\_day\_udf):

This UDF will use a regular expression with the pattern `(r'(\b\d(4)-\d(2))- \d(2)\b')` to identify date string in the format YYYY-MM-DD and than use replacement pattern `r'\1'` to remove the day part from the date because it would be easy for me to analyze the dates by removing unnecessary day part.

### Second UDF (format\_date):

This UDF will attempt to parse the string(date values stored as a string in start and completion column) using `datetime.strptime` with the format `%Y-%m`. If the parsing is successful than this UDF create a `date_obj` and format it using `strftime` with the format `%b%Y`. So in this way this UDF standardizes the date format.

Cmd 6

```

1
2 from pyspark.sql.types import StructType, StructField, StringType, DateType
3 from pyspark.sql.functions import col, udf, lit
4 import re
5
6 def delete_day(date_str):
7     """Removes the day part from a date string (if present) using regular expressions."""
8     return re.sub(r'(\b\d(4)-\d(2))- \d(2)\b', r'\1', str(date_str))
9
10 def format_date(date_str):
11     """Formats a date string (YYYY-MM) into a more human-readable format (Month Year)."""
12     if date_str:
13         try:
14             date_obj = datetime.strptime(date_str, "%Y-%m")
15             return date_obj.strftime("%b %Y")
16         except ValueError:
17             return date_str
18     else:
19         return date_str
20

```

Command took 0.11 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 6:40:59 PM on fdscsx

Now I am going to apply UDFs in correct order so for this purpose I apply **delete\_day\_udf** to “Start” and “completion” column, and create a new column with the name Start\_Cleaned and Completion\_cleaned respectively. Then I Apply the **format\_date\_udf** to the 'Start\_Cleaned' column and 'Completion\_Cleaned' column, attempting to parse and format the date string.

```

Cmd 6
Python
1
2 # Define UDFs
3 delete_day_udf = udf(delete_day, StringType())
4 format_date_udf = udf(format_date, StringType())
5 # Apply UDFs in the correct order
6 df = df.withColumn("Start_Cleaned", delete_day_udf(col("Start")))
7 df = df.withColumn("Completion_Cleaned", delete_day_udf(col("Completion")))
8 df = df.withColumn("Start_Date", format_date_udf(col("Start_Cleaned")))
9 df.withColumn("Completion_Date", format_date_udf(col("Completion_Cleaned")))

Command took 0.11 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 6:40:59 PM on fdscsx

```

Now in this way we have corrected the date format using user defined function. But the issue is that we have now multiple columns for start and completion, so to resolve this issue I am going to delete the previous(original) start and completion column to avoid the confusion.

```

Cmd 7
Python
1
2 df = df.drop("Start")
3 df = df.drop("Completion")
4

df: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string ... 12 more fields]
Command took 0.13 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 6:30:33 PM on fdscsx

```

## 6. Temporary View:-

I also define the temporary view of above data frame so that I can view the data clearly and apply SQL queries on this temporary view. I use **createOrReplaceTempView** command for making this temporary view.

```

Cmd 8
Python
1
2 df.createOrReplaceTempView("clinical_trials")
3
4 # Show the cleaned DataFrame
5 df.show()

(1) Spark Jobs

```

**Output:-**

► (1) Spark Jobs

Id	Study Title	Acronym	Status	Conditions	Interventions	Sponsor	Collaborators	Enrollment	Funder
Type	Type	Study Design	Start_Cleaned	Completion_Cleaned					
NCT03630471	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 Monoe...		RECRUITING	Cerebrovascular F...	OTHER: Placebo DI...	McMaster University	Alzheimer's Socie...	30.0	
OTHER INTERVENTIONAL	Allocation: RANDO...		2023-10-25	2024-08					
NCT00237471	Impact of Tight G...		TERMINATED	Myocardial Infarc...	DRUG: Insulin (ti...	Melbourne Health	National Health a...	40.0	
OTHER INTERVENTIONAL	Allocation: RANDO...		2005-10	2006-05					
NCT03820271	New Prognostic Pr...	SUPERMELD	RECRUITING	Decompensated Cir...	OTHER: SuperMELD	Assistance Publiq...		500.0	
OTHER INTERVENTIONAL	Allocation: NA In...		2020-10-01	2023-10-01					
NCT06229171	InTake Care: Deve...	InTakeCare	NOT_YET_RECRUITING	Hypertension Trea...	OTHER: adherence ...	Istituto Auxologi...	Istituti Clinici ...	206.0	
OTHER INTERVENTIONAL	Allocation: RANDO...		2024-10-01	2026-04-01					
NCT02945371	Tailored Inhibito...	REV	COMPLETED	Smoking Alcohol D...	BEHAVIORAL: Perso...	University of Oregon		103.0	
OTHER INTERVENTIONAL	Allocation: RANDO...		2014-09	2016-05					
NCT01055171	Neuromodulation o...		COMPLETED	Alcohol Dependenc...	DRUG: Propranolol...	Medical Universit...	National Institut...	44.0	
OTHER INTERVENTIONAL	Allocation: RANDO...		2010-01	2012-08					
NCT01125371	Computerized Brie...		COMPLETED	Alcohol: Harmful ...	BEHAVIORAL: Comou...	Johns Hopkins Uni...	National Institut...	439.0	

Command took 4.57 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/23/2024, 6:30:43 PM on fdcsx

Now from above we can see that the data is now completely cleaned is prepared for further analyzing which I am going to do so.

## Question#1:

**The number of studies in the dataset. You must ensure that you explicitly check distinct studies.**

## Assumptions:-

The assumptions that I made before answering this question are

1. Each row in the dataset is a unique clinical trial.
2. The column name "ID" is used to represent each clinical trial.
3. The dataset is assumed to be properly cleaned and formatted.

## SQL implementation:-

### Answer:-

I used the count function to count the distinct rows from the temporary view named "Clinical\_trial" by applying the distinct keyword. This will ensures that only the rows will be counted which are distinct.

### Code:-

Cmd 9

```

1
2 SELECT DISTINCT count(*) FROM clinical_trials

```

▶ (2) Spark Jobs

Table ▾ +

	count(1)
1	483422

1 row | 17.58 seconds runtime

Refreshed 10 days ago

Command took 17.58 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/14/2024, 3:17:54 PM on gtfdivc

## Result:-

From above code we get the result that there are 483422 distinct rows which means that there are 483422 numbers of distinct studies are present in the dataset.

## DF Implementation

### Answer:-

To find the distinct number of studies I use df.distinct command to remove all duplicate rows and then I use count function to only count distinct row.

### Code:-

Cmd 3

```

1
2 df.distinct().count()

```

▶ (3) Spark Jobs

Out[15]: 483422

Command took 45.12 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/14/2024, 3:20:15 PM on gtfdivc

## Result:-

From above you can see that distinct number of studies is 483422.

## RDD implementation

To find the distinct number of studies I use lambda function which takes every row as an input and extract only the first column and then I use distinct() function to eliminate the duplicate elements and then I use count function to count the distinct elements. As in the assumption that ID is used to represent each clinical trial so I did this process on ID which will be equal to distinct number of studies.

```
Cmd 6
Python ▶ ⌵ ⌵ ⌵ ✕
1
2
3 # Assuming the first column of your RDD is what you're interested in for distinct counts
4 distinct_first_elements_count = data_rdd.map(lambda row: row[0]).distinct().count()
5
6 print(f"Number of distinct studies: {distinct_first_elements_count}")
7

▶ (1) Spark Jobs
Number of distinct studies: 483422
Command took 11.52 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/27/2024, 4:10:58 AM on fcrsz
```

## Result:-

Distinct number of study is 483422.

## Discussion on results:-

In all three implementations, regardless of whether using SQL, DataFrames, or RDDs, the analysis consistently shows that there are 483,422 distinct studies present in the dataset.

## Question#2:

**You should list all the types (as contained in the Type column) of studies in the dataset along with the frequencies of each type. These should be ordered from most frequent to least frequent.**

## Assumptions:-

The assumptions that I made before answering this question are

1. I assumed that fields in data are separated by /t delimiter.
2. The header line is considered to be the initial line of the file.
3. The dataset is assumed to be properly cleaned and formatted.
4. The type of study is contained in the "Type" column of the data set.

## Answer:

### SQL Implementation

I extract the "Type" column from the temporary view "Clinical\_trial" and utilize the count function to determine the frequency of each type. Subsequently, I employ the groupby clause to group the rows by the values in the "Type" column, followed by sorting them in descending order with the OrderBy clause

## Code:-

Cmd 10

```

1
2 SELECT clinical_trials.Type, count(*) as count FROM clinical_trials
3 GROUP BY clinical_trials.Type
4 ORDER BY count DESC

```

► (2) Spark Jobs

Table ▾ +

	Type	count
1	INTERVENTIONAL	371382
2	OBSERVATIONAL	110221
3	EXPANDED_ACCESS	928
4		891

4 rows | 18.41 seconds runtime

Refreshed 10 days ago

Command took 18.41 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/14/2024, 3:17:54 PM on gtf5vc

## Result:-

From the result set we can see that it showcases the various types of studies present in the dataset, along with their respective frequencies arranged in descending order, as outlined below.

Table ▾ +

	Type	count
1	INTERVENTIONAL	371382
2	OBSERVATIONAL	110221
3	EXPANDED_ACCESS	928
4		891

## DF Implementation

To enumerate all study types, I employ the DataFrame's **groupby()** method on the "Type" column, followed by the **count()** function to tally the occurrences of each type. Finally, I arrange the results in ascending order.

## Code:-

Cmd 4

```

1 df.groupby('Type').count().orderBy('count', ascending=False).show(4)
2

```

► (2) Spark Jobs

```

+-----+
|      Type| count|
+-----+
| INTERVENTIONAL| 371382|
| OBSERVATIONAL| 110221|
| EXPANDED_ACCESS| 928|
|              | 891|
+-----+

```

Command took 19.68 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/14/2024, 3:20:15 PM on gtf5vc

## Result:-

From above SS you can see all the type of studies with their respective count.

## RDD Implementation

In this code I extract the "Type" values from each row and removes entries with missing values(None). Then I count how many times each unique type appears and sort them by their frequency. At the end I print a table showing the type and its frequency.

```
1 # Extract the "Type" column from the RDD and filter out None values
2 type_rdd = processed_rdd.map(lambda row: row[10]).filter(lambda x: x is not None)
3
4 # Count the occurrences of each type
5 type_counts_rdd = type_rdd.map(lambda type: (type, 1)).reduceByKey(lambda a, b: a + b)
6
7 # Sort the counts in descending order
8 sorted_type_counts = type_counts_rdd.sortBy(lambda x: x[1], ascending=False)
9
10 # Collect the sorted counts
11 sorted_type_counts_list = sorted_type_counts.collect()
12
13 # Print the results
14 print("{:<30} {:<10}".format('Type', 'Frequency'))
15 for type, count in sorted_type_counts_list:
16     print("{:<30} {:<10}".format(type, count))
17
```

► (3) Spark Jobs

Type	Frequency
INTERVENTIONAL	371382
OBSERVATIONAL	110221
EXPANDED_ACCESS	928
	891
Type	1

## Result:-

```
► (3) Spark Jobs
```

Type	Frequency
INTERVENTIONAL	371382
OBSERVATIONAL	110221
EXPANDED_ACCESS	928
	891
Type	1

Command took 9.81 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 5:17:14 PM on gfvx

## Discussion on result:-

Irrespective of the approach employed (be it SQL, DataFrame, or RDD), the analysis furnishes an exhaustive inventory of the various types of studies present in the dataset, delineating their occurrences from highest to lowest frequency as

```
► (3) Spark Jobs
```

Type	Frequency
INTERVENTIONAL	371382
OBSERVATIONAL	110221
EXPANDED_ACCESS	928
	891

Command took 8.18 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/27/2024, 3:50:26 AM on fcxz



### Question#3:

The top 5 conditions (from Conditions) with their frequencies.

#### Assumptions:-

The assumptions that I made before answering this question are

1. The condition column is assumed to be the 5<sup>th</sup> column (index 4) in the dataset.
2. Conditions are extracted by splitting the “Conditions” column by the delimiter “|”.
3. The dataset is properly cleaned and formatted.

### SQL Implementation

I created a temporary view with the name “all conditions” .It has explode function which split the values in the condition column of temporary view “clinical\_trial 2023” by a delimiter “|”.

Then I query the data by selecting condition column from temporary view “all condition” and count the occurrences of each condition and then I use ‘Where’ clause to filter out any empty conditions. Then the results are grouped by the condition and sorted in descending order of count to identify the most common conditions. Then finally I limit the output to the top 5 most conditions.

#### Code:-



```
1
2 CREATE OR REPLACE TEMP VIEW all_conditions AS
3 SELECT explode(split(Conditions, '\\|')) AS condition
4 FROM clinical_trials;
5
6 SELECT TRIM(condition) AS condition, COUNT(*) as count
7 FROM all_conditions
8 WHERE condition != ''
9 GROUP BY condition
10 ORDER BY count DESC
11 LIMIT 5;
12
```

▶ (2) Spark Jobs

#### Result:-

▶ (2) Spark Jobs

Table ▾ +

	condition	count
1	Healthy	9731
2	Breast Cancer	7502
3	Obesity	6549
4	Stroke	4073
5	Hypertension	4022

5 rows | 20.05 seconds runtime

Refreshed 10 days ago

Command took 20.05 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/14/2024, 3:17:54 PM on gtf0vc

## DF Implementation

To identify the top 5 conditions along with their frequencies, I extract the "Condition" column from the DataFrame. After splitting this column by the delimiter "|", I utilize the groupby operation to aggregate the data based on the condition column. By applying the count() function, I determine the frequency of each unique condition. Finally, I sort the results in descending order to obtain the top 5 conditions.

## Code:-

```

1 from pyspark.sql.functions import explode, split, col, trim
2
3 # Define the delimiter for the conditions column
4 conditions_delimiter = {
5     "clinicaltrial_2023": "\\|",
6 }
7
8 # Replace special characters in the delimiter string
9 delimiter = conditions_delimiter.get("clinicaltrial_2023", "\\|")
10
11 # Explode the Conditions column by the delimiter "|", handle empty conditions, and clean whitespace
12 exploded_df = df.withColumn('Conditions', explode(split(trim(col('Conditions')), delimiter))) \
13     .filter(col('Conditions') != '')
14
15 # Group by Conditions and count occurrences
16 condition_counts_df = exploded_df.groupBy('Conditions').count()
17
18 # Order by count in descending order and show the top 5
19 condition_counts_df.orderBy(condition_counts_df['count'].desc()).show(5, truncate=False)
20

```

▶ (2) Spark Jobs

▶ exploded\_df: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string ... 14 more fields]

▶ condition\_counts\_df: pyspark.sql.dataframe.DataFrame = [Conditions: string, count: long]

## Result:-

```

▶ (2) Spark Jobs
▶ exploded_df: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string ... 14 more fields]
▶ condition_counts_df: pyspark.sql.dataframe.DataFrame = [Conditions: string, count: long]
+-----+-----+
|Conditions|count|
+-----+-----+
|Healthy|9731|
|Breast Cancer|7502|
|Obesity|6549|
|Stroke|4073|
|Hypertension|4022|
+-----+-----+
only showing top 5 rows

```

Command took 20.73 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 4:35:14 PM on ny Cluster

## RDD Implementation

In this code I split each condition into separate entries. Then I counts how many times each unique condition appears and then finds the top 5 most frequent conditions. At the end I print the table showing the condition with its respective frequency.

```

Cmd 6
Python ▶ ▼ - ✕

1 # Explode the Conditions column and filter out empty conditions
2 exploded_rdd = processed_rdd.flatMap(lambda row:
3     # Assuming Conditions is at index 4, check if it's a string before splitting
4     [(tuple(cond.split("|")),) for cond in (row[4] or "").split("|") if cond]) \
5     .filter(lambda x: x[0] != '')
6
7 # Count occurrences of each condition
8 condition_counts = exploded_rdd.map(lambda x: (x[0], 1)) \
9     .reduceByKey(lambda a, b: a + b)
10
11 # Sort by count in descending order and take the top 5
12 top_5_conditions = condition_counts.sortBy(lambda x: x[1], ascending=False).take(5)
13
14 # Print the top 5 conditions and their frequencies
15 for condition, count in top_5_conditions:
16     print(f"Condition: {condition}, Count: {count}")
17
▶ (3) Spark Jobs

```

Result:-

```

▶ (3) Spark Jobs

Condition: ('Healthy',), Count: 9731
Condition: ('Breast Cancer',), Count: 7502
Condition: ('Obesity',), Count: 6549
Condition: ('Stroke',), Count: 4073
Condition: ('Hypertension',), Count: 4022

Command took 13.33 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 5:31:57 PM on gfvxcx

```

## Discussion on Result:-

From all three methods RDD, data frame and SQL we got the same answer.

## Question#4:

**Find the 10 most common sponsors that are not pharmaceutical companies, along with the number of clinical trials they have sponsored. Hint: For a basic implementation, you can assume that the Parent Company column contains all possible pharmaceutical companies.**

## Answer:-

To answer this question I read the CSV file name `pharmaceutical_data.csv` into a DF and then I only extract the column name "Parent\_Company" and make its temporary view.

```
Cmd 14
Python ▶ ▼ - x

1 # Assuming your data is stored in a CSV file named "pharmaceutical_data.csv"
2 file_path = "/FileStore/tables/pharma/pharma.csv" # Update the file path accordingly
3
4 # Read the CSV file into a DataFrame
5 pharma= spark.read.option("header", "true").csv(file_path)
6
7 # Select only the "Parent_Company" column
8 parent_companies = pharma.select("Parent_Company")
9
10 # Show the first few rows of the DataFrame
11
12 parent_companies.createOrReplaceTempView("parent_companies_temp")

▶ (1) Spark Jobs
▶ pharma: pyspark.sql.dataframe.DataFrame = [Company: string, Parent_Company: string ... 32 more fields]
▶ parent_companies: pyspark.sql.dataframe.DataFrame = [Parent_Company: string]
Command took 4.33 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 5:08:37 AM on gdfcx
```

## Assumptions:-

The assumptions that I made before answering this question are

1. The pharmaceutical dataset has a column name "Parent\_company" which include all pharmaceutical companies.
2. The clinical trial data set has a column name "sponsor" which contain the names of all the sponsors of clinical trials.
3. All the data is properly cleaned and well formatted as in schema.
4. The sponsors not present in pharmaceutical companies list are considered to be non-pharmaceutical companies.

## SQL Implementation

```
Cmd 14
SQL ▶ ▼ - x

1
2 CREATE OR REPLACE TEMP VIEW non_pharma_sponsor AS SELECT Sponsor FROM clinical_trials WHERE Sponsor NOT IN (SELECT Parent_Company FROM parent_companies_temp);
3 SELECT Sponsor, count(*) as count FROM non_pharma_sponsor
4 GROUP BY Sponsor
5 ORDER BY count DESC
6 LIMIT 10

▶ (3) Spark Jobs
Table ▼ +
```

I created a temporary view with the name non\_pharma\_sponsor and select the Sponsor column from another temp view "Clinical\_trial" using the select command and then use the filter to filter only those Sponsors which are not present in temporary view

Parent\_Company using the Where clause. Because our aim is to find only non pharmaceutical companies.

Then I count the sponsor column from the temporary view non\_pharma\_sponsor.

Then I groupby the column with sponsor name.

Then sort the count column in descending order and limit the column to 10.

## Result:-

▶ (3) Spark Jobs

Table ▾ +

	Sponsor	count
1	National Cancer Institute (NCI)	3410
2	Assiut University	3335
3	Cairo University	3023
4	Assistance Publique - Hôpitaux de Paris	2951
5	Mayo Clinic	2766
6	M.D. Anderson Cancer Center	2702
7	Novartis Pharmaceuticals	2393
8	National Institute of Allergy and Infectious Diseases (NIAID)	2340
9	Massachusetts General Hospital	2263
10	National Taiwan University Hospital	2181

10 rows | 20.09 seconds runtime

Refreshed 12 days ago

Command took 20.09 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/24/2024, 3:17:54 PM on gtfdev

## DF Implementation

Now to find the non-pharmaceutical companies I write a code to find the most frequent sponsors in the clinical\_trial 2023 data set and the filter out those companies which are present in pharmaceutical companies data.

```
Cmd 12
```

```
1 pharma_list = pharma.select("Parent_Company").rdd.flatMap(lambda x: x).collect()
2 dfs = df.select("Sponsor")
3
4 non_pharma_sponsors = dfs.groupBy("Sponsor").count().orderBy("count", ascending=False).filter(~dfs.Sponsor.isin(pharma_list)).show(10)
```

▶ (3) Spark Jobs

▶ dfs: pyspark.sql.dataframe.DataFrame = [Sponsor: string]

```
+-----+-----+
|      Sponsor|count|
+-----+-----+
|National Cancer I...| 3410|
|  Assiut University| 3335|
|  Cairo University| 3023|
|Assistance Publiq...| 2951|
|      Mayo Clinic| 2766|
|M.D. Anderson Can...| 2702|
|Novartis Pharmace...| 2393|
|National Institut...| 2340|
|Massachusetts Gen...| 2263|
|National Taiwan U...| 2181|
+-----+-----+
only showing top 10 rows
```

Command took 31.29 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 4:40:50 PM on My Cluster

## RDD Implementation

To find the non-pharmaceutical sponsor companies, first of all I clean the parent-company column in pharma.csv. Then I apply the logic using rdds that will use clinical trial data set to fetch all the sponsor companies and then subtract those companies which are present in pharma.csv file. In this way we are only left with non-pharma sponsor companies.

```

1 from pyspark.sql import SparkSession
2
3 # Initialize Spark session
4 spark = SparkSession.builder.appName("RDD Example").getOrCreate()
5
6 # Define file path
7 file_path = "/FileStore/tables/pharma/pharma.csv"
8
9 # Load the CSV file as an RDD
10 rdd = spark.sparkContext.textFile(file_path)
11
12 # Extract the header
13 header = rdd.first()
14 |
15 # Split each row by comma
16 rdd = rdd.map(lambda row: row.split(','))
17
18 # Extract the Parent_Company column
19 parent_companies_rdd = rdd.map(lambda row: (row[1],)) # Assuming Parent_Company is at index 1
20
21 # Show the first few rows of the RDD
22 parent_companies_rdd.take(5) # You can change 5 to any number to view more or fewer rows
23

```

▶ (2) Spark Jobs

Out[9]: (('Parent Company', ...)

```

1 ct_sponsor_col_index = processed_rdd.first().index('Sponsor')
2
3 parent_pharm_comp = rdd.map(lambda x: x[1].replace(' ', ''))
4
5 processed_rdd.map(lambda x: x[ct_sponsor_col_index]).filter(lambda row: row != 'Sponsor').subtract(rdd.map(lambda x: x[1].replace(' ', ''))).map(lambda x: (x, 1)).reduceByKey(lambda x, y: x + y).sortBy(lambda x: x[1], ascending=False).take(10)

```

▶ (4) Spark Jobs

Out[9]: [('National Cancer Institute (NCI)', 3410),  
('Assiut University', 3335),  
('Cairo University', 3023),  
('Assistance Publique - Hôpitaux de Paris', 2951),  
('Mayo Clinic', 2766),  
('M.D. Anderson Cancer Center', 2702),  
('Novartis Pharmaceuticals', 2393),  
('National Institute of Allergy and Infectious Diseases (NIAID)', 2340),  
('Massachusetts General Hospital', 2263),  
('National Taiwan University Hospital', 2181)]

Command took 16.70 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/27/2024, 3:55:55 AM on fci32z

## Discussion of result:-

I got the same answer by applying all three implementations (SQL, DataFrame, and RDD) effectively, meeting the requirements of the question.

### Question # 05

Plot number of completed studies for each month in 2023. You need to include your visualization as well as a table of all the values you have plotted for each month.

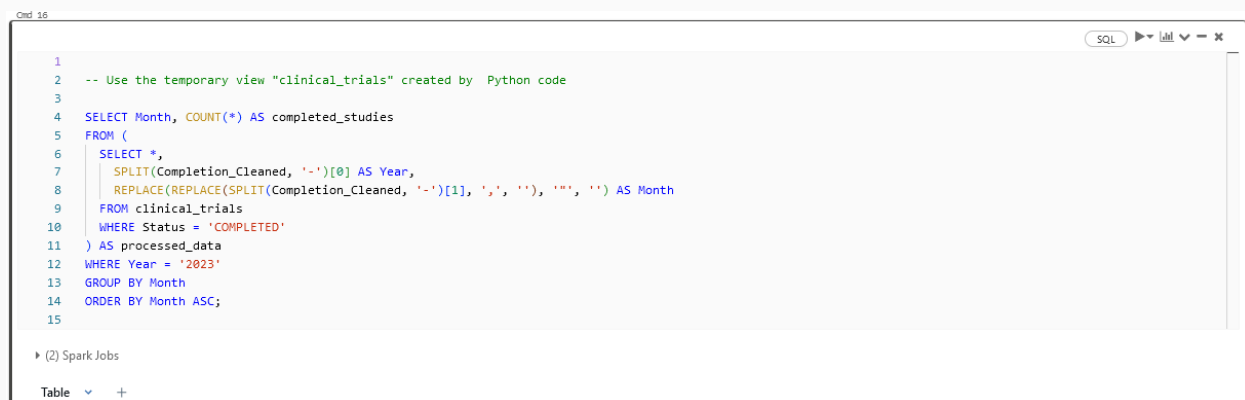
#### Assumption:-

1. Completion dates of the clinical trial data set are stored in the "Completion\_Cleaned" column.
2. The format of completion dates is YY-MM.
3. The status of each clinical trial is stored in "Status" column where "COMPLETED" indicates that the trial has been completed.
4. The data is properly cleaned and well formatted as in schema.

#### Answer:-

#### SQL Implementation

First of all I clean up the completion\_cleaned date format and extracts the month only, then I only consider those trials with completed status and then group the data by months and count the number of completed trials for each month. I apply the filter using the where clause to only filter the result with year "2023". Then at the end sort the results with months in chronological order.



```
1
2 -- Use the temporary view "clinical_trials" created by Python code
3
4 SELECT Month, COUNT(*) AS completed_studies
5 FROM (
6     SELECT *,
7         SPLIT(Completion_Cleaned, '-') [0] AS Year,
8         REPLACE(REPLACE(SPLIT(Completion_Cleaned, '--')[1], ','), '') AS Month
9     FROM clinical_trials
10    WHERE Status = 'COMPLETED'
11 ) AS processed_data
12 WHERE Year = '2023'
13 GROUP BY Month
14 ORDER BY Month ASC;
15
```

▶ (2) Spark Jobs

Table ▾ +

#### Result:-

Table		
	Month	completed_studies
1	01	1494
2	02	1272
3	03	1552
4	04	1324
5	05	1415
6	06	1619
7	07	1360
8	08	1230
9	09	1152
10	10	1058
11	11	909
12	12	1082

12 rows | 26.32 seconds runtime

Refreshed 12 days ago

Command took 26.32 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/14/2024, 3:17:54 PM on gtf4vc

## Plotting:-

For plotting I convert the cell from SQL to python and read the data from temporary view "clinical\_trial" into a spark data frame.

```

1 %python
2 from pyspark.sql import SparkSession
3
4 spark = SparkSession.builder.appName("ClinicalTrialsPlot").getOrCreate()
5
6 # Assuming "clinical_trials" is the temporary view name created by your Spark code
7
8 # Read data from the temporary view into a Spark DataFrame
9 df = spark.sql("SELECT * FROM clinical_trials")
10

```

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

Command took 0.89 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 5:18:09 AM on gdfcx

Then I clean the completion\_cleaned column and apply filter where status is completed and year is 2023 and order it and group it by months.

```

1 %python
2 from pyspark.sql.functions import split, regexp_replace
3 import matplotlib.pyplot as plt
4
5 # Ensure columns are properly formatted
6 for col in df.columns:
7     df = df.withColumnRenamed(col, col.strip(",").strip("'"))
8
9 # Extract data for completed clinical trials in 2023
10 completed_cd = df.withColumn('Year', split('Completion_Cleaned', "-")[0]) \
11     .withColumn('Month', split('Completion_Cleaned', "-")[1]) \
12     .withColumn('Month', regexp_replace("Month", ",", "")) \
13     .withColumn('Month', regexp_replace("Month", "'", "")) \
14     .filter(df.Status.isin(["COMPLETED"])) \
15     .filter(df.Completion_Cleaned.startswith("2023")) \
16     .select("Month", "Year", "Status")
17
18
19 completed_cd.filter(completed_cd.Year.isin(["2023"])).groupBy("Month").count().orderBy("Month", ascending=True).show()
20

```

(2) Spark Jobs

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

Output:-



```
▶ (2) Spark Jobs
▶ df: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string ... 12 more fields]
▶ completed_cd: pyspark.sql.dataframe.DataFrame = [Month: string, Year: string ... 1 more field]
```

```
+-----+
|Month|count|
+-----+
| 01 | 1494 |
| 02 | 1272 |
| 03 | 1552 |
| 04 | 1324 |
| 05 | 1415 |
| 06 | 1619 |
| 07 | 1360 |
| 08 | 1230 |
| 09 | 1152 |
| 10 | 1058 |
| 11 | 909 |
| 12 | 1082 |
+-----+
```

Command took 23.02 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 5:24:48 AM on gdfcx

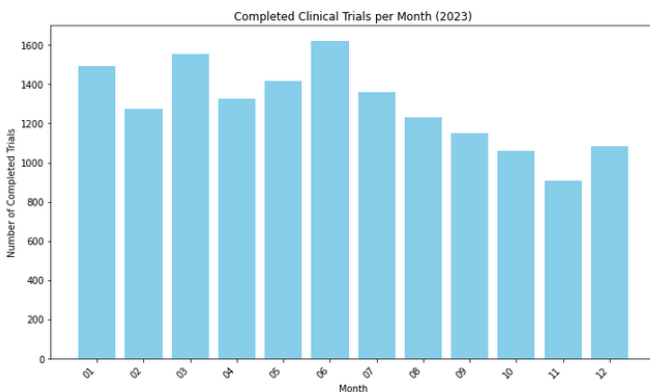
Then I use matplotlib library to plot it.

```
1 %python
2 # Import necessary libraries for plotting
3 from pyspark.sql.functions import col
4 import matplotlib.pyplot as plt
5
6
7 # Extract months and counts into separate lists
8 months = [row["Month"] for row in monthly_counts]
9 counts = [row["count"] for row in monthly_counts]
10
11 # Create the plot
12 plt.figure(figsize=(10, 6)) # Adjust figure size as desired
13 plt.bar(months, counts, color='skyblue')
14 plt.xlabel("Month")
15 plt.ylabel("Number of Completed Trials")
16 plt.title("Completed Clinical Trials per Month (2023)")
17 plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
18 plt.tight_layout()
19
20 # Display the plot
21 plt.show()
22
```

rm/firefox/7utm medium-firefox-dec dvc-toolhar8utm campaign-new-users8utm content--global

Output:-

22



Command took 0.41 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 5:43:14 AM on gdfcx

DF Implementation

I clean the completion date column using the format YY-MM. Then I filter only those clinical trials whose status is completed in 2023 and group it by “month” and order it by “month” and the plot it using matplotlib library.

```

1 from pyspark.sql.functions import split, regexp_replace
2 import matplotlib.pyplot as plt
3
4 # Ensure columns are properly formatted
5 for col in df.columns:
6     df = df.withColumnRenamed(col, col.strip(",").strip(' '))
7
8 # Extract data for completed clinical trials in 2023
9 completed_cd = df.withColumn('Year', split('Completion', "-")[0]) \
10     .withColumn('Month', split('Completion', "-")[1]) \
11     .withColumn('Month', regexp_replace("Month", "-", "")) \
12     .withColumn('Month', regexp_replace("Month", "", "")) \
13     .filter(df.Status.isin(["COMPLETED"])) \
14     .filter(df.Completion.startswith("2023")) \
15     .select("Month", "Year", "Status")
16
17
18 completed_cd.filter(completed_cd.Year.isin(["2023"])).groupBy("Month").count().orderBy("Month", ascending=True).show()
19

```

▶ (2) Spark Jobs

df: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string ... 14 more fields]

Output:-

```

▶ (2) Spark Jobs
df: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string ... 14 more fields]
completed_cd: pyspark.sql.dataframe.DataFrame = [Month: string, Year: string ... 1 more field]
+-----+
|Month|count|
+-----+
| 01 | 1494 |
| 02 | 1272 |
| 03 | 1552 |
| 04 | 1324 |
| 05 | 1415 |
| 06 | 1619 |
| 07 | 1360 |
| 08 | 1230 |
| 09 | 1152 |
| 10 | 1058 |
| 11 | 909 |
| 12 | 1082 |
+-----+

```

Command took 21.29 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 4:54:51 PM on my Cluster

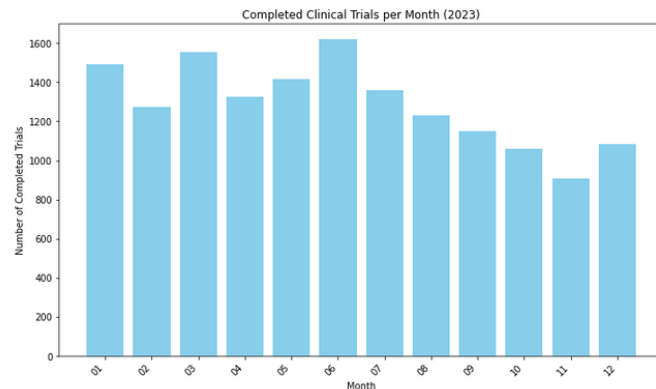
Now I plot the above table using matplotlib library.

```

1 from pyspark.sql.functions import col
2 import matplotlib.pyplot as plt
3
4 # Get monthly counts for completed trials in 2023
5 monthly_counts = (
6     completed_cd.filter(completed_cd.Year.isin(["2023"]))
7     .groupBy("Month")
8     .count()
9     .orderBy("Month", ascending=True)
10    .collect()
11 )
12 months = [row["Month"] for row in monthly_counts]
13 counts = [row["count"] for row in monthly_counts]
14
15 # Create the plot
16 plt.figure(figsize=(10, 6)) # Adjust figure size as desired
17 plt.bar(months, counts, color='skyblue')
18 plt.xlabel("Month")
19 plt.ylabel("Number of Completed Trials")
20 plt.title("Completed Clinical Trials per Month (2023)")
21 plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
22 plt.tight_layout()
23
24 # Display the plot
25 plt.show()
26

```

Output:-



Command took 37.60 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/26/2024, 5:05:43 PM on My Cluster

## RDD Implementation

I parse the month, year, and status fields from each record, handling potential errors. Then, I filter for completed studies in 2023, tallying the count of completed studies for each month. Next, I organize the data into a DataFrame with month and count columns, sorting it by month, and present it as a tabular display. Finally, utilizing the matplotlib library, I generate a bar chart illustrating the distribution of completed studies by month.

```
1 from pyspark.sql import Row
2 # Extract necessary fields assuming indexes: Status - 3, Completion - 13
3 # Completion date expected format "YYYY-MM"
4 def extract_fields(row):
5     try:
6         completion = row[13].split("-")
7         year = completion[0] if len(completion) > 0 else None
8         month = completion[1] if len(completion) > 1 else None
9         status = row[3]
10        return (month, year, status)
11    except IndexError:
12        # Handle the error: you can choose to return None or a default value
13        return (None, None, None)
14
15 extracted_rdd = data_rdd.map(extract_fields).filter(lambda x: None not in x)
16
17 # Filter completed studies and specific year
18 completed_rdd = extracted_rdd.filter(lambda x: x[2] == "COMPLETED" and x[1] == "2023")
```

Command took 0.12 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 6:06:11 PM on gfvx

```
Cmd 10

1
2
3 # Count by month
4 month_counts = completed_rdd.map(lambda x: (x[0], 1)).reduceByKey(lambda a, b: a + b)
5
6 # Convert the results into a DataFrame
7 results_rdd = month_counts.map(lambda x: Row(month=int(x[0]), count=x[1])) if month_counts.first()[0].isdigit() \
8 | else month_counts.map(lambda x: Row(month=x[0], count=x[1]))
9 results_df = spark.createDataFrame(results_rdd)
10 sorted_results_df = results_df.orderBy("month")
11
12 # Show the DataFrame as a table
13 sorted_results_df.show()
14
15
16

▶ (3) Spark Jobs

▶ results_df: pyspark.sql.dataframe.DataFrame = [month: long, count: long]
```

## Output:-

```
▶ (3) spark jobs

▶ results_df: pyspark.sql.dataframe.DataFrame = [month: long, count: long]
▶ sorted_results_df: pyspark.sql.dataframe.DataFrame = [month: long, count: long]

+-----+
|month|count|
+-----+
|  1 | 1494 |
|  2 | 1272 |
|  3 | 1552 |
|  4 | 1324 |
|  5 | 1415 |
|  6 | 1619 |
|  7 | 1360 |
|  8 | 1230 |
|  9 | 1152 |
| 10 | 1058 |
| 11 |  909 |
| 12 | 1082 |
+-----+

Command took 10.13 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 6:07:30 PM on gfvex
```

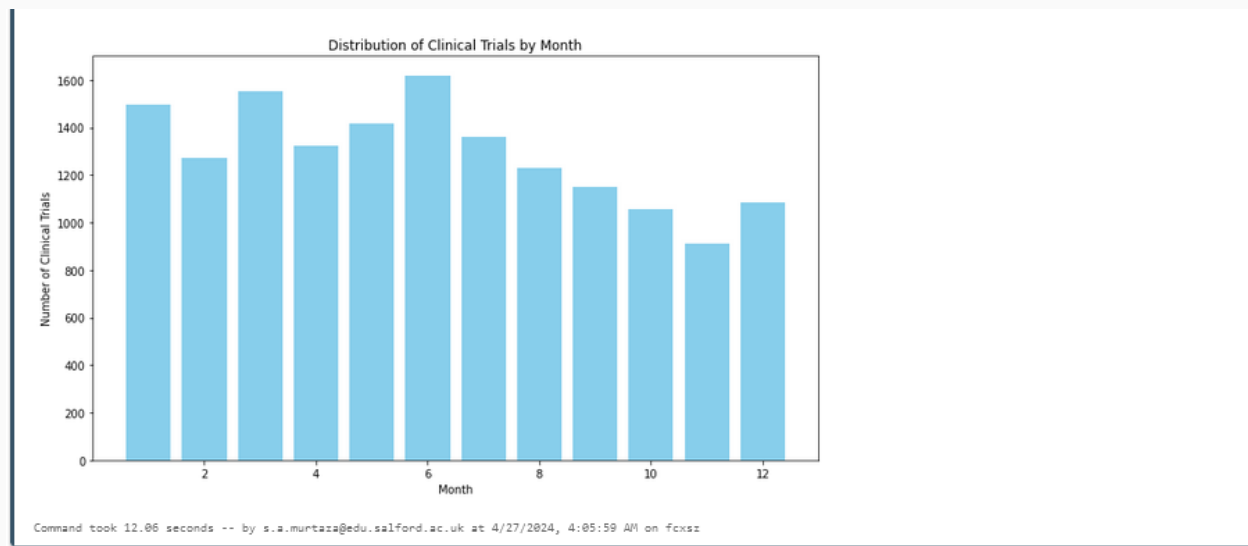
## Plotting

```
Cmd 11

1 # Plotting (using matplotlib, assumed to be imported)
2 import matplotlib.pyplot as plt
3
4 # Extract month and count from DataFrame
5 months = [row["month"] for row in sorted_results_df.collect()] # Collect results as list of dictionaries
6 counts = [row["count"] for row in sorted_results_df.collect()]
7
8 plt.figure(figsize=(10, 6))
9 plt.bar(months, counts, color='skyblue')
10 plt.xlabel("Month")
11 plt.ylabel("Number of Clinical Trials")
12 plt.title("Distribution of Clinical Trials by Month")
13 plt.xticks(rotation=0)
14 plt.tight_layout()
15 plt.show()

▶ (4) Spark Jobs
```

Output:-



Discussion of Result:-

From above in all three implementations we got the same answer.

### Extra features

**Question:-**

**Write a general and reusable code for example for clinicaltrial\_2020 and clinicaltrial\_2021 datasets.**

**Answer:-**

I write a general code for dataset clinicaltrial\_2020 and clinical\_trial\_2021. In this code I read the data and make its RDD named as raw\_rdd. I defined the schema which determine, fields of the data frame along with their data types. Then I define a function which split the data by the delimiter "/" and pads it with the empty string if the elements are lesser than 9(because 9 elements are in the defined schema) . Then I map that function with raw\_rdd. Then I create a data frame using rdd and specified schema.

```
Cmd 2
1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import StructType, StructField, StringType
3 import zipfile
4
5 # Initialize Spark session
6 spark = SparkSession.builder.appName("Data Cleaning").getOrCreate()
```

Command took 0.08 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/28/2024, 5:18:07 PM on gfdvdfvc

Cmd 3

```

1 # Define the schema
2 schema = StructType([
3     StructField("Id", StringType(), True),
4     StructField("Sponsor", StringType(), True),
5     StructField("Status", StringType(), True),
6     StructField("Start", StringType(), True),
7     StructField("Completion", StringType(), True),
8     StructField("Type", StringType(), True),
9     StructField("Submission", StringType(), True),
10    StructField("Conditions", StringType(), True),
11    StructField("Interventions", StringType(), True),
12 ])
13
14 # Load data
15 # Load data
16 file_path = "/FileStore/tables/clinicaltrial_2020"
17 raw_rdd = spark.sparkContext.textFile(file_path)
18

```

Command took 0.12 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/28/2024, 5:10:48 PM on gfdvdfvc

Cmd 4

```

1 # Process and clean data
2 def clean_and_pad(line):
3     # Split line by pipe character
4     parts = line.split("|")
5     # Pad the row if it has fewer elements than expected
6     if len(parts) < 9:
7         parts += [""] * (9 - len(parts))
8     return parts
9
10 processed_rdd = raw_rdd.map(clean_and_pad)
11 # Filter out the header if it's the first row and matches expected headers
12 header = processed_rdd.first() # Assuming the first row is the header
13 data_rdd = processed_rdd.filter(lambda row: row != header and len(row) == 9) # Ensure all rows have exactly 9 elements
14
15 # Create DataFrame
16 df = spark.createDataFrame(data_rdd, schema=schema)
17 # Show the cleaned DataFrame
18 df.show()
19

```

▶ (2) Spark Jobs

▶ (2) Spark Jobs

▶ df: pyspark.sql.dataframe.DataFrame = [Id: string, Sponsor: string ... 7 more fields]

Id	Sponsor	Status	Start	Completion	Type	Submission	Conditions	Interventions
NCT02758028	The University of...	Recruiting	Aug 2005	Nov 2021	Interventional	Apr 2016		
NCT02751957	Duke University	Completed	Jul 2016	Jul 2020	Interventional	Apr 2016	Autistic Disorder...	
NCT02758483	Universidade Fede...	Completed	Mar 2017	Jan 2018	Interventional	Apr 2016	Diabetes Mellitus	
NCT02759848	Istanbul Medeniye...	Completed	Jan 2012	Dec 2014	Observational	May 2016	Tuberculosis,Lung...	
NCT02758860	University of Rom...	Active, not recru...	Jun 2016	Sep 2020	Observational [Pa...	Apr 2016	Diverticular Dise...	
NCT02757209	Consorzio Futuro ...	Completed	Apr 2016	Jan 2018	Interventional	Apr 2016	Asthma Fluticasone,Xhanc...	
NCT02752438	Ankara University	Unknown status	May 2016	Jul 2017	Observational [Pa...	Apr 2016	Hypoventilation	
NCT02753543	Ruijin Hospital	Unknown status	Nov 2015	Nov 2019	Interventional	Apr 2016	Lymphoma	
NCT02757508	Washington Univer...	Completed	Mar 2016	Jul 2017	Interventional	Apr 2016		Vitamins
NCT02753530	Orphazyme	Completed	Aug 2017	Jan 2021	Interventional	Apr 2016	Myositis	
NCT02754817	Novo Nordisk A/S	Completed	Apr 2016	Oct 2016	Observational	Apr 2016	Diabetes Mellitus Liraglutide,Xultophy	
NCT02759276	Daniel Alexandre ...	Completed	May 2015	Dec 2015	Observational	Apr 2016	Hypertension	
NCT02750956	Bulent Ecevit Uni...	Completed	Jun 2015	Mar 2016	Observational	Apr 2016	Periodontal Diseases	
NCT02752113	Institut für Phar...	Completed	Apr 2016	May 2019	Interventional	Apr 2016	Diabetes Mellitus Metformin,Empagli...	
NCT02752698	The Third Xiangya...	Active, not recru...	Jan 2015	Dec 2021	Interventional	Jun 2015	Appendicitis,Stom...	
NCT02755779	Tel Aviv Medical ...	Unknown status	Jun 2016	Jun 2017	Observational	Apr 2016		
NCT02750384	Medicines for Mal...	Terminated	May 2016	Jul 2016	Interventional	Apr 2016		
NCT02754609	James Cook Univer...	Completed	Sen 2016	Oct 2019	Interventional	Apr 2016	Honkwanrm_Infectio...	

Command took 5.35 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/28/2024, 5:14:59 PM on gfdvdfvc

## Additional Analysis 1: Top Sponsors by Enrollment (SQL):

In this sql query I am going to group the data by sponsor column and sum the enrollment values and then order the results in ascending order and shows the only top 10 sponsors.

Cmd 2

```

1
2 SELECT sponsor, SUM(CAST(enrollment AS INT)) AS total_enrollment
3 FROM clinical_trials
4 GROUP BY sponsor
5 ORDER BY total_enrollment DESC
6 LIMIT 10;
7

```

▶ (2) Spark Jobs

▶ `_sqldf: pyspark.sql.dataframe.DataFrame = [sponsor: string, total_enrollment: long]`

Table ▾ +

	sponsor	total_enrollment
1	AstraZeneca	288811380
2	KU Leuven	200238787
3	Center for International Blood and Marrow Transplant Research	200192109
4	RWTH Aachen University	100266992
5	Centre d'études et d'expertise sur les risques" l'environnement" la mobilité et l'aménagement	80000000
6	University Hospital" Caen	76483961
7	Taipei Medical University WanFang Hospital	67170653
8	ModernaTX" Inc.	61362977
9	Assistance Publique - Hôpitaux de Paris	43892196
10	University of Pennsylvania	40923451

10 rows | 21.82 seconds runtime Refreshed 7 minutes ago

Command took 21.82 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/28/2024, 5:40:24 PM on gfdvdfvc

## Additional Analysis 2: Find the number of trials initiated by each collaborator.(using df)

### Answer:-

In this code I group by the data frame by the collaborator column and count the occurrence of each unique collaborator in descending order and limit the output to 10 to see only top 10 trial by collaborator.

```

1 from pyspark.sql.functions import desc
2
3 collaborator_count = df.groupBy("Collaborators").count().orderBy(desc("count")).limit(10)
4 collaborator_count.show()
5

```

▶ (2) Spark Jobs

collaborator\_count: pyspark.sql.dataframe.DataFrame = [Collaborators: string, count: long]

Collaborators	count
National Cancer I...	7918
National Heart" L...	1689
National Institut...	1590
National Institut...	1394
Merck Sharp & Doh...	1334
National Institut...	1119
Pfizer	940
National Institut...	872
GlaxoSmithKline	805

Command took 20.88 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 5:02:03 AM on dscdsxc

### Additional Analysis 3: Find the number of clinical trials per year of start date (using RDD)

**Answer:-**

In this code I calculate the number of clinical trials per year based on the start date. I maps each trial's start year with a value of 1, then aggregates the counts for each year using a reduce operation. Finally, I collect and returns the result as a list of tuples, where each tuple represents a year and the count of trials starting in that year.

```

1 # additional analysis find the number of clinical trials per year of start date
2 trials_per_year_rdd = data_rdd.map(lambda row: (row[12].split("-")[0], 1)).reduceByKey(lambda a, b: a + b)
3 trials_per_year_rdd.collect()
4

```

▶ (1) Spark Jobs

Out[5]: [(['2023', 33676),  
 (['2011', 18321),  
 (['2006', 11314),  
 (['', 5146),  
 (['2022', 34460),  
 (['2008', 14924),  
 (['1999', 1730),  
 (['2001', 2902),  
 (['1996', 647),  
 (['1984', 44),  
 (['1993', 304),  
 (['1992', 224),  
 (['1979', 16),  
 (['2020', 29),  
 (['1977', 24),  
 (['1988', 95),  
 (['1987', 63),  
 (['1989', 20),  
 (['1976', 19),  
 (['2029', 1),  
 (['1945', 1)]]

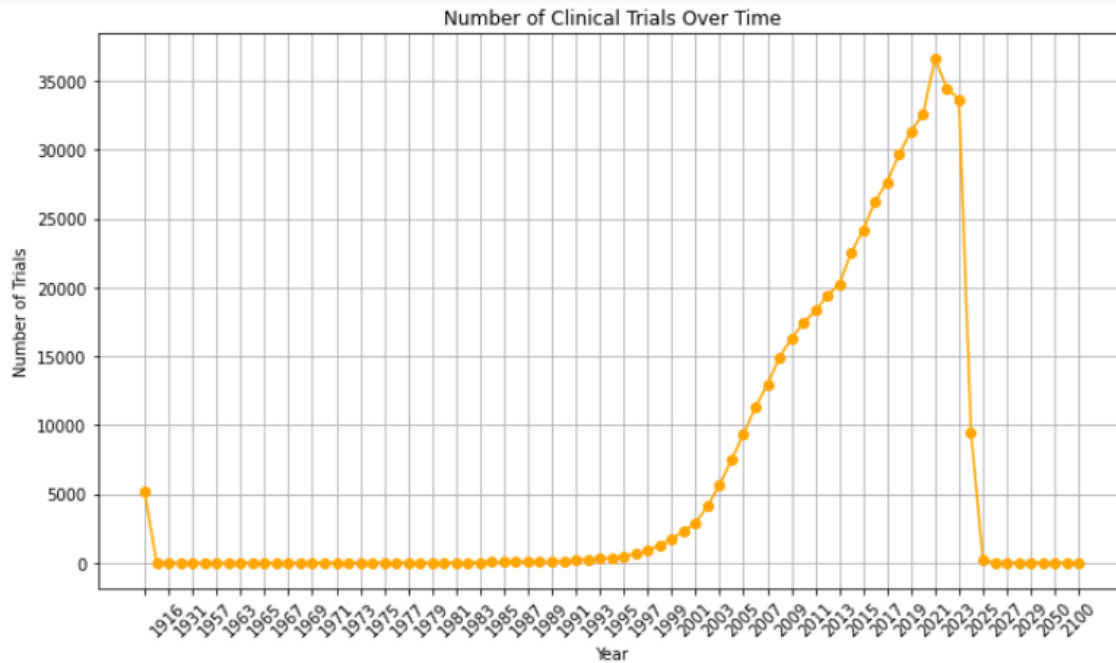
Command took 10.34 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 5:03:33 AM on dscdsxc

**Additional analysis:- Creation of additional visualizations presenting useful information based on your own exploration**

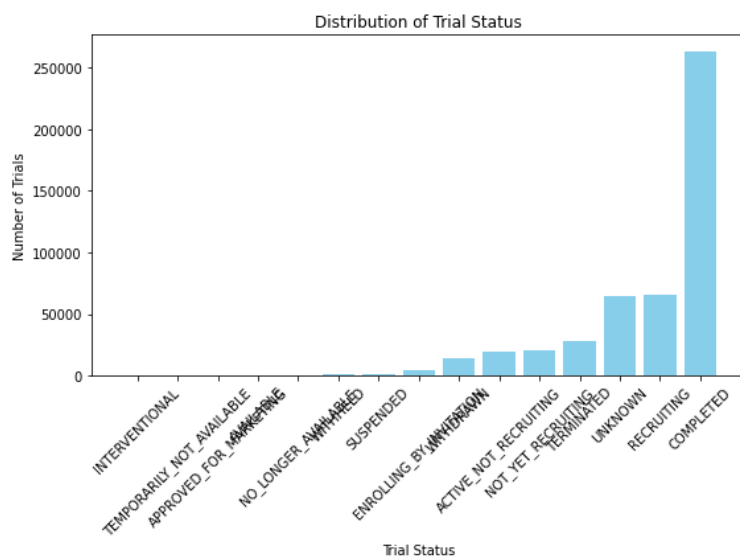
**Answer:-**

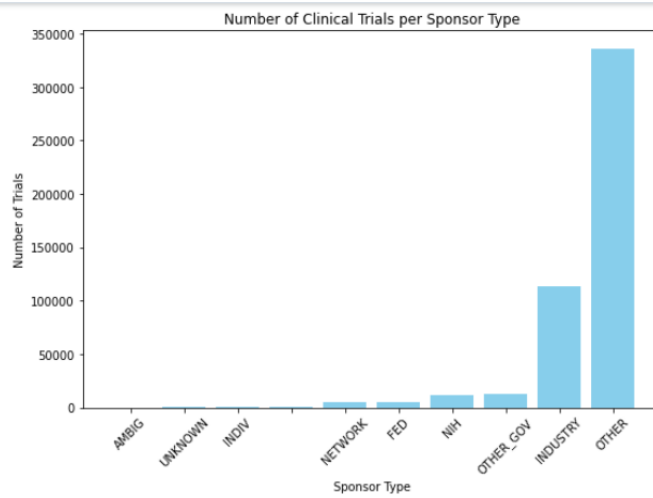
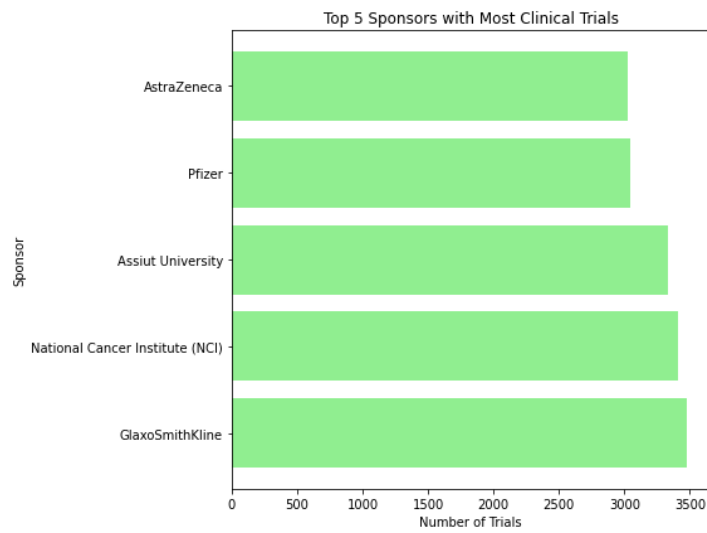
Here are some of my additional visualizations based on my own exploration





Command took 0.64 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 5:39:16 AM on dscdsc



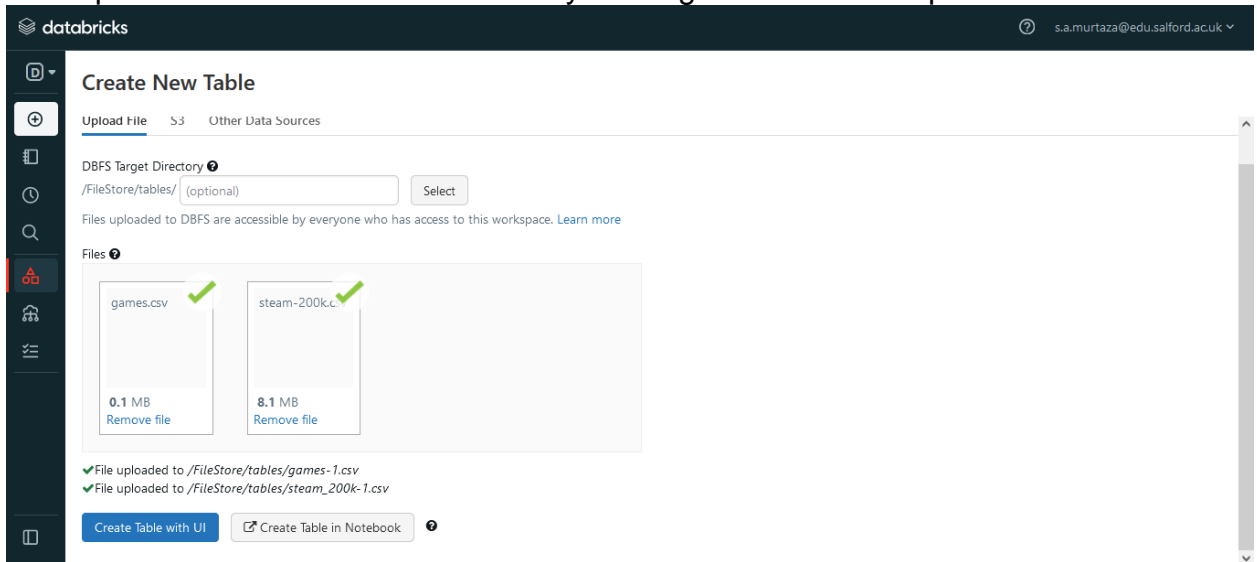


Command took 34.95 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/29/2024, 6:41:09 AM on dscdsxc

## Assignment Part # 02

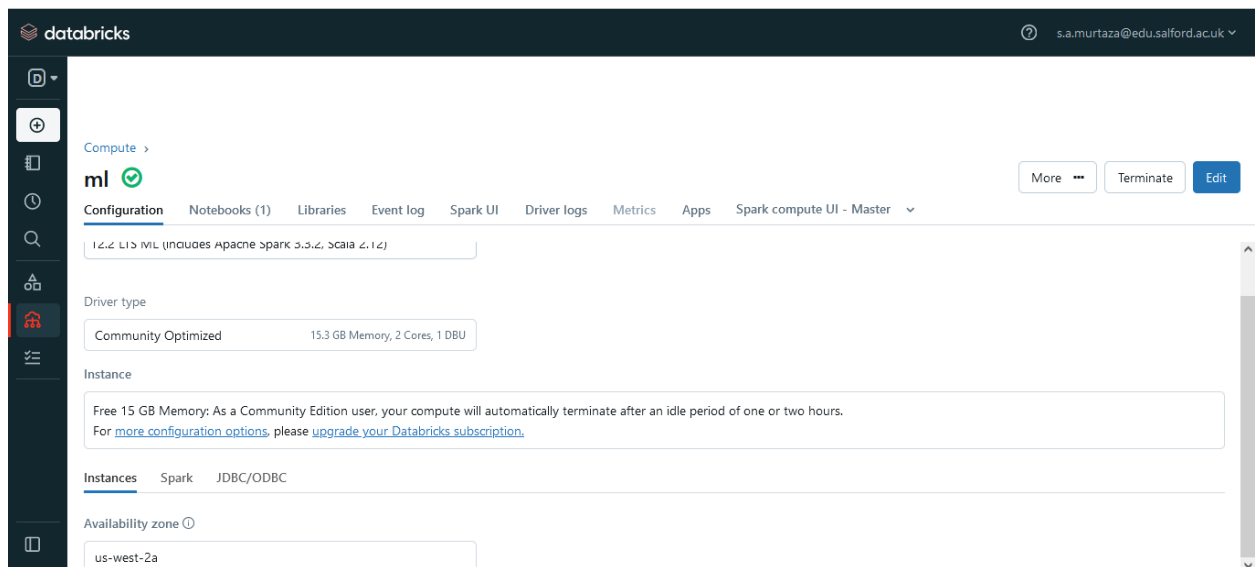
### 1. Description of the necessary configuration or prerequisites needed to successfully accomplish this task:-

I logged into data bricks community edition account. Then I download the data files which were given for the analysis in this assignment and upload them into data bricks plat form from the home screen by clicking on browse file option.



After uploading the files, I transitioned from the Data Science & Engineering role to the Machine Learning role.

Then I create a cluster using ml runtime by clicking on create compute button.



After creating the cluster, I click on “create” and select “Notebook”. In this way a python notebook appears to me. I name this notebook and set the default language as python.

## 2. Preparing data for ingestion into a Spark DataFrame, including any preliminary examination or visualization performed to understand the data structure and characteristics before initiating model training.

Then I use Ls command to check whether the file is uploaded successfully or not

Cmd 1

```

1 %fs
2 ls /FileStore/tables/steam_200k.csv
3

```

Table

	path	name	size	modificationTime
1	dbfs:/FileStore/tables/steam_200k.csv	steam_200k.csv	8059447	1713225488000

1 row | 20.67 seconds runtime

Refreshed 1 hour ago

Command took 20.67 seconds -- by s.a.murtaza@edu.salford.ac.uk at 5/1/2024, 5:05:53 AM on fxdcd

From above we can see that the required file is uploaded successfully and is ready to use for analysis purposes.

Now I am going to import the ML flow library and enable autologging.

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

Command took 0.17 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 6:13:35 AM on m1

When I looked into the file I feel its better to define my own schema with my own concept and do further analysis. So I defined schema using StructType and StructFeild. StructFeild is used to define the column names of the data frame while StructType is used to define the data type of the column. I define “member ID” as double, “name of game” as string, “member behaviour” as string and “Status” as double type. I read the data file into spark dataframe using spark.read.csv()

```
1 from pyspark.sql.types import *
2 inputPath = "/FileStore/tables/steam_200k.csv"
3 myschema = StructType([
4     StructField("member Id", DoubleType(), True),
5     StructField("Name of the Game", StringType(), True),
6     StructField("member behaviour", StringType(), True),
7     StructField("Status", DoubleType(), True),
8 ])
9 # Read the data as a DataFrame with the specified schema
10 OccupancyDF = (
11     spark.read
12     .option("header", True)
13     .schema(myschema)
14     .csv(inputPath)
15 )
16
17 display(OccupancyDF)
18
```

► (1) Spark Jobs

▼ (1) Spark Jobs

▼ OccupancyDF: pyspark.sql.dataframe.DataFrame

```
member Id: double
Name of the Game: string
member behaviour: string
Status: double
```

Table Visualization 1 +

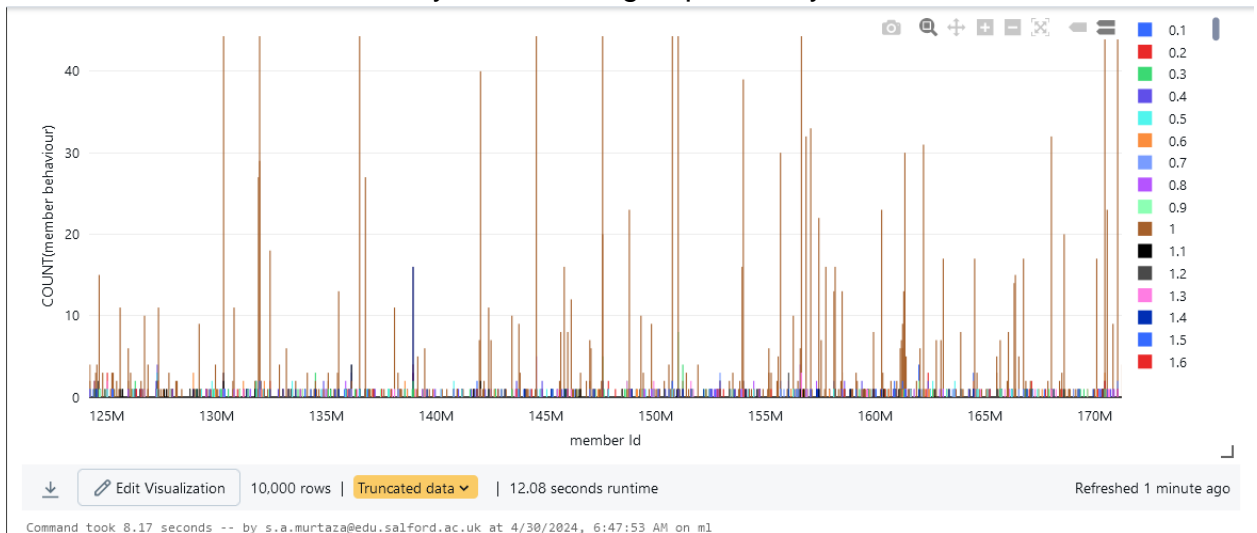
	member Id	Name of the Game	member behaviour	Status
1	151603712	The Elder Scrolls V Skyrim	play	273
2	151603712	Fallout 4	purchase	1
3	151603712	Fallout 4	play	87
4	151603712	Spore	purchase	1
5	151603712	Spore	play	14.9
6	151603712	Fallout New Vegas	purchase	1
7	151603712	Fallout New Vegas	play	12.1

10,000 rows | Truncated data | 1.28 seconds runtime

Refreshed 30 minutes ago

Command took 1.28 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 6:13:35 AM on m1

For better understanding of data I create visualization by setting member ID on x-axis and member behavior on y-axes and group them by status.



### 3. Data preprocessing conducted to prepare the data for model training.

For the training of data I have to apply ALS matrix factorization, so I need to have integer ID values for both users and items but the dataset does not contain the ID for the games. So first of all I will create a string Indexer to generate integer IDs for the game names and then fit that string indexer to the data and transform the data frame

```

Cmd 4
Python
1 from pyspark.ml.feature import StringIndexer
2
3 # Create a StringIndexer to generate integer IDs for the game names
4 indexer = StringIndexer(inputCol="Name of the Game", outputCol="gameId")
5
6 # Fit the StringIndexer to the data and transform the DataFrame
7 indexed_df = indexer.fit(OccupancyDF).transform(OccupancyDF)
8
9 # Show the DataFrame with the generated integer IDs
10 indexed_df.show()
11
(3) Spark Jobs
(1) MLflow run
  Logged 1 run to an experiment in MLflow. Learn more

```

member Id	Name of the Game	member behaviour	Status	gameId
151603712	The Elder Scrolls...	play	273.0	8.0
151603712	Fallout 4	purchase	1.0	100.0
151603712	Fallout 4	play	87.0	100.0
151603712	Spore	purchase	1.0	332.0
151603712	Spore	play	14.9	332.0
151603712	Fallout New Vegas	purchase	1.0	29.0
151603712	Fallout New Vegas	play	12.1	29.0
151603712	Left 4 Dead 2	purchase	1.0	4.0
151603712	Left 4 Dead 2	play	8.9	4.0
151603712	HuniePop	purchase	1.0	867.0
151603712	HuniePop	play	8.5	867.0
151603712	Path of Exile	purchase	1.0	39.0
151603712	Path of Exile	play	8.1	39.0
151603712	Poly Bridge	purchase	1.0	1347.0
151603712	Poly Bridge	play	7.5	1347.0
151603712	Left 4 Dead	purchase	1.0	49.0
151603712	Left 4 Dead	play	3.3	49.0
151603712	Team Fortress 2	purchase	1.0	1.0

Command took 6.82 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 6:13:35 AM on m1

Now my next step is data pre-processing. Data preprocessing mean the data should be in a correct format to train an MLlib model

So particularly in task like recommendation system when we are trying to predict a numeric values we ensures that the model treats all the ratings equally. In this system I considered to include both behavior “play” and “purchase” collectively in the single column name as status. But there is too much variation in status column from 1 to thousands. So I am going to apply normalization on the status column using the formula.

```
indexed_df = (col("Status") - min_status) / (max_status - min_status))
```

This formula will ensure that minimum status maps to Zero and maximum status maps to 1.

But if we do not apply normalization to status column than model does not treat all the ratings equally and there would be large value of RMSE(root mean square error) which is not acceptable in our case. So it is important to apply normalization.

So in this code I calculated the minimum value and maximum value of the status column and then apply the normalization formula

```

Cmd 7

1  from pyspark.sql.functions import col, min as pyspark_min, max as pyspark_max
2  from pyspark.ml.feature import StringIndexer, MinMaxScaler
3  from pyspark.ml.recommendation import ALS
4  from pyspark.ml.evaluation import RegressionEvaluator
5
6  # Calculate the minimum and maximum values of the "Status" column
7  min_status = indexed_df.select(pyspark_min("Status")).first()[0]
8  max_status = indexed_df.select(pyspark_max("Status")).first()[0]
9
10 # Normalize the "Status" column to a range between 0 and 1
11 indexed_df = indexed_df.withColumn("Normalized_Status", (col("Status") - min_status) / (max_status - min_status))

▶ (4) Spark Jobs

Command took 8.38 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 8:01:54 AM on m1

```

Next, I am going to divide the data into two sets: a training set and a test set because it is crucial for supervised machine learning, where we reserve a portion of the data for evaluating the model's performance. So I split the data randomly, allocating 80% for training and 20% for testing. To ensure reproducibility, I set the seed to 100, ensuring consistent results each time the code is run.

```
Cmd 8
Python ▶ ▾ - ✕
1 # Split the data into training and test sets
2 (training, test) = indexed_df.randomSplit([0.8, 0.2], seed=100)

Command took 0.59 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 8:07:29 AM on ml
```

#### 4. Identifying optimal hyperparameters, training the model, evaluating its performance, and utilizing MLflow for experiment tracking

I applied ALS matrix factorization to this recommender system. I used the following parameters

maxIter=10 (maximum number of iterations)  
regParam=0.01(for controlling regularization strenght)  
userCol: for identifying user  
itemCol: for identifying items  
ratingCol: for ratings associated with user-item interaction

```
Cmd 8
Python ▶ ▾ - ✕
1 # Define the ALS model with the required parameters
2 als = ALS(maxIter=10, regParam=0.01, userCol="member Id", itemCol="gameId", ratingCol="Normalized_Status",
3         coldStartStrategy="drop")

Command took 0.11 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 10:57:31 AM on dsx
```

Then I fit the ALS matrix factorization to the training dataset

```
Cmd 9
Python ▶ ▾ - ✕
1 # Fit the ALS model to the training data
2 model = als.fit(training)

▶ (7) Spark Jobs
▼ (1) MLflow run
  Logged 1 run to an experiment in MLflow. Learn more

2024/04/30 17:57:50 INFO mlflow.utils.autologging_utils: Created MLflow autologging run with ID '8fb8faa73e874a27914191189355ac88', which will track hyperparameters, performance metrics, model artifacts, and lineage information for the current pyspark.ml workflow
2024/04/30 17:58:18 WARNING mlflow.pyspark.ml: Model ALS_089fa6688ad1 will not be autologged because it is not allowlisted or 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 29.59 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 10:57:49 AM on dsx
```



Then I generate prediction on the test data using the already trained model named as “model” which applies the model to the test dataset and adds a new column named as “Prediction” to the data frame containing the predicted values.

Then I initiate a RegressionEvaluator with the following parameters

metricName="rmse"

labelCol="Normalized\_Status"

predictionCol="prediction"

Then I use that regression evaluator to calculate the RMSE value between Normalized\_Status values and the predicted values

```
1 predictions = model.transform(test)
2 evaluator = RegressionEvaluator(metricName="rmse", labelCol="Normalized_Status", predictionCol="prediction")
3 rmse = evaluator.evaluate(predictions)
4 print("Root Mean Squared Error (RMSE) on test data after normalization = {:.2f}".format(rmse))
5
6 # Generate predictions using the trained model
7 predictions.show()
```

► (5) Spark Jobs

► predictions: pyspark.sql.dataframe.DataFrame = [member Id: integer, Name of the Game: string ... 5 more fields]

2024/04/30 18:00:18 WARNING mlflow.utils.autologging\_utils: MLflow autologging encountered a warning: "/databricks/python/lib/python3.10/site-packages/mlflow/data/spark\_dataset.py:159: UserWarning: Hint: Inferred schema contains integer column(s). Integer columns in Python cannot represent missing values. If your input data contains missing values at inference time, it will be encoded as floats and will cause a schema enforcement error. The best way to avoid this problem is to infer the model schema based on a realistic data sample (training dataset) that includes missing values. Alternatively, you can declare integer columns as doubles (float64) whenever these columns may have missing values. See 'Handling Integers With Missing Values' <<https://www.mlflow.org/docs/latest/models.html#handling-integers-with-missing-values>>' for more details."

Root Mean Squared Error (RMSE) on test data after normalization = 0.01

ains missing values at inference time, it will be encoded as floats and will cause a schema enforcement error. The best way to avoid this problem is to infer the model schema based on a realistic data sample (training dataset) that includes missing values. Alternatively, you can declare integer columns as doubles (float64) whenever these columns may have missing values. See 'Handling Integers With Missing Values' <<https://www.mlflow.org/docs/latest/models.html#handling-integers-with-missing-values>>' for more details."

Root Mean Squared Error (RMSE) on test data after normalization = 0.01

member Id	Name of the Game	member behaviour	Status	gameId	Normalized_Status	prediction
5250	Cities Skylines	purchase	1.0	158	7.657032984796535E-5	1.0907977E-8
5250	Dota 2	play	0.2	0	8.507814427551706E-6	7.487604E-6
5250	Half-Life 2 Death...	purchase	1.0	13	7.657032984796535E-5	1.3883077E-9
5250	Portal	purchase	1.0	14	7.657032984796535E-5	2.1084414E-9
5250	Team Fortress 2	purchase	1.0	1	7.657032984796535E-5	5.0977206E-8
76767	Alien Swarm	play	0.8	32	5.955470099286195E-5	1.6275358E-7
76767	Banished	play	24.0	214	0.002033367648184...	1.3941788E-7
76767	Banished	purchase	1.0	214	7.657032984796535E-5	1.3941788E-7
76767	Call of Duty Blac...	purchase	1.0	50	7.657032984796535E-5	2.3781568E-7
76767	Call of Duty Blac...	play	12.5	57	0.001054968989016...	1.4741174E-6
76767	Call of Duty Mode...	play	65.0	25	0.005521571563481058	4.414076E-7
76767	Call of Duty Mode...	purchase	1.0	25	7.657032984796535E-5	4.414076E-7
76767	Counter-Strike Gl...	play	3.5	2	2.892656690536758E-4	2.9329345E-5
76767	Half-Life	play	1.2	45	9.358595870306876E-5	6.80719E-8
76767	Half-Life	purchase	1.0	45	7.657032984796535E-5	6.80719E-8
76767	Rise of Nations E...	play	5.7	691	4.764376079428956E-4	5.274042E-8
76767	Thief - Opportunist	purchase	1.0	994	7.657032984796535E-5	2.9432597E-8

Command took 29.99 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 11:00:05 AM on dsx

## Result:-

From above we can see that our RMSE value=0.01 which mean that the predicted values are very close to the actual values that's good.

## Hyperparameter tuning:-

Then I create a parameter grid for hyperparameter tuning. I define different hyperparameters as maximum iterations, Rank, and regularization parameter. Then I uses ParamGridBuilder to create a grid containing combinations of these hyperparameters.

```

1 from pyspark.ml.tuning import ParamGridBuilder
2
3 # Define the values for each parameter
4 maxIter_values = [5, 10, 15]
5 rank_values = [5, 10, 15]
6 regParam_values = [0.01, 0.1, 0.2]
7
8 # Create a parameter grid builder
9 paramGrid = (ParamGridBuilder()
10             .addGrid(als.maxIter, maxIter_values)
11             .addGrid(als.rank, rank_values)
12             .addGrid(als.regParam, regParam_values)
13             .build())
14
Command took 0.10 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 11:15:48 AM on dsx

```

Then I perform hyperparameter tuning using tvs. I configure TVS with an ALS estimator, a parameter grid, a evaluator for model performance and a specific train-validation split ratio of 80:20

```

1 from pyspark.ml.tuning import TrainValidationSplit
2
3 # Instantiate TrainValidationSplit
4 tvs = TrainValidationSplit(estimator=als,
5                           estimatorParamMaps=paramGrid,
6                           evaluator=evaluator,
7                           # 80% of the data will be used for training, 20% for validation
8                           trainRatio=0.8)
9
Command took 0.06 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 11:15:52 AM on dsx

```

Then I train the model using train validation split method.

```

1 # Train the model using grid search
2 model = tvs.fit(training)

```

(5) Spark Jobs  
 (28) MLflow runs  
 Logged 28 runs to an experiment in MLflow. [Learn more](#)

2024/04/30 18:15:54 INFO mlflow.utils.autologging\_utils: Created MLflow autologging run with ID '2df3d785fc184fad9019ac871e156c93', which will track hyperparameters, performance metrics, model artifacts, and lineage information for the current pyspark.ml workflow  
 2024/04/30 18:32:18 WARNING mlflow.pyspark.ml: Model TrainValidationSplitModel\_a00c04b16ce1 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 16.42 minutes -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 11:15:54 AM on dsx

## Best model:-

Then I retrieve the best model and access its parameters using getter method.

```
Cmd 13

1
2 # Get the best model from TrainValidationSplit using its getter method
3 best_model = model.getEstimator()
4
5
6 # Explore the public API of the best_model (usually ALS) to access parameters
7 print("Parameters for Best Model:")
8
9 # Example assuming ALS model:
10 print("maxIter:", best_model.getMaxIter())
11 print("rank:", best_model.getRank())
12 print("regParam:", best_model.getRegParam())
13 print("RMSE:", rmse)

Parameters for Best Model:
maxIter: 10
rank: 10
regParam: 0.01
RMSE: 0.010672728321225879

Command took 0.11 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 3:39:49 PM on fdc
```

## Multiple Runs:-

Then I uses multiple values of hyperparameters to train the model and check which parameter best fit with the model with lowest RMSE value.

## Second run:-

In my 2nd run I defined the hyperparameters as

```
maxIter_values = [5, 10, 20]
rank_values = [4, 10, 18]
regParam_values = [0.01, 0.001, 0.2]
```

and then train them by the procedure I explained earlier but here I have not use getter method for getting best model instead I use bestmodel attribute for getting best model and the rest procedure is same as above.

```
Cmd 14

1 from pyspark.ml.tuning import ParamGridBuilder, TrainValidationSplit
2 from pyspark.ml.evaluation import RegressionEvaluator
3 # Define the values for each parameter
4 maxIter_values = [5, 10, 20]
5 rank_values = [4, 10, 18]
6 regParam_values = [0.01, 0.001, 0.2]
7 # Create a parameter grid builder
8 paramGrid = ParamGridBuilder() \
9     .addGrid(als.maxIter, maxIter_values) \
10    .addGrid(als.rank, rank_values) \
11    .addGrid(als.regParam, regParam_values) \
12    .build()

Command took 0.12 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 4:52:32 PM on fdc
```

```
cmd 15

1
2 # Define the evaluator
3 evaluator = RegressionEvaluator(metricName="rmse", labelCol="Normalized_Status", predictionCol="prediction")
4 # Create TrainValidationSplit
5 tvs = TrainValidationSplit(estimator=als,
6                             estimatorParamMaps=paramGrid,
7                             evaluator=evaluator,
8                             trainRatio=0.8)
9 # Fit TrainValidationSplit to the training data
10 model = tvs.fit(training)
11 # Get the best model from TrainValidationSplit
12 best_model = model.bestModel
13 # Make predictions on the test data
14 predictions = best_model.transform(test)
15 # Evaluate the predictions using the evaluation metric
16 rmse = evaluator.evaluate(predictions)
17 print("Best Model RMSE:", rmse)
18

Best Model RMSE: 0.010672728321225879

Command took 0.08 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 3:44:33 PM on fdc
```

## Result:-

From this result we can see that we got RMSE value of 0.01 from the best hyperparameters which are

MaxIter=10

Rank=10

regParam=0.01

```
cmd 16

1 # Print the hyperparameters and RMSE
2 print("Parameters for Best Model:")
3 print("maxIter:", best_model._java_obj.getMaxIter())
4 print("rank:", best_model._java_obj.getRank())
5 print("regParam:", best_model._java_obj.getRegParam())
6 print("Best Model RMSE:", rmse)
7

Parameters for Best Model:
maxIter: 10
rank: 10
regParam: 0.01
Best Model RMSE: 0.010672728321225879

Command took 0.08 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 3:45:25 PM on fdc
```

## Third Run:-

In my 3rd run I changed the hyperparameters and defined the hyper parameters as

maxIter\_values = [10, 20, 30]

rank\_values = [8, 13, 20]

regParam\_values = [0.1, 0.001, 0.2]

and then train these hyper parameters on the model using the above procedure as in Second run.

```
Cmd 17
Python ▶ ▼ - X

1 from pyspark.ml.tuning import ParamGridBuilder, TrainValidationSplit
2 from pyspark.ml.evaluation import RegressionEvaluator
3
4 # Define the values for each parameter
5 maxIter_values = [10, 20, 30]
6 rank_values = [8, 13, 20]
7 regParam_values = [0.1, 0.001, 0.2]
8
9 # Create a parameter grid builder
10 paramGrid = ParamGridBuilder() \
11     .addGrid(als.maxIter, maxIter_values) \
12     .addGrid(als.rank, rank_values) \
13     .addGrid(als.regParam, regParam_values) \
14     .build()
15

Command took 0.08 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 4:57:02 PM on fdc

1 # Define the evaluator
2 evaluator = RegressionEvaluator(metricName="rmse", labelCol="Normalized_Status", predictionCol="prediction")
3 # Create TrainValidationSplit
4 tvs = TrainValidationSplit(estimator=als,
5                             estimatorParamMaps=paramGrid,
6                             evaluator=evaluator,
7                             trainRatio=0.8)
8 # Fit TrainValidationSplit to the training data
9 model = tvs.fit(training)
10 # Get the best model from TrainValidationSplit
11 best_model = model.bestModel
12 # Make predictions on the test data
13 predictions = best_model.transform(test)
14 # Evaluate the predictions using the evaluation metric
15 rmse = evaluator.evaluate(predictions)
16 print("Best model RMSE:", rmse)

(6) Spark Jobs
▼ (28) MLflow runs
  Logged 28 runs to an experiment in MLflow. Learn more
  predictions: pyspark.sql.dataframe.DataFrame = [member id: integer, Name of the Game: string ... 5 more fields]
2024/04/30 22:45:57 INFO mlflow.utils.autologging_utils: Created MLflow autologging run with ID 'd16e746bbae4e6a5427699353b62', which will track hyperparameters, performance metrics, model artifacts, and lineage information for the current pyspark.ml workflow
2024/04/30 23:39:36 WARNING mlflow.pyspark.ml: Model TrainValidationSplitModel_6b38c5837982 will not be autologged because it is not allowedlisted or because one or more of its nested models are not allowedlisted. 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).
Best model RMSE: 0.010672728321225879

Command took 53.88 minutes -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 3:45:56 PM on fdc
```

## Result:-

Here from the result we can see that we got the RMSE value of 0.0106 from the best hyperparameters which are

Maxiter=10

Rank=8

Regparam=0.1

```
Cmd 19
Python ▶ ▼ - X

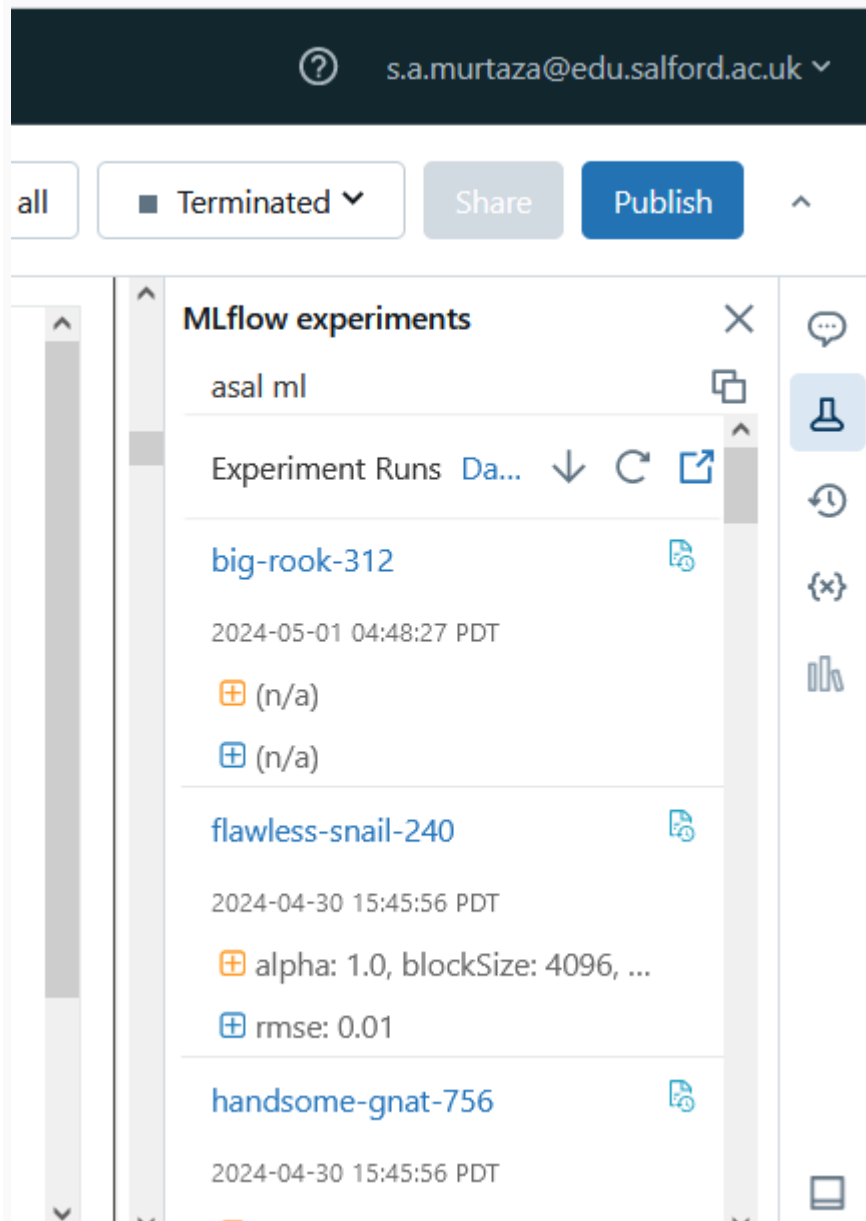
1 # Print the parameters of the best model
2 print("Best Model Parameters:")
3 print("MaxIter:", best_model._java_obj.parent().getMaxIter())
4 print("Rank:", best_model._java_obj.parent().getRank())
5 print("RegParam:", best_model._java_obj.parent().getRegParam())
6 print("Best Model RMSE:", rmse)

Best Model Parameters:
MaxIter: 10
Rank: 8
RegParam: 0.1
Best Model RMSE: 0.010672728321225879

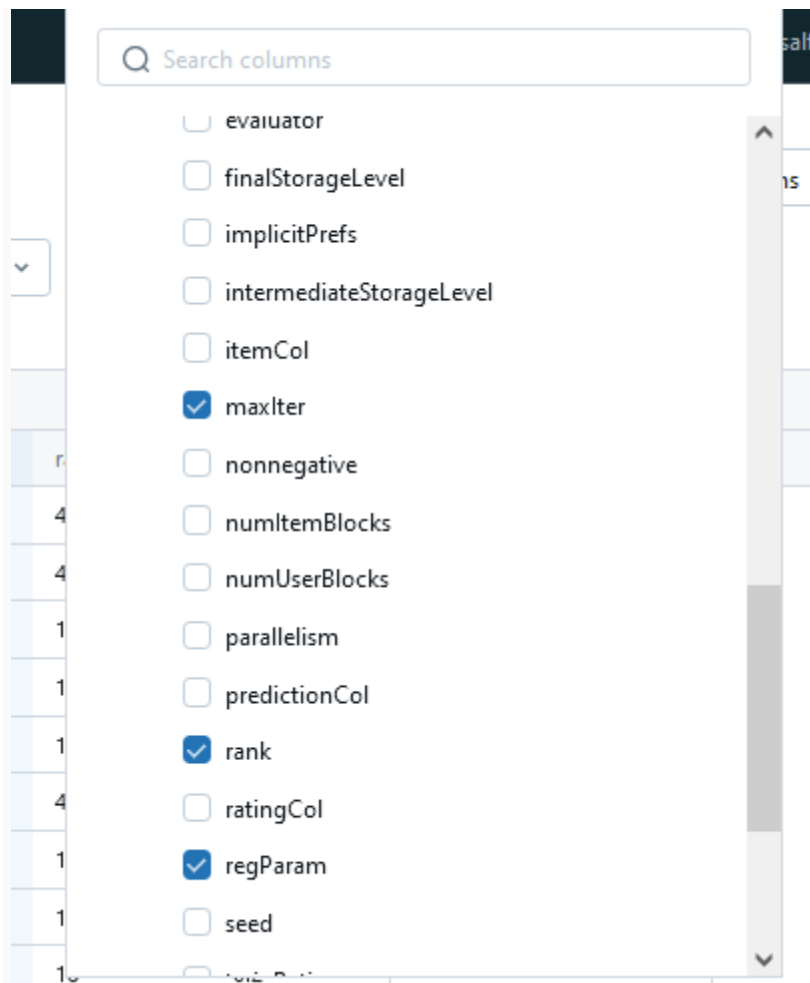
Command took 0.06 seconds -- by s.a.murtaza@edu.salford.ac.uk at 4/30/2024, 4:47:57 PM on fdc
```

## Using ML flow:-

Now as we performed multiple runs so I am going to use the experiment UI to view more information on all the runs that I performed earlier as part of the grid search. So for this purpose I clicked on the torch icon and select experiment.



I select all the hyperparameters by clicking on the column button.



Here you can see that I have performed 87 runs and result of each run is displayed. I clicked on RMSE column to sort the result in ascending order.

**databricks**

Experiments >

## Untitled Notebook 2024-04-30 14:14:37 [Provide Feedback](#)

Search: metrics.rmse < 1 and params.model = "tree"

State: created | Active | Datasets | Sort: rmse | Columns | Expand rows

						Metrics		Parameters	
Run Name	Created	Duration	Source	rmse	maxIter	rank	regParam		
secretive-panda-804	2 hours ago	21.2min	Untitled...	0.01208519...	10	10	0.001	Show more columns (35 total)	
capricious-toad-18	1 hour ago	53.4min	Untitled...	0.01208072...	20	20	0.001		
handsome-gnat-756	1 hour ago	53.4min	Untitled...	0.01206437...	10	13	0.001		
amazing-foal-195	2 hours ago	21.2min	Untitled...	0.01205262...	10	18	0.001		
stylish-turtle-242	1 hour ago	53.4min	Untitled...	0.01204640...	10	20	0.001		
exultant-newt-583	2 hours ago	21.2min	Untitled...	0.01204379...	5	10	0.001		
traveling-fowl-177	2 hours ago	21.2min	Untitled...	0.01191808...	5	18	0.001		
nimble-pig-878	2 hours ago	21.2min	Untitled...	0.01031156...	10	10	0.01		
orderly-toad-849	2 hours ago	16.3min	Untitled...	0.01031156...	10	10	0.01		
charming-pug-33	2 hours ago	16.3min	Untitled...	0.01031139...	10	15	0.01		
capricious-hawk-838	2 hours ago	21.2min	Untitled...	0.01031120...	10	18	0.01		
whimsical-flea-731	2 hours ago	16.3min	Untitled...	0.01031119...	15	10	0.01		

87 matching runs



1

Untitled Notebook 2024-04-30 14:14:37

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

Time created | State: Active | Datasets | Sort: rmse | Columns | Expand rows

Run Name	Created	Duration	Source	rmse	maxIter	rank	regParam
skillful-cod-315	1 hour ago	53.4min	Untitled...	0.010093417523583349	10	8	0.1
burly-goat-230	1 hour ago	53.4min	Untitled...	0.010093417523583349	20	8	0.1
rogue-shrike-73	1 hour ago	53.4min	Untitled...	0.010093417523583349	30	20	0.2
placid-goat-78	1 hour ago	53.4min	Untitled...	0.010093417523583349	10	13	0.2
sedate-cow-771	1 hour ago	53.4min	Untitled...	0.010093417523583349	10	20	0.2
nervous-grouse-683	1 hour ago	53.4min	Untitled...	0.010093417523583349	20	13	0.2
capricious-lamb-296	1 hour ago	53.4min	Untitled...	0.010093417523583349	30	20	0.1
exultant-hare-927	2 hours ago	21.2min	Untitled...	0.010093417523583349	20	4	0.2
nosy-snail-506	2 hours ago	21.2min	Untitled...	0.010093417523583349	5	10	0.2
bold-moth-385	2 hours ago	21.2min	Untitled...	0.010093417523583349	5	4	0.2
thundering-steed-860	2 hours ago	21.2min	Untitled...	0.010093417523583349	10	4	0.2
capricious-mare-206	2 hours ago	21.2min	Untitled...	0.010093417523583349	20	18	0.2

87 matching runs

## Description of results:-

I have trained the model with different hyperparameters as I described earlier. The RMSE value I got is 0.01 with a little bit difference when changing the hyperparameters value. A low RMSE value of 0.01 indicates that model is performing well and making accurate predictions.

Furthermore I got the lowest RMSE value of 0.010093417523583349 from the hyperparameters as

MaxIter=5  
Rank=15  
regParam=0.1

**databricks** s.a.murtaza@edu.salford.ac.uk

Experiments > **asal ml** [Provide Feedback](#)

Search: metrics.rmse < 1 and params.model = "tree" Time created State: Active Datasets Sort: rmse Columns Expand rows

**Table** Chart Evaluation Preview

					Metrics	Parameters			
	Run Name	Created	Duration	Source	rmse	maxIter	rank	regParam	
<input type="checkbox"/>	orderly-eel-756	16 hours ago	16.3min	Untitled ...	0.01009341...	15	5	0.2	
<input type="checkbox"/>	rare-ox-376	16 hours ago	16.3min	Untitled ...	0.01009341...	15	10	0.1	
<input type="checkbox"/>	nosy-sponge-916	16 hours ago	16.3min	Untitled ...	0.01009341...	5	15	0.1	
<input type="checkbox"/>	big-rook-312	2 hours ago	7.0s	asal ml	-	-	-	-	
<input type="checkbox"/>	glamorous-lamb-581	15 hours ago	53.7min	Untitled ...	-	-	-	-	
<input type="checkbox"/>	adorable-pig-168	15 hours ago	30.9min	Untitled ...	-	-	-	-	

88 matching runs

## Further Analysis

For the further analysis of data I write the following codes.

### 1. Number of actions per status and member behavior:-

I write a code to find the number of actions per status and member behavior combination by grouping the data frame by "Status," "Member\_Behavior," "Name of the Game," and "member id" column in the **OccupancyDF** and count the distinct entries

```

1
2
3 from pyspark.sql.functions import count
4
5 # Group by Status and Member_Behavior and count occurrences
6 action_counts = OccupancyDF.groupBy("member id","Name of the Game","Status", "member behaviour").count()
7
8 # Display the results (number of actions per status and member behavior combination)
9 action_counts.show()
10

```

▶ (2) Spark Jobs

▶ action\_counts: pyspark.sql.dataframe.DataFrame = [member id: double, Name of the Game: string ... 3 more fields]

member id	Name of the Game	Status	member behaviour	count
5.9945701E7	The Elder Scrolls...	58.0	play	1
5.9945701E7	Sid Meier's Civil...	22.0	play	1
5.9945701E7	L.A. Noire	13.8	play	1
5.9945701E7	Company of Heroes...	1.0	purchase	1
5.3875128E7	Arma 2	1.0	play	1
5.3875128E7	Transistor	1.0	purchase	1
6.2923086E7	Half-Life Blue Shift	1.0	purchase	1

▶ (2) Spark Jobs

▶ action\_counts: pyspark.sql.dataframe.DataFrame = [member id: double, Name of the Game: string ... 3 more fields]

member id	Name of the Game	Status	member behaviour	count
5.9945701E7	The Elder Scrolls...	58.0	play	1
5.9945701E7	Sid Meier's Civil...	22.0	play	1
5.9945701E7	L.A. Noire	13.8	play	1
5.9945701E7	Company of Heroes...	1.0	purchase	1
5.3875128E7	Arma 2	1.0	play	1
5.3875128E7	Transistor	1.0	purchase	1
6.2923086E7	Half-Life Blue Shift	1.0	purchase	1
6.5117175E7	BIT.TRIP BEAT	1.0	purchase	1
6.5117175E7	NightSky	1.0	purchase	1
1.1373749E7	The Beginner's Guide	1.9	play	1
1.1373749E7	BEEP	1.0	purchase	1
1.1373749E7	Samorost 2	1.0	purchase	1
5.6038151E7	Crysis	1.0	purchase	1
9823354.0	3DMark	1.0	purchase	1
1.9586574E7	Dynasty Warriors ...	11.0	play	1
1.98346312E8	Dragon Age Origins	1.0	purchase	1
2.6802825E7	Prince of Persia ...	1.0	purchase	1
4.2681063E7	Football Manager ...	3.0	play	1

Command took 4.18 seconds -- by s.a.murtaza@edu.salford.ac.uk at 5/1/2024, 4:48:48 AM on fxdcid

## 2. Top played games:-

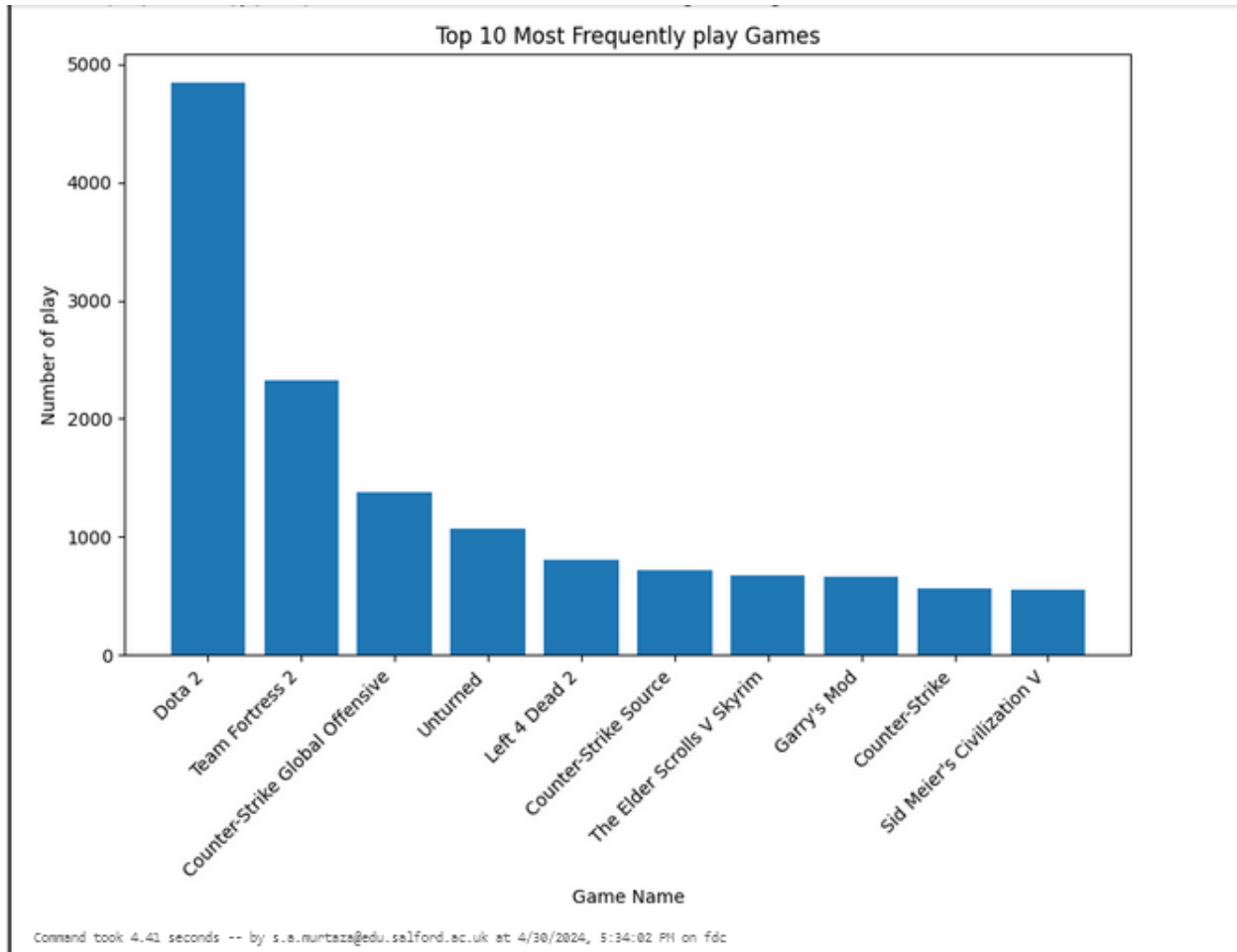
I write a code to define the number of top played games to visualize from the data frame indexed\_df and plot them using bar chart

```
1 # Import necessary libraries
2 import matplotlib.pyplot as plt
3
4 # Define the number of top games to visualize
5 N = 10 # You can change this value to visualize a different number of top games
6
7 # Summarize the most frequently purchased games
8 most_purchased = indexed_df.filter(indexed_df['member behaviour'] == 'play') \
9     .groupBy('Name of the Game') \
10     .count() \
11     .orderBy('count', ascending=False)
12
13 # Visualize the top N most frequently purchased games
14 top_n_purchased = most_purchased.limit(N)
15 top_n_purchased_df = top_n_purchased.toPandas()
16 plt.figure(figsize=(10, 6))
17 plt.bar(top_n_purchased_df['Name of the Game'], top_n_purchased_df['count'])
18 plt.xlabel('Game Name')
19 plt.ylabel('Number of play')
20 plt.title('Top {} Most Frequently play Games'.format(N))
21 plt.xticks(rotation=45, ha='right')
22 plt.show()
23
```

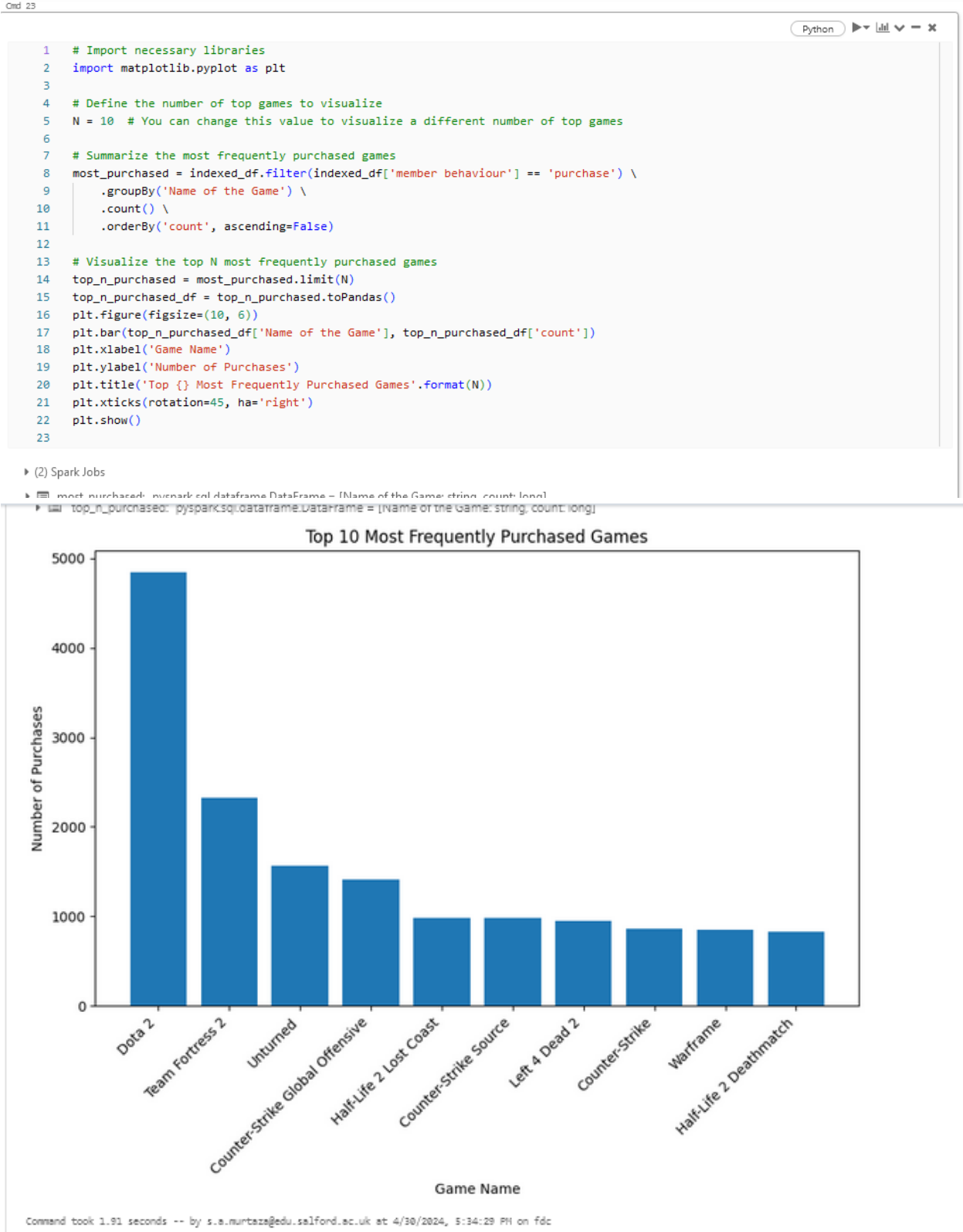
▶ (2) Spark Jobs

▶ `most_purchased:` pyspark.sql.dataframe.DataFrame = [Name of the Game: string, count: long]

**Result:-**



I write a code to define the number of top purchased games to visualize from the data frame indexed\_df and plot them using bar chart.



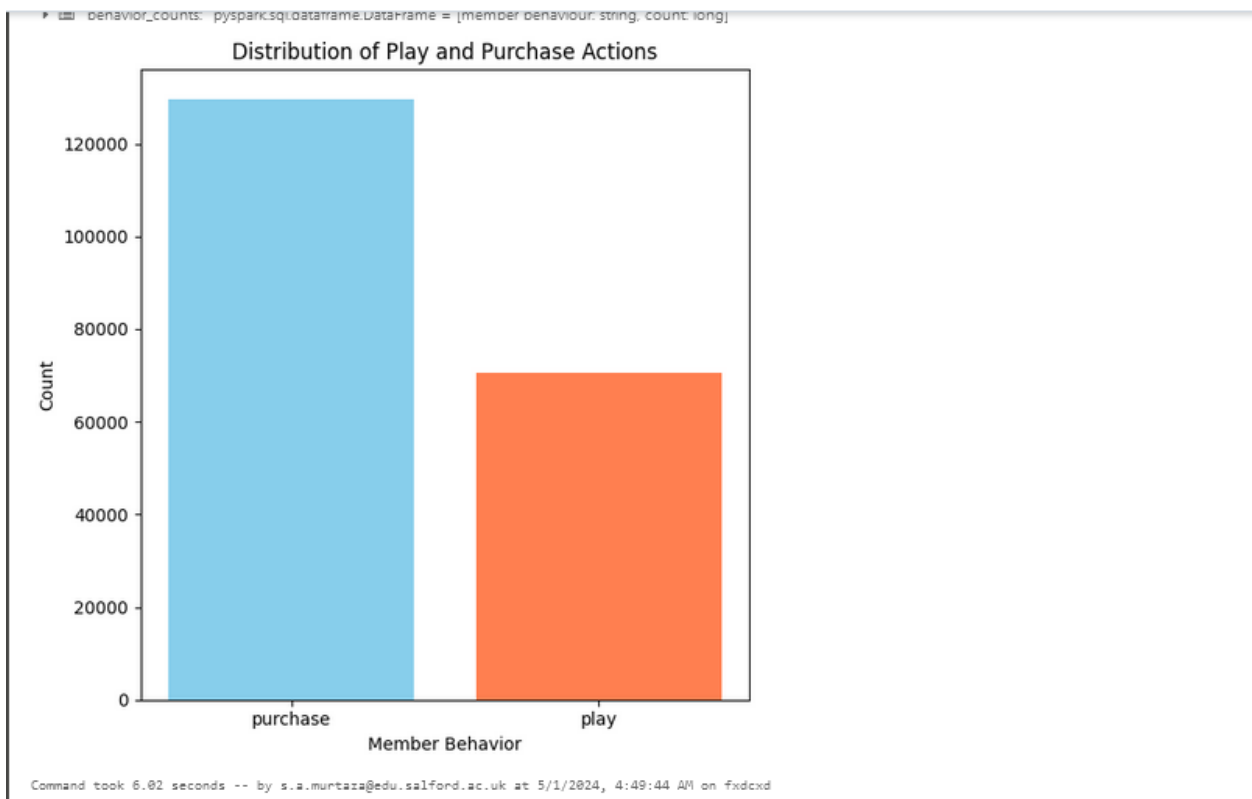
### 3. Distribution of Play and Purchase actions:-

Then I write a code to calculates the count of actions by member behaviour("purchase" or "play") from the data Frame "OccupancyDF" and visualizes the distribution using a bar chart. I first group the dataframe by member behavior and count the occurrences. Then I extract the data for the chart and create a bar chart with appropriate labels and color for each behavior.

```
1 # Assuming you have already loaded your data as 'staticInputDF' with appropriate schema
2
3 from pyspark.sql.functions import count
4
5 # Option 2: Count Actions by Member Behavior ("purchase" or "play")
6 behavior_counts = OccupancyDF.groupBy("member behaviour").count()
7 # Import libraries for plotting
8 import matplotlib.pyplot as plt
9
10 # Assuming 'behavior_counts' contains data for member behavior counts
11
12 # Extract data for the chart (assuming columns are named "member behaviour" and "count")
13 labels = list(behavior_counts.select("member behaviour").rdd.flatMap(lambda x: x).collect())
14 counts = list(behavior_counts.select("count").rdd.flatMap(lambda x: x).collect())
15
16 # Create a bar chart
17 plt.figure(figsize=(6, 6))
18 plt.bar(labels, counts, color=['skyblue', 'coral']) # Assign different colors for each behavior
19 plt.xlabel("Member Behaviour")
20 plt.ylabel("Count")
21 plt.title("Distribution of Play and Purchase Actions")
22 plt.xticks(rotation=0) # No rotation for member behavior labels
23 plt.tight_layout()
24 plt.show()
25
```

▶ (4) Spark Jobs  
▶ behavior\_counts: pyspark.sql.dataframe.DataFrame = [member behaviour: string, count: long]

**Result:-**



Then I use data bricks inbuilt tool for visualizing the data frame by mapping member id on axes member behavior on y-axes and group by them by status.

