In this task, we will explore the use of Apache Spark on Databricks to analyse clinical trial data and pharmaceutical company information. The goal is to learn about different facets of clinical trials—such as study types, popular diseases, leading sponsors, and concluded studies.

The following datasets are available for analysis:

**Clinical Trial Data:** This dataset includes study types, conditions, sponsors, and completion statuses for individual clinical trials.

**Pharmaceutical Company Information:** This dataset offers details on pharmaceutical companies, including parent companies and violations.

DataFrames, Spark SQL, and Spark RDDs will be used in the analysis.

Our goal in performing these analyses is to offer insightful information about the clinical trial landscape and pharmaceutical involvement, information that will help guide future research initiatives and decision-making processes. We will guarantee data accuracy throughout the analysis, address any possible inconsistencies, and clearly present the findings.

## SETUP

Using the Databricks community edition and the 12.2 LTS runtime version, the task was accomplished. The "Assessment" compute, which is displayed below, was utilised.
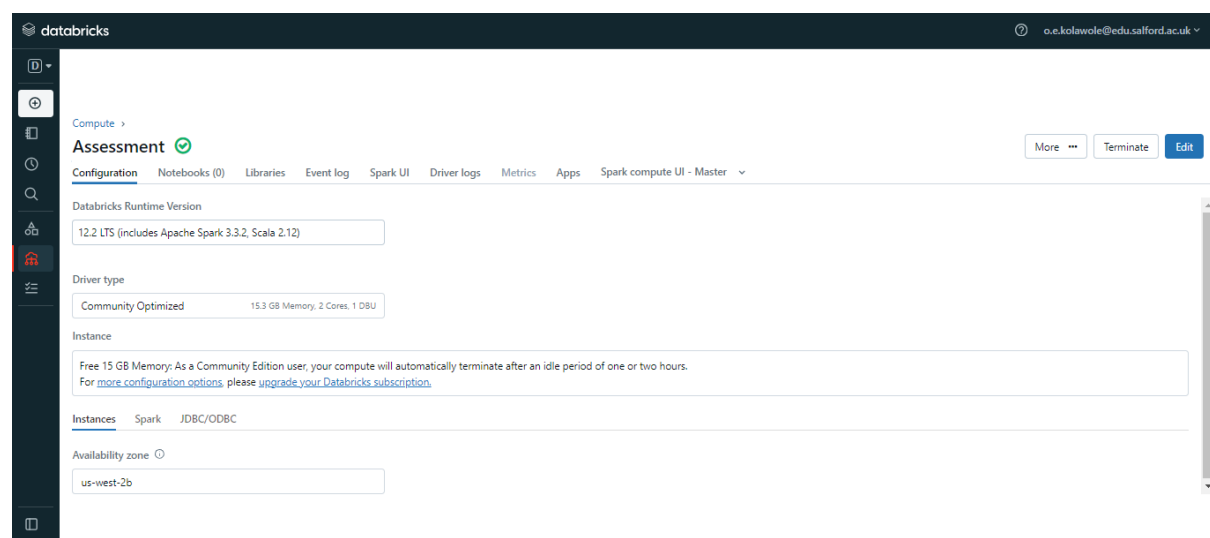


**Fig 1.0: Compute**

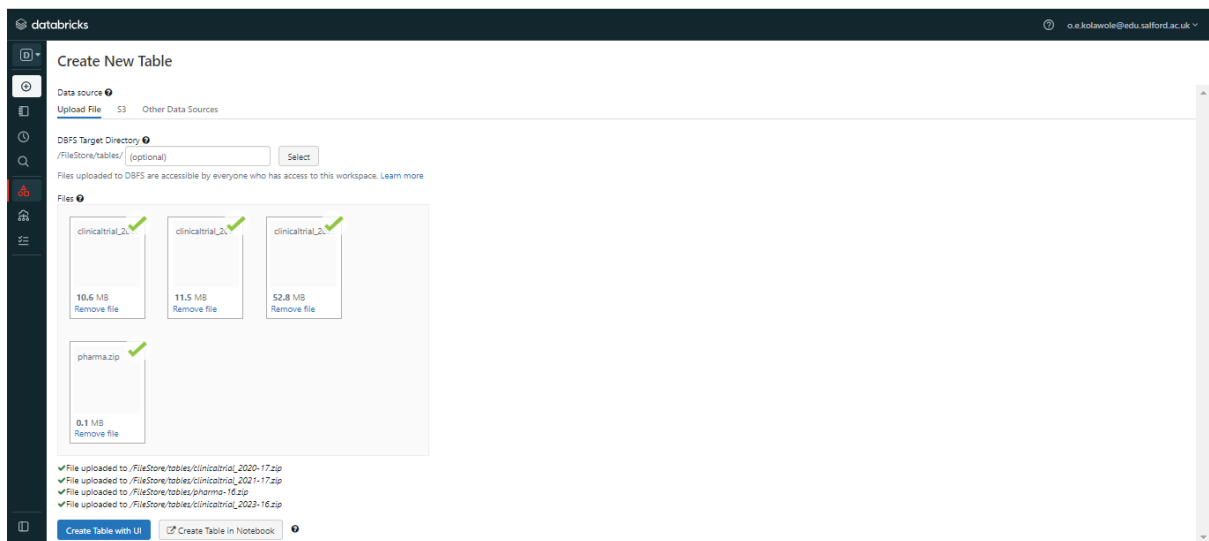The zip files were uploaded in the FileStore/tables directory. A screenshot of this process is shown below.



**Fig 1.1: Dataset uploading**
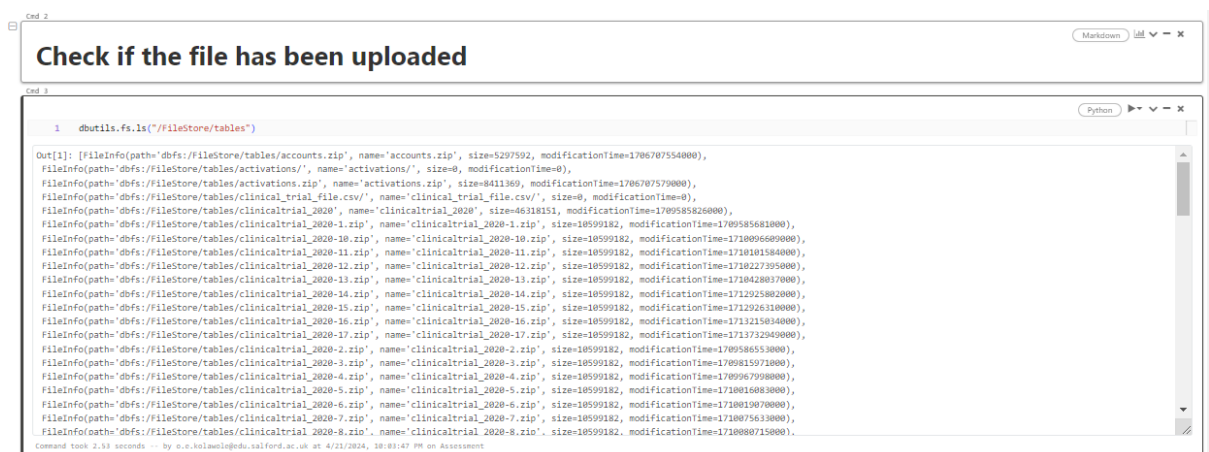
Checking if the file has been uploaded is done below



**Fig 1.2**

Inserting a reusable variable named 'fileroot' for the clinical trial data and 'pharma' for the pharma file
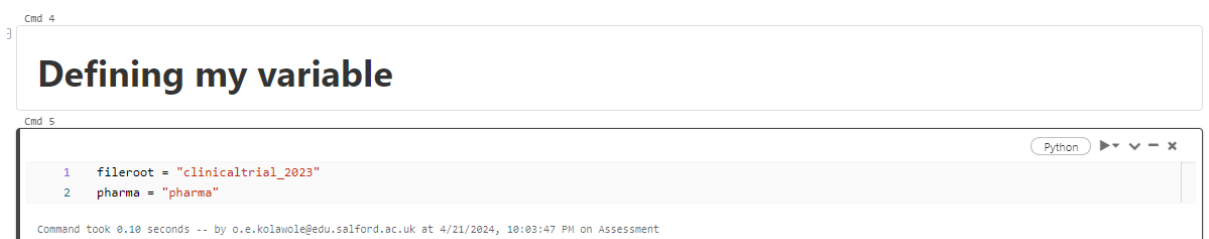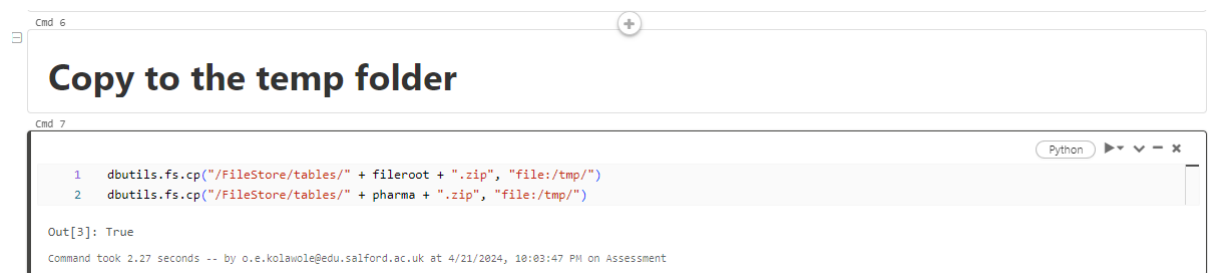


**Fig 1.3 Variable naming**

copying the zipped files to the tmp directory on my driver node

Cmd 6

## Copy to the temp folder

Cmd 7

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

Out[3]: True
```

Command took 2.27 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
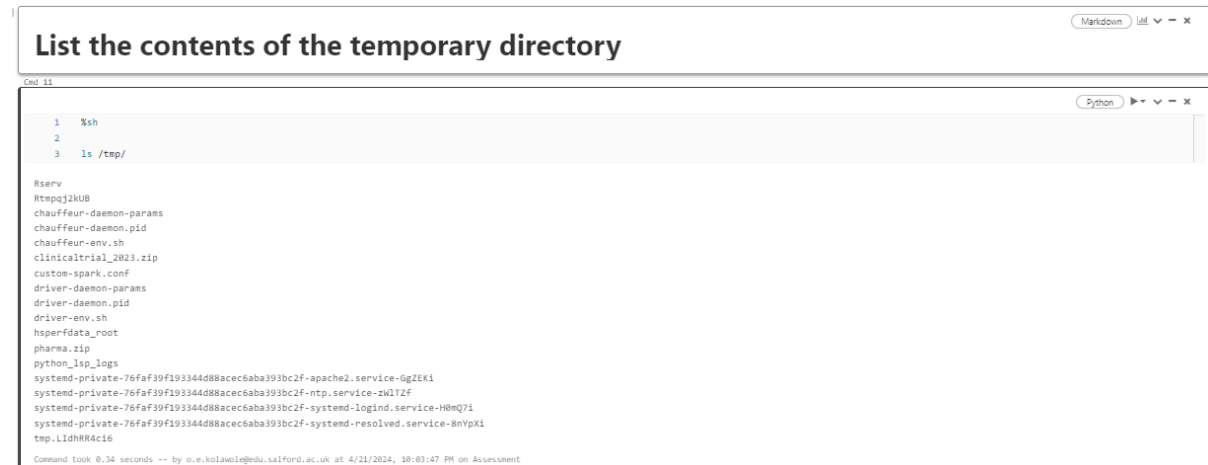
**Fig 1.4**

Cmd 8

## Making the variable accessible by the command line

Cmd 9

```python
1   import os
2   os.environ['fileroot'] = fileroot
3
4   import os
5   os.environ['pharma'] = pharma
```

Command took 0.06 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

**Fig 1.5**

Check the tmp directory

## List the contents of the temporary directory

Cmd 11

```python
1   %sh
2
3   ls /tmp/

Rserv
Rtmpqj2kUB
chauffeur-daemon-params
chauffeur-daemon.pid
chauffeur-env.sh
clinicaltrial_2023.zip
custom-spark.conf
driver-daemon-params
driver-daemon.pid
driver-env.sh
hsperfdata_root
pharma.zip
python_lsp_logs
systemd-private-76faf39f193344d88acec6aba393bc2f-apache2.service-GgZEKi
systemd-private-76faf39f193344d88acec6aba393bc2f-ntp.service-zW1TZf
systemd-private-76faf39f193344d88acec6aba393bc2f-systemd-logind.service-H0mQ7i
systemd-private-76faf39f193344d88acec6aba393bc2f-systemd-resolved.service-8nYpXi
tmp.LIdhRR4ci6
```

Command took 0.34 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

**Fig 1.6**

Unzipping the zip files into the tmp directory

## Cmd 12

## Unzip the files

### Cmd 13

```python
1   %sh
2
3   unzip -d /tmp/ /tmp/$fileroot.zip
4   unzip -d /tmp/ /tmp/$pharma.zip
```

```
Archive:  /tmp/clinicaltrial_2023.zip
  inflating: /tmp/clinicaltrial_2023.csv
Archive:  /tmp/pharma.zip
  inflating: /tmp/pharma.csv
Command took 2.65 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
```

**Fig 1.7**

### Cmd 14

## Checking the unzipped files

### Cmd 15

```python
1   %sh
2
3   ls /tmp/$fileroot.csv
4   ls /tmp/$pharma.csv
```

```
/tmp/clinicaltrial_2023.csv
/tmp/pharma.csv
Command took 0.09 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
```

**Fig 1.8**

Create a DBFS directory for the files

### Cmd 16

## Creating the new DBFS directory

### Cmd 17

```python
1   dbutils.fs.mkdirs("/FileStore/tables/" + fileroot)
2   dbutils.fs.mkdirs("/FileStore/tables/" + pharma)
```

```
Out[8]: True
Command took 0.18 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
```

**Fig 1.9**

Moving the unzipped CSV files from the tmp directory to the DBFS directory

### Cmd 18

## Moving the file to the new DBFS

### Cmd 19

```python
1   dbutils.fs.mv("file:/tmp/" + fileroot + ".csv", "/FileStore/tables/" + fileroot + ".csv", True)
2   dbutils.fs.mv("file:/tmp/" + pharma + ".csv", "/FileStore/tables/" + pharma + ".csv", True)
```

```
Out[9]: True
Command took 19.64 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
```

**Fig 2.0**

Confirming if the CSV files have been moved to the filestore/table



**Fig 2.1**

Checking the first few rows of the clinical trial dataset



**Fig 2.2**

## DATA CLEANING AND PREPARATION

## RDD

Cmd 24



### CREATE RDD FOR CLINICAL TRIALS

### defines a function to read and convert the clinicaltrial data into RDD

Cmd 25

```python
def load_clinical_trial_data(clinicaltrial_2023):
    """
    Load clinical trial data from a CSV file into an RDD.

    Parameters:
        clinicaltrial_2023 (str): The name of the CSV file to load.

    Returns:
        pyspark.rdd.RDD: An RDD containing the clinical trial data.
    """
    crdd = sc.textFile("/FileStore/tables/" + clinicaltrial_2023 + ".csv")
    return crdd
```

Command took 0.08 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

Fig 2.3 Clinical file RDD Creation

The code in the figure above creates RDD for clinical trials.

The function reads the data from the CSV file (Clinicaltrial 2023 + ".csv" + "/FileStore/tables/") into an RDD using the sc.textFile() method.

The RDD (crdd) with the clinical trial data is then returned by the function.

Cmd 26

### CREATE RDD FOR pharmaceutical data

### defines a function to read and convert the pharmaceutical data into RDD

Cmd 27

```python
def load_pharma_data(pharma):
    """
    Load pharmaceutical data from a CSV file into an RDD.

    Parameters:
        pharma (str): The name of the pharmaceutical data CSV file.

    Returns:
        pyspark.rdd.RDD: An RDD containing the pharmaceutical data.
    """
    prdd = sc.textFile("/FileStore/tables/" + pharma + ".csv")
    return prdd
```

Command took 0.09 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

Fig 2.4 Pharma file rdd creation

The function reads the contents of the CSV file ("/FileStore/tables/" + pharma + ".csv") supplied by pharma into an RDD using the sc.textFile() method. The RDD with the pharmaceutical data is then returned by the function.

## Removing the delimeters

```python
1   delimiter_selector = {
2       "clinicaltrial_2023": "\t",
3       "clinicaltrial_2021": "|",
4       "clinicaltrial_2020": "|",
5       "pharma": ","
6   }
7
8   def clean_clinical_rdd(crdd, clinicaltrial_2023):
9       """
10      Clean the clinical trial RDD by splitting each line using the delimiter
11      specified by the file type and removing unwanted characters.
12
13      Parameters:
14          crdd (pyspark.rdd.RDD): The RDD containing clinical trial data.
15          clinicaltrial_2023 (str): The name of the clinical trial file.
16
17      Returns:
18          pyspark.rdd.RDD: The cleaned RDD.
19      """
20      CLEAN_CRDD = crdd.map(lambda x: x.split(delimiter_selector[clinicaltrial_2023])).map(lambda x: [i.replace(",",'').replace('"','')for i in x])
21      return CLEAN_CRDD
22
23  def clean_pharma_rdd(prdd, pharma):
24      """
25      Clean the pharmaceutical RDD by splitting each line using the delimiter
26      specified by the file type and removing unwanted characters.
27
28      Parameters:
29          prdd (pyspark.rdd.RDD): The RDD containing pharmaceutical data.
30          pharma_file (str): The name of the pharmaceutical file.
31
32      Returns:
33          pyspark.rdd.RDD: The cleaned RDD.
34      """
35      CLEAN_PRDD = prdd.map(lambda x: x.split(delimiter_selector[pharma])).map(lambda x: [i.replace(",",'').replace('"','')for i in x])
36      return CLEAN_PRDD
```

Command took 0.09 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

**Fig 2.5 Removing the delimiters**

The two functions clean_pharma_rdd and clean_clinical_rdd, which are defined in this Python code, clean the data in RDDs that contain pharmaceutical and clinical trial data, respectively.

crdd.map(lambda x: x.split(delimiter_selector[clinicaltrial_2023])) splits each line of the RDD (crdd) into a list of elements using the delimiter specified by delimiter. The delimiter is chosen based on the clinical trial year.

.map(lambda x: [i.replace(",",'').replace('"','') for i in x]) removes unwanted characters (commas and double quotes) from each element. This also done for the pharma rdd.

## Creating new RDD after cleaning the data

```
1    # Creating new RDD  after cleaning the data
2
3    # Creating the new RDD
4    Clinical_RDD = load_clinical_trial_data(fileroot)
5    PHARMA_RDD = load_pharma_data(pharma)
6
7    # Cleaning the RDD
8    Clinical_RDD_Clean = clean_clinical_rdd(Clinical_RDD, fileroot)
9    PHARMA_RDD_Clean = clean_pharma_rdd(PHARMA_RDD, pharma)
10
11   # Take first 5 elements from cleaned Clinical RDD
12   Clinical_RDD_Clean.take(5)
```

▶ (1) Spark Jobs

```
Out[15]: [['Id',
 'Study Title',
 'Acronym',
 'Status',
 'Conditions',
 'Interventions',
 'Sponsor',
 'Collaborators',
 'Enrollment',
 'Funder Type',
 'Type',
 'Study Design',
 'Start',
 'Completion'],
 ['NCT03638471',
 'Effectiveness of a Problem-solving Intervention for Common Adolescent Mental Health Problems in India',
 'PRIDE',
 'COMPLETED',
 'Mental Health Issue (E.G. Depression Psychosis Personality Disorder Substance Abuse)',
 "BEHAVIORAL: PRIDE 'Step 1' problem-solving intervention|BEHAVIORAL: Enhanced usual care",
 'Sangath',
```

Command took 6.01 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

**Fig 2.6 Creating new rdd**

The load_clinical_trial_data and load_pharma_data is use to load the previously cleaned and then renamed. Thereafter, the cleaned clinical trial rdd is then display to show the first five rows.

## DATAFRAME

```
1    fileroot = "clinicaltrial_2023"
2    pharma = "pharma"
3
4    dbutils.fs.head("/FileStore/tables/" + fileroot + ".csv")
```

[Truncated to first 65536 bytes]
```
Out[1]: '"Id\tStudy Title\tAcronym\tStatus\tConditions\tInterventions\tSponsor\tCollaborators\tEnrollment\tFunder Type\tType\tStudy Design\tStart\tCompletio
n",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,\r\n"NCT03630471\tEffectiveness of a Problem-solving Intervention for Common Adolescent Mental Health Problems in India\tPRIDE\tCOMPLETED\tMen
tal Health Issue (E.G.", Depression, Psychosis, Personality Disorder," Substance Abuse)\tBEHAVIORAL: PRIDE \'Step 1\' problem-solving intervention|BEHAVIORAL: Enhanced usual car
e\tSangath\tHarvard Medical School (HMS and HSDM)|London School of Hygiene and Tropical Medicine\t250.0\tOTHER\tINTERVENTIONAL\tAllocation: RANDOMIZED|Intervention Model: PARALL
EL|Masking: 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|Cognition\tOTH
ER: Placebo|DIETARY_SUPPLEMENT: β-OHB\tMcMaster University\tAlzheimer\'s Society of Brant, Haldimand Norfolk," Hamilton Halton\t30.0\tOTHER\tINTERVENTIONAL\tAllocation: RANDOMI
ZED|Intervention Model: CROSSOVER|Masking: TRIPLE (PARTICIPANT," INVESTIGATOR," OUTCOMES_ASSESSOR)|Primary Purpose: BASIC_SCIENCE\t2023-10-25\t2024-0
8",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,,,,,,,,,,,,,,,,,\r\n"NCT00237471\tImpact of Tight Glycaemic Control in Acute Myocardial Infarction\t\tTERMINATED\tMyocardial Infarct|Hyperglycemia\tDRUG: Insulin
(tight blood glucose control)\tMelbourne Health\tNational Health and Medical Research Council," Australia|Bristol-Myers Squibb\t40.0\tOTHER\tINTERVENTIONAL\tAllocation: RANDOMI
ZED|Intervention Model: PARALLEL|Masking: NONE|Primary Purpose: TREATMENT\t2005-10\t2006-0
5",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,,,,,,,,,,,,,,,\r\n"NCT03820271\tNew Prognostic Predictive Models of Mortality of Decompensated Cirrhotic Patients Waiting for Liver Transplantation\tSUPERMEL
D\tRECRUITING\tDecompensated Cirrhosis|Liver Transplantation\tOTHER: SuperMELD\tAssistance Publique - Hôpitaux de Paris\t\t500.0\tOTHER\tINTERVENTIONAL\tAllocation: NA|Intervent
ion Model: SINGLE_GROUP|Masking: NONE|Primary Purpose: OTHER\t2020-10-01\t2023-10-0
1",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,...............................\r\n"NCT06229171\tInTake Care: Development and Validation of an Innovative"." Personalized Digital Health Solution for Medication Adherence Su
```

Command took 0.39 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment

**Fig 2.7**

The first few rows of the clinical trial file or file root was checked.



**Fig 2.8**

This code defines the create_dataframe function, which takes a CSV file containing clinical trial data and uses it to create a DataFrame. It Import 'monotonically_increasing_id' from pyspark.sql.functions which makes every row in a DataFrame has a unique ID. The delimiter selector was then defined which act as a key for each clinical trial file. The dataframe is created in response to the clinicaltrial_2023 by dividing the rows after reading the data as an RDD and ensuring every row has the same number of columns as the header even after extracting the header. Then the first 20 rows of the created dataframe is viewed.

## SQL

```
1   fileroot = "clinicaltrial_2023"
2   pharma = "pharma"
3
4   dbutils.fs.head("/FileStore/tables/" + fileroot + ".csv")
```

```
0",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,\r\n"NCT00368771\tA Six-month Study to Compare Outcome Differences and Visceral Response ... Irritable Bowel Syndrome\t\tCOMPLETED\tIrritable Bowel Syndrome\tBEHAVIORAL: IBS Stress Management|BEHAVI
ORAL: IBS Symptom Management|BEHAVIORAL: IBS Educational Training\tUniversity of California," Los Angeles\tNational Institute of Nursing Research (NINR)\t163.0\tOTHER\tINTERVENTIONAL\tAllocation: RANDOMIZ
ED|Intervention Model: PARALLEL|Masking: NONE|Primary Purpose: TREATMENT\t2002-07\t2010-0
5",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,\r\n"NCT04979871\tSARS-CoV-2 Antibodies and Virus Neutralisation in a Cohort Vaccinted Against COVID-19\tDER-CoV2-001\tCOMPLETED\tVaccine Reaction\tPROCEDURE: Venous bleeding\tUniversity of Zurich\t
\t50.0\tOTHER\tOBSERVATIONAL\tObservational Model: |Time Perspective: p\t2021-07-22\t2021-12-3
1",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,\r\n"NCT04097171\tThe Effect of Diet Composition on Performance", Expenditure, Blood Lipids," and Appetite Hormones in Highly Trained Cyclists\tDCAP\tCOMPLETED\tEndurance Cycling Performance\tOTHER:
Diet\tTexas Christian University\t\t34.0\tOTHER\tINTERVENTIONAL\tAllocation: RANDOMIZED|Intervention Model: CROSSOVER|Masking: NONE|Primary Purpose: BASIC_SCIENCE\t2019-11-11\t2020-03-3
0",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,\r\n"NCT06199271\tNeoadjuvant Adebrelimab Plus Dalpiciclib in Head and Neck Squamous Cell Carcinoma\t\tNOT_YET_RECRUITING\tHead and Neck Squamous Cell Carcinoma\tDRUG: Adebrelimab and dalpiciclib\tZhon
gzheng Xiang\t\t30.0\tOTHER\tINTERVENTIONAL\tAllocation: NA|Intervention Model: SINGLE_GROUP|Masking: NONE|Primary Purpose: TREATMENT\t2024-01-31\t2026-12-3
1",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,,,,,,,\r\n"NCT04032171\tStudy of Evobrutinib in Participants With RMS\t\tTERMINATED\tRelapsing-remitting Multiple Sclerosis\tDRUG: Evobrutinib|DRUG: Avonex®|DRUG: Avonex® matched Placebo|DRUG: Evobrutinib
matched Placebo\tEMD Serono Research & Development Institute"," Inc.\tMerck KGaA", Darmstadt," Germany\t1.0\tINDUSTRY\tINTERVENTIONAL\tAllocation: RANDOMIZED|Intervention Model: PARALLEL|Masking: QUADRUPLE
(PARTICIPANT", CARE_PROVIDER, INVESTIGATOR," OUTCOMES_ASSESSOR)|Primary Purpose: TREATMENT\t2019-09-10\t2020-05-2
0",,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
,\r\n"NCT05674071\tApplying an Osteopathic Intervention to Improve Mental Health Symptoms: a Mixed-methods Feasibility Study Protocol.\t\tCOMPLETED\tMental Health Issue\tBEHAVIORAL: Articulation/HVT|BEHAVI
ORAL: Soft-tissue massage|BEHAVIORAL: Craniosacral techniques|BEHAVIORAL: Combination of the three interventions: HVT"," soft-tissue and craniosacral techniques\tSwansea University\tOsteopathic Foundation|
University College of Osteopathy\t32.0\tOTHER\tINTERVENTIONAL\tAllocation: RANDOMIZED|Intervention Model: PARALLEL|Masking: SINGLE (OUTCOMES_ASSESSOR)|Primary Purpose: TREATMENT\t2022-12-20\t2023-08-0
```

Command took 0.49 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment

## Fig 2.9

The first few rows of the clinical trial file or file root were checked.

```
26          rdd = sc.textFile(f"/FileStore/tables/{clinicaltrial_2023}.csv").map(lambda row: row.replace(',', '').replace('"', '').split(delimiter_selector
            [clinicaltrial_2023]))
27          head = rdd.first()
28          rdd = rdd.map(lambda row: row + [" " for i in range(len(head) - len(row))] if len(row) < len(head) else row )
29          df = rdd.toDF()
30          first = df.first()
31          for col in range(0, len(list(first))):
32              df = df.withColumnRenamed(f"_{col + 1}", list(first)[col])
33          df = df.withColumn('index', monotonically_increasing_id())
34          return df.filter(~df.index.isin([0])).drop('index')
35      else:
36          return spark.read.csv(f"/FileStore/tables/clinicaltrial_2023.csv", sep=delimiter_selector[clinicaltrial_2023], header = True)
37
38  clinical_dataframe = create_clinical_dataframe(fileroot)
39  clinical_dataframe.show(20)
```

▶ (6) Spark Jobs

▶ ▦ pharma_dataframe: pyspark.sql.dataframe.DataFrame = [Company: string, Parent_Company: string … 32 more fields]

▶ ▦ clinical_dataframe: pyspark.sql.dataframe.DataFrame = [Id: string, Study Title: string … 12 more fields]

```
+-----------+--------------------+----------+----------+----------------+--------------------+--------------------+--------------------+--------------------+----------+-----------+
|         Id|         Study Title|   Acronym|    Status|      Conditions|       Interventions|             Sponsor|       Collaborators|Enrollment|Funder Type|
|Type|        Study Design|     Start|Completion|
+-----------+--------------------+----------+----------+----------------+--------------------+--------------------+--------------------+----------+-----------+
|NCT03630471|Effectiveness of ...|     PRIDE|   COMPLETED|Mental Health Iss...|BEHAVIORAL: PRIDE...|            Sangath|Harvard Medical S...|     250.0|     OTHER|INTERVEN
TIONAL|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|INTERVEN
TIONAL|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|INTERVEN
TIONAL|Allocation: RANDO...|   2005-10|   2006-05|
|NCT03820271|New Prognostic Pr...| SUPERMELD|   RECRUITING|Decompensated Cir...|    OTHER: SuperMELD|Assistance Publiq...|                    |     500.0|     OTHER|INTERVEN
TIONAL|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|INTERVEN
TIONAL|Allocation: RANDO...|2024-10-01|2026-04-01|
|NCT02945371|Tailored Inhibito...|       REV|   COMPLETED|Smoking|Alcohol D...|BEHAVIORAL: Perso...|University of Oregon|                    |     103.0|     OTHER|INTERVEN
TIONAL|Allocation: RANDO...|   2014-09|   2016-05|
|NCT01055171|Neuromodulation o...|          |   COMPLETED|Alcohol Dependenc...|DRUG: Propranolol...|Medical Universit...|National Institut...|      44.0|     OTHER|INTERVEN
TIONAL|Allocation: RANDO...|   2010-01|   2012-08|
|NCT01125371|Computerized Brie...|          |   COMPLETED|Alcohol: Harmful ...|BEHAVIORAL: Compu...|Johns Hopkins Uni...|National Institut...|     439.0|     OTHER|INTERVEN
```

Command took 13.73 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment

## Fig 3.0

The "import * from pyspark.sql.types and import * from pyspark.sql.functions" imports the libraries needed to transform and manipulate the data

The Purpose of the Pharmaceutical DataFrame: create_pharma_dataframe and create_clinical_dataframe returns the dataframe after receiving the pharma and clinicaltrial_2023 dataset as input.

Then views are created for each of the dataframe.



**Fig 3.1**

SQL queries can be run against the clinical trial and pharmaceutical datasets now that temporary views have been created for them. It will be simpler to carry out different analyses and respond to the questions as a result.



**Fig 3.2**

A SQL query to extract every column from the 2023 clinical trial dataset. This will show every row and every column in the dataset.

**Fig 3.3**

A SQL query to select all columns from the pharmaceutical dataset. This will display all the rows and columns of the dataset.

## PROBLEM ANSWERS

## QUESTION 1
**Assumptions**

- Each row in the study represents a distinct study.
- In the clinical_data, there are no duplicate ID rows.

**RDD**

The code uses the ".first()" to extract the header from the Clinical_RDD. The header row is then filtered out, duplicate rows are removed using the distinct() transformation, and the number of remaining rows is finally counted to determine the number of distinct studies in the dataset.

```python
# QUESTION 1: NUMBER OF STUDIES IN THE DATASET

header = Clinical_RDD.first()

# Extract study names and count distinct studies
num_studies = (
    Clinical_RDD
    .filter(lambda row: row != header)
    .distinct()
    .count()
)

print("Number of studies in the dataset:", num_studies)
```

> (2) Spark Jobs

Number of studies in the dataset: 483422

Command took 14.80 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

## DATA FRAME

Duplicate rows are eliminated by using the DataFrame's distinct() function, and the number of unique rows that remain is then counted using the count() function.

```
Cmd 3
                                                                              Python ▶▾ ∨ — ✕
    1   # QUESTION 1: NUMBER OF STUDIES IN THE DATASET
    2
    3   # Count the number of distinct rows in the DataFrame, effectively counting the number of unique rows.
    4   clinical_dataframe.distinct().count()
    5

▸ (3) Spark Jobs
Out[3]: 483422
Command took 17.46 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment
```

## SQL

The DISTINCT count in the SQL will count the number of unique rows in the clinicaltrial_2023 table using the 'SELECT' and 'FROM' statements.

```
Cmd 6
                                                                              SQL ▶▾ ⊔⊔ ∨ — ✕
    1   %sql
    2   SELECT DISTINCT count(*)
    3   FROM clinicaltrial_2023

▸ (2) Spark Jobs
▸ ▭ _sqldf: pyspark.sql.dataframe.DataFrame = [count(1): long]

 Table  ∨    +

        count(1)   ▲
    1   483422

  ↓  1 row   | 12.90 seconds runtime                                  Refreshed 15 minutes ago
Command took 12.90 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment
```

## Discussion of Result

The total number of distinct studies on the 2023 clinical trial is represented by 483,422 rows in the data. This also sheds light on the quantity and range of the clinical trial that were carried out.

## QUESTION 2
### Assumptions

The "Type" column in the "clinicaltrial_2023" contains the different types of clinical trials performed and are likely to be replicated.

## RDD

```python
# QUESTION TWO: ALL TYPES OF TRIALS & THEIR FREQUENCY

type_of_trial = Clinical_RDD_Clean.first().index('Type')
Clinical_RDD_Clean.filter(lambda x: len(x) > type_of_trial + 1).map(lambda x: (x[type_of_trial], 1)).filter(lambda row: row[0] != 'Type').reduceByKey
(lambda a,b: a + b).filter(lambda x: x[0] != '').sortBy(lambda x: x[1], ascending=False).collect()
```

▸ (4) Spark Jobs

```
Out[18]: [('INTERVENTIONAL', 371382),
 ('OBSERVATIONAL', 110221),
 ('EXPANDED_ACCESS', 928)]
```

Command took 13.79 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

The code first finds the index of the 'Type' column in your RDD then removes any rows that don't have the 'Type' column then maps each row to a tuple where the first element is the value in the 'Type' column and the second element is 1, indicating a count of one. It then filters out any tuples where the value in the 'Type' column is 'Type'. There sums up the counts for each type and filters out any empty values in the 'Type' column and sorts the results by the count of each type in descending order.

The result is then displayed.

## DATA FRAME

```python
# QUESTION TWO: ALL TYPES OF TRIALS & THEIR FREQUENCY
clinical_dataframe.groupBy('Type').count().orderBy('count', ascending=False).show(3)
```

▸ (2) Spark Jobs

```
+---------------+------+
|           Type| count|
+---------------+------+
| INTERVENTIONAL|371382|
|  OBSERVATIONAL|110221|
|EXPANDED_ACCESS|   928|
+---------------+------+
only showing top 3 rows
```

Command took 12.24 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment

This query groups the dataframe by the 'Type' column , counts the number of occurrences of each distinct value in the column within each group, orders it in descending order then shows the top three rows.

## SQL



```sql
%sql
SELECT clinicaltrial_2023.Type, count(*) as count FROM clinicaltrial_2023
GROUP BY clinicaltrial_2023.Type
ORDER BY count DESC
LIMIT 3
```

| | Type | count |
|---|---|---|
| 1 | INTERVENTIONAL | 371382 |
| 2 | OBSERVATIONAL | 110221 |
| 3 | EXPANDED_ACCESS | 928 |

3 rows | 13.21 seconds runtime          Refreshed 15 minutes ago

Command took 13.21 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment

The query will group the clinicaltrial_2023 file by their type, count the number of trials for each type, order them by count in descending order, and then limit the result to the top 3 types from the select statements.

## Discussion of Result

It appears that the most common type of clinical trial is "INTERVENTIONAL," with a frequency of 371,382 trials. This suggests that a significant portion of the dataset consists of trials involving interventions or treatments, where researchers actively intervene to study their effects on participants.

The second most common type is "OBSERVATIONAL," with a frequency of 110,221 trials. Observational trials involve observing participants in their natural environment without any intervention.

"EXPANDED_ACCESS" trials are the least common, with a frequency of 928 trials. Expanded access trials, also known as compassionate use or early access programs, provide investigational treatments to patients with serious diseases or conditions outside of clinical trials when no satisfactory alternative treatments are available.

## QUESTION 3

## Assumptions

The "Conditions" column in the "clinicaltrial_2023" contains the different conditions or diseases being studied in each clinical trial performed and are likely to be reoccurring.

## RDD

```
% Cmd 36

1    # QUESTION 3: TOP 5 CONDITIONS WITH THEIR FREQUENCIES
2
3    conditions_delimeter = {
4        "clinicaltrial_2023": "\t",
5        "clinicaltrial_2021": ",",
6        "clinicaltrial_2020": ",",
7    }
8    conditions_column_index = Clinical_RDD_Clean.first().index('Conditions')
9
10   Clinical_RDD_Clean.flatMap(lambda x: x[conditions_column_index].split(conditions_delimeter[fileroot])).filter(lambda row: row != 'Conditions').filter
     (lambda row: row != '').map(lambda x: (x, 1)).reduceByKey(lambda a,b: a + b).sortBy(lambda x: x[1], ascending=False).take(5)
```

▶ (4) Spark Jobs

```
Out[19]: [('Healthy', 7997),
('Breast Cancer', 4556),
('Prostate Cancer', 2650),
('Asthma', 2309),
('Obesity', 2284)]
```

Command took 16.20 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment

In this code, the index of the 'Conditions' column was derived from the header of the RDD then the column was split for each row using the appropriate delimiter. The header and empty rows were filtered out then the remaining rows were counted by the occurrence of each status in descending order. The top 5 most frequent conditions of the sorted RDD were then displayed.

## DATA FRAME

```
Cmd 5
                                                                                               Python   ▶ ∨ ─ ✕

1    # QUESTION 3: TOP 5 CONDITIONS WITH THEIR FREQUENCIES
2
3    from pyspark.sql.functions import split, explode, trim, col
4
5    # Define delimiter for the Conditions column based on the clinical trial data
6    conditions_delimiter = {
7        "clinicaltrial_2023": "\t",
8        "clinicaltrial_2021": ",",
9        "clinicaltrial_2020": ",",
10   }
11
12   # Split the Conditions column based on the delimiter for the specific clinical trial data
13   split_conditions_df = clinical_dataframe \
14       .withColumn('Conditions', explode(split(trim(col('Conditions')), conditions_delimiter[fileroot])))
15
16   # Group by Conditions, count occurrences, and filter out empty values
17   result_df = split_conditions_df \
18       .groupBy('Conditions') \
19       .count() \
20       .orderBy('count', ascending=False) \
21       .filter("Conditions != ''")
22
23   # Show top 5 results
24   result_df.show(5, truncate=False)
```

▶ (2) Spark Jobs

▶ ▤ split_conditions_df: pyspark.sql.dataframe.DataFrame = ["Id: string, Study Title: string ... 12 more fields]
▶ ▤ result_df: pyspark.sql.dataframe.DataFrame = [Conditions: string, count: long]

```
+---------------+-----+
|Conditions     |count|
+---------------+-----+
|Healthy        |7997 |
|Breast Cancer  |4556 |
|Prostate Cancer|2650 |
|Asthma         |2309 |
|Obesity        |2284 |
+---------------+-----+
only showing top 5 rows
```

Command took 13.59 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment

The 'Conditions' column was split for each row using the delimiter corresponding to the fileroot variable and a new row was created using the explode function for each condition any empty row is also removed from the condition value. The result was then grouped by 'Conditions' column and the number of the occurrences was counted, and it is then ordered in descending order. The top 5 conditions is then displayed with its full name and number of occurrences.

**SQL**



Create a temporary view called all_conditions by splitting the clinicaltrial_2023.conditions column by comma delimiter and creating a new row for each condition. Then, select the 'conditions' column and count the number of occurrences using count(*), group results by conditions, order results in descending order, and limit output to top 5 rows.

**Discussion of result**

With a total of 7,997 occurrences, the results show that the majority of the clinical trial condition is "Healthy." When compared to the total number of trials conducted, the frequency of the top five conditions does not offer significant insights into the clinical trials.

**QUESTION 4**

**Assumptions**

Organizations listed in the "Parent Company" column of the "pharma" dataset are assumed to be pharmaceutical companies.

The "Sponsor " column in the "clinicaltrial_2023" comprises the names of companies that are both pharmaceutical and non pharmaceutical.

## RDD

```
Cmd 37                                                          +

  1    # QUESTION 4,  TOP 10 NON - PHARMA COMPANIES
  2
  3    clinical_trial_sponsor_col_index = Clinical_RDD_Clean.first().index('Sponsor')
  4
  5    parent_pharm_comp = PHARMA_RDD_Clean.map(lambda x: x[1].replace('"', ''))
  6
  7    Clinical_RDD_Clean.map(lambda x: x[clinical_trial_sponsor_col_index]).filter(lambda row: row != 'Sponsor').subtract(PHARMA_RDD_Clean.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[20]: [('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 18.87 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
```

The code retrieves the 'Sponsor' column index from the Clinical_RDD_Clean DataFrame, creates an RDD for pharmaceutical companies, and removes double quotes using the replace function. Procedures include mapping the 'Sponsor' column, filtering out 'Sponsor' headers, and subtracting the RDD.

## DATA FRAME

```
Cmd 6                                                                                     Python  ▶▾ ∨ − ×

  1    # QUESTION 4: RETRIEVE THE TOP 10 SPONSORS THAT ARE NOT PHARMACEUTICAL COMPANIES
  2    pharma_list = create_dataframe(pharma).select("Parent_Company").rdd.flatMap(lambda x: x).collect()
  3    clinical_sponsor_dataframe = clinical_dataframe.select("Sponsor")
  4
  5    non_pharma_sponsors = clinical_sponsor_dataframe.groupBy("Sponsor").count().orderBy("count", ascending=False).filter(~clinical_sponsor_dataframe.Sponsor.isin(pharma_list)).show(10)

 ▸ (4) Spark Jobs
 ▸ ▦ clinical_sponsor_dataframe: 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 20.59 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment
```

The pharma dataset is converted into a DataFrame using the create_dataframe(pharma) function. The 'Parent_Company' column is selected and flattened into a single list using the rdd.flatMap(lambda x: x) function. The RDD is materialized and all values are retrieved. A new DataFrame named clinical_sponsor_dataframe is created, containing the 'Sponsor' column. The

dataframe is manipulated using groupBy("Sponsor"), count(), orderBy("count"), filter(pharma_list), and show(10).

**SQL**



```sql
%sql
CREATE OR REPLACE TEMP VIEW non_pharma_sponsor AS SELECT Sponsor FROM clinicaltrial_2023 WHERE Sponsor NOT IN (SELECT Parent_Company FROM pharma);
SELECT Sponsor, count(*) as count FROM non_pharma_sponsor
GROUP BY Sponsor
ORDER BY count DESC
LIMIT 10
```

▶ (3) Spark Jobs

▶ ▤ _sqldf: pyspark.sql.dataframe.DataFrame = [Sponsor: string, count: long]

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 |

⬇ 10 rows | 14.25 seconds runtime          Refreshed 17 minutes ago

Command took 14.25 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment

A temporary view named non_pharma_sponsor is created to filter out non-pharmaceutical sponsors from the clinicaltrial_2023 table. The WHERE clause removes sponsors from the pharma table, allowing only non-pharmaceutical sponsors. The SQL query uses the 'Sponsor' column to determine the number of sponsor appearances, group results, sort outcomes in descending order, and restrict output to the top ten rows. The results are sorted in descending order.

**Discussion of result**

The top non-pharmaceutical company sponsoring clinical trials is National Cancer Institute (NCI) having a frequency of 3410.

These results highlight a diverse range of organizations, including research institutions, hospitals, and government agencies, actively sponsoring clinical trials.

**QUESTION 5**

**Assumptions**

The "Start" column in the "clinicaltrial_2023 " dataset contains the start date of each clinical trial.

The "Completion" column contains the completed studies per month.

## RDD

```
1   # QUESTION 5 The completed trial for 2023
2   # Assuming month_selector, ct_completion_col_index, ct_status_col_index, year, and CLINICAL_RDD_Clean are defined
3   year = {
4       "clinicaltrial_2020": "2020",
5       "clinicaltrial_2021": "2021",
6       "clinicaltrial_2023": "2023",
7   }
8   Clinical_RDD_Clean = Clinical_RDD_Clean.map(lambda x: [i.replace(',', '').replace('"', '') for i in x])
9
10  # Define a dictionary to map month abbreviations to their corresponding full names
11  month_selector = {"01": "Jan", "02":"Feb", "03":"Mar", "04":"Apr", "05":"May", "06":"Jun", "07":"Jul", "08":"Aug", "09":"Sept", "10":"Oct", "11":"Nov", "12":"Dec"}
12
13  # Convert month abbreviations to their corresponding numerical representations
14  month_numerical = {month: int(num) for num, month in month_selector.items()}
15
16  ct_completion_col_index = Clinical_RDD_Clean.first().index('Completion')
17  ct_status_col_index = Clinical_RDD_Clean.first().index('Status')
18
19  # Filter and process the data to find completed studies in 2023
20  completed_studies = Clinical_RDD_Clean \
21      .filter(lambda x: len(x) > ct_completion_col_index and len(x) > ct_status_col_index) \
22      .map(lambda x: (x[ct_completion_col_index], x[ct_status_col_index])) \
23      .filter(lambda x: x[0] != 'Completion' and (x[1] == 'COMPLETED' or x[1] == 'Completed')) \
24      .map(lambda x: (x[0][5:7], x[0][0:4])) \
25      .filter(lambda x: x[1] == year["clinicaltrial_2023"]) \
26      .map(lambda x: (month_selector[x[0]], 1)) \
27      .reduceByKey(lambda a, b: a + b) \
28      .sortBy(lambda x: month_numerical[x[0]])
29
30  # Convert RDD to the specified format
31  result = completed_studies.collect()
32
33  for month, count in result:
34      print(f'({month}, {count})')
```

```
▶ (5) Spark Jobs
(Jan, 1494)
(Feb, 1272)
(Mar, 1552)
(Apr, 1324)
(May, 1415)
(Jun, 1619)
(Jul, 1360)
(Aug, 1230)
(Sept, 1152)
(Oct, 1058)
(Nov, 909)
(Dec, 1082)
Command took 15.16 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:03:47 PM on Assessment
```

The Clinical RDD Clean process maps the year of clinical trial data to the fileroot variable. It removes quotes and commas from values, associates month abbreviations with numerical values, and locates the Clinical RDD_Clean RDD's index for the 'Completion' column. The primary processing involves applying the Clinical RDD Clean filter, extracting the month and year from the 'Completion' column, filtering the data to include records from 2023, and mapping month values to abbreviations. The function reduces and sorts the findings in ascending order, and the resultant (month, count) pairs are printed.

## DATAFRAME

```python
# Question 5:PLOTTING THE NUMBER OF COMPLETED STUDIES EACH MONTH IN A GIVEN YEAR

from pyspark.sql.functions import split, regexp_replace, when

# Rename columns to remove any leading or trailing commas and quotes
for col in clinical_dataframe.columns:
    clinical_dataframe = clinical_dataframe.withColumnRenamed(col, col.strip(",").strip('"'))

# Define a dictionary to map numerical month values to month names
month_names = {
    "01": "January", "02": "February", "03": "March", "04": "April", "05": "May", "06": "June",
    "07": "July", "08": "August", "09": "September", "10": "October", "11": "November", "12": "December"
}

# Extract year and month from the 'Completion' column and clean month format
completed_cd = clinical_dataframe \
    .withColumn('Year', split('Completion', "-")[0]) \
    .withColumn('Month', split('Completion', "-")[1]) \
    .withColumn('Month', regexp_replace("Month", ",", "")) \
    .withColumn('Month', regexp_replace("Month", '"', "")) \
    .filter(clinical_dataframe.Status.isin(["COMPLETED"])) \
    .select("Month", "Year", "Status")

# Filter for the year 2023 and group by month
completed_cd_2023 = completed_cd.filter(completed_cd.Year.isin(["2023"])) \
    .groupBy("Month").count().orderBy("Month", ascending=True)

# Map numerical month values to month names
completed_cd_2023 = completed_cd_2023.withColumn("Month",
                    when(completed_cd_2023["Month"] == "01", "January")
                    .when(completed_cd_2023["Month"] == "02", "February")
                    .when(completed_cd_2023["Month"] == "03", "March")
                    .when(completed_cd_2023["Month"] == "04", "April")
                    .when(completed_cd_2023["Month"] == "05", "May")
                    .when(completed_cd_2023["Month"] == "06", "June")
                    .when(completed_cd_2023["Month"] == "07", "July")
                    .when(completed_cd_2023["Month"] == "08", "August")
                    .when(completed_cd_2023["Month"] == "09", "September")
                    .when(completed_cd_2023["Month"] == "10", "October")
                    .when(completed_cd_2023["Month"] == "11", "November")
                    .when(completed_cd_2023["Month"] == "12", "December"))

# Show the result
completed_cd_2023.show()
```

▶ (2) Spark Jobs
▶ ☰ clinical_dataframe: 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]
▶ ☰ completed_cd_2023: pyspark.sql.dataframe.DataFrame = [Month: string, count: long]

```
+---------+-----+
|    Month|count|
+---------+-----+
|  January| 1494|
| February| 1272|
|    March| 1552|
|    April| 1324|
|      May| 1415|
|     June| 1619|
|     July| 1360|
|   August| 1230|
|September| 1152|
|  October| 1058|
| November|  909|
| December| 1082|
+---------+-----+
```

Command took 12.28 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment

The query extracts year and month values from the 'Completion' column using the SUBSTRING function. The month name is assigned to the MonthName column. The total number of completed studies for each month is counted. The data is filtered to include records from 2023. Results are grouped by year, month, and month name. The results are sorted by month in ascending order. The table displays year, month, month name, and finalized_studies.

## SQL

```sql
1   %sql
2
3   SELECT SUBSTRING(Completion, 1, 4) AS Year, SUBSTRING(Completion, 6, 2) as Month,
4   CASE SUBSTRING(Completion, 6, 2)
5       WHEN '01' THEN 'Jan'
6       WHEN '02' THEN 'Feb'
7       WHEN '03' THEN 'Mar'
8       WHEN '04' THEN 'Apr'
9       WHEN '05' THEN 'May'
10      WHEN '06' THEN 'Jun'
11      WHEN '07' THEN 'Jul'
12      WHEN '08' THEN 'Aug'
13      WHEN '09' THEN 'Sep'
14      WHEN '10' THEN 'Oct'
15      WHEN '11' THEN 'Nov'
16      WHEN '12' THEN 'Dec'
17    END AS MonthName,
18   COUNT(Status) as Completed_Studies
19   FROM clinicaltrial_2023
20   WHERE SUBSTRING(Completion, 1, 4) = 2023 AND Status = 'COMPLETED'
21   GROUP BY 1, 2, 3
22   ORDER BY 2 ASC ;
```

▶ (2) Spark Jobs
▶ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [Year: string, Month: string … 2 more fields]

| | Year | Month | MonthName | Completed_Studies |
|---|---|---|---|---|
| 1 | 2023 | 01 | Jan | 1494 |
| 2 | 2023 | 02 | Feb | 1272 |
| 3 | 2023 | 03 | Mar | 1552 |
| 4 | 2023 | 04 | Apr | 1324 |
| 5 | 2023 | 05 | May | 1415 |
| 6 | 2023 | 06 | Jun | 1619 |
| 7 | 2023 | 07 | Jul | 1360 |

⬇ 12 rows | 12.96 seconds runtime                    Refreshed 18 minutes ago

Command took 12.96 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment

The query extracts year and month values from the 'Completion' column using the SUBSTRING function. The month name is assigned to the MonthName column. The total number of completed studies for each month is counted. The data is filtered to include records from 2023. Results are grouped by year, month, and month name. The results are sorted by month in ascending order. The table displays year, month, month name, and finalized_studies.

## Discussion of Result

```python
import matplotlib.pyplot as plt

# Assuming completed_studies is the RDD containing the result
data = completed_studies.collect()  # Collect the data from RDD

# Separate months and counts from the collected data
months = [item[0] for item in data]
counts = [item[1] for item in data]

# Create a line chart
plt.figure(figsize=(10, 6))
plt.plot(months, counts, marker='o', linestyle='-')

# Add data points on the line
for x, y in zip(months, counts):
    plt.text(x, y, f'{y}', ha='right', va='bottom', fontsize=10)

# Add title and labels
plt.title('Number of Completed Studies Over Time (2023)')
plt.xlabel('Month')
plt.ylabel('Number of Completed Studies')

# Rotate x-axis labels for better readability
plt.xticks(rotation=45)

# Show plot
plt.tight_layout()
plt.grid(True)  # Add grid for better visualization
plt.show()
```

▸ (1) Spark Jobs



Command took 1.37 seconds -- by o.o.kolawole@edu.salford.ac.uk at 4/30/2024, 4:42:50 PM on Assessment

The code creates a line chart using **Phyton's Matplotlib** to visualize the number of completed studies over time. It assumes `completed_studies` is an RDD containing the result and collects data into lists. The chart is created using `plt.plot()` with `months` as the x-axis and `counts` as the y-axis. Data points are added, titles, x-labels, and y-labels set, and the x-axis labels rotated for better readability.

From the result June has the highest number of completed studies of 1,619 clinical trials and the month has the lowest number of completed studies of November is 909 clinical trials. This can be due to factors such as seasonal trends, study duration, funding availability, study duration etc.

**FURTHER ANALYSIS**

**QUESTION 6**

The top Status (from Status) with their frequencies.

**RDD**

**Assumptions**

The "Status" column in the "clinical trial 2023" contains multiple status values based on the level the trial has progressed.

```python
# FURTHER ANALYSIS 6
# TOP 5 CLINICAL TRIAL STATUS WITH THEIR FREQUENCIES

Status_delimeter = {
    "clinicaltrial_2023": "\t",
    "clinicaltrial_2021": ",",
    "clinicaltrial_2020": ",",
}
Status_column_index = Clinical_RDD_Clean.first().index('Status')

Clinical_RDD_Clean.flatMap(lambda x: x[Status_column_index].split(Status_delimeter[fileroot])).filter(lambda row: row != 'Status').filter(lambda row: row != '').map(lambda x: (x, 1)).reduceByKey(lambda a,b: a + b).sortBy(lambda x: x[1], ascending=False).take(5)
```

▸ (5) Spark Jobs

Out[25]: [('COMPLETED', 263498),
 ('RECRUITING', 66158),
 ('UNKNOWN', 64813),
 ('TERMINATED', 28022),
 ('NOT_YET_RECRUITING', 20098)]

Command took 14.03 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 10:48:24 PM on Assessment

In this code, the index of the 'Status' column was derived from the header of the RDD then the column is split for each row using the appropriate delimiter. The header and empty rows were filtered out then the remaining rows was counted by the occurrence of each status in descending order. The top 15 most frequent status of the sorted RDD was then displayed.

**Discussion of result**



Top 5 Statuses with Their Frequencies

The result shows the distribution of different statuses within the clinical trial dataset, along with their respective frequencies. A considerable portion of the trials in the dataset have been completed. These trials have finished their data collection, analysis, and reporting phases.

**QUESTION 7**

Find the 10 most common sponsors that are pharmaceutical companies, along with the number of clinical trials they have sponsored.

**DATA FRAME**

**Assumptions**

- Organizations listed in the "Parent Company" column of the "pharma" dataset are assumed to be pharmaceutical companies.

- The "Sponsor " column in the "clinicaltrial_2023" comprises the names of companies that are both pharmaceutical and non-pharmaceutical.

```python
# QUESTION 7: Companies that are pharmaceutical companies
from pyspark.sql.functions import col

# Retrieve the list of pharmaceutical companies
pharma_list = create_dataframe(pharma).select("Parent_Company").rdd.flatMap(lambda x: x).collect()

# Filter out pharmaceutical companies from the clinical trial sponsors
pharma_sponsors = clinical_dataframe.select("Sponsor").groupBy("Sponsor").count().filter(col("Sponsor").isin(pharma_list)).orderBy("count", ascending=False).limit(10)

# Show the top 10 pharmaceutical sponsors
pharma_sponsors.show()
```

▸ (4) Spark Jobs

▸ ▦ pharma_sponsors: pyspark.sql.dataframe.DataFrame = [Sponsor: string, count: long]

```
+--------------------+-----+
|             Sponsor|count|
+--------------------+-----+
|     GlaxoSmithKline| 3482|
|              Pfizer| 3045|
|          AstraZeneca| 3024|
|Boehringer Ingelheim| 2146|
|              Sanofi| 1404|
|Bristol-Myers Squibb| 1383|
|               Amgen|  851|
|              AbbVie|  728|
|            Novartis|  697|
|     Gilead Sciences|  625|
+--------------------+-----+
```

Command took 19.49 seconds -- by o.o.kolawole@edu.salford.ac.uk at 4/21/2024, 11:07:25 PM on Assessment

A pharma_list is created from the pharma dataset, containing unique parent company names. The dataframe is filtered by the sponsor column, ensuring only sponsors on the pharma list are included, and displayed in descending order.

**Discussion of result**

The top -pharmaceutical company sponsoring clinical trials is GlaxoSmithKline having a frequency of 3482.

The output can be used to compare the sponsorship efforts of various pharmaceutical companies or to assess how dominant pharmaceutical companies are in the field of clinical trials.

**QUESTION 8**

Find the top offense commited by the pharmaceutical companies and frequencies

**SQL**

**Assumption**

The offense column is in the pharma file.

```sql
%sql
SELECT offense_group, COUNT(*) AS count
FROM pharma
GROUP BY offense_group
ORDER BY count DESC;
```

▶ (2) Spark Jobs

▶ ▦ _sqldf: pyspark.sql.dataframe.DataFrame = [offense_group: string, count: long]

Table ∨ +

| | offense_group | count |
|---|---|---|
| 1 | government-contracting-related offenses | 281 |
| 2 | environment-related offenses | 202 |
| 3 | safety-related offenses | 128 |
| 4 | healthcare-related offenses | 128 |
| 5 | competition-related offenses | 103 |
| 6 | employment-related offenses | 72 |
| 7 | consumer-protection-related offenses | 29 |

⬇ 9 rows | 1.00 second runtime                                    Refreshed 19 minutes ago

ⓘ SQL cell result stored as PySpark data frame _sqldf . Learn more

Command took 1.00 second -- by o.e.kolawole@edu.salford.ac.uk at 4/21/2024, 11:56:28 PM on Assessment

The query will group the pharma file by its offense group, count the number for each offense, order them by count in descending order, and show its result using the 'FROM' and 'SELECT' Statements.

**Discussion of result**

The most common offense the pharmaceutical companies commit is "government contracting related offenses" having a frequency of 281 occurrences.