

# **BIG DATA TOOLS AND TECHNIQUES**

**STUDENT ID: @00704794**

**MSC DATA SCIENCE  
School of Computing, Science & Engineering  
University of Salford, Manchester.**

**KOSEMANI ABIOLA M.  
APRIL 11, 2023**

# **TASK ONE**

## **ABSTRACT.**

This task is aimed at providing a broad overview of the general area of big data systems through developing specific knowledge in areas demonstrating an interaction and synergy between ongoing research and practical use.

The use of common Big Data tools and techniques, the ability to perform a standard data analysis process which includes loading information, cleansing it, analysing it, and reporting with Python, SQL or Linux terminal commands are essential skills that will be evaluated in this evaluation.

## **INTRODUCTION**

I'm going to use clinical trial data in this work and combine it with the list of pharmaceutical companies. To verify the core solution to these problems, I was provided with answers on two main implementations from historic data sets.

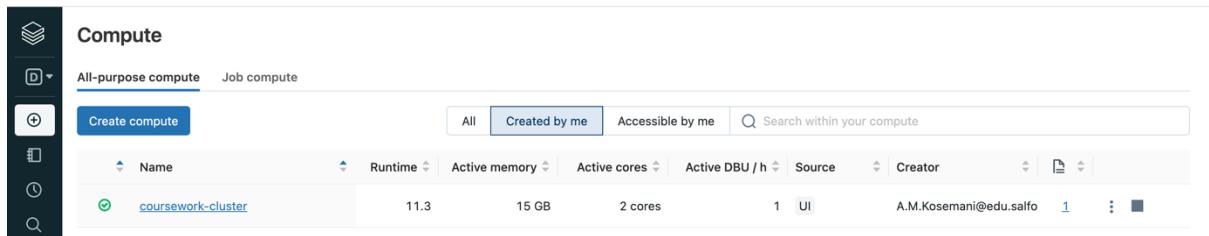
The dataset used for this assessment was downloaded as .csv files, The .csv files have a header describing the file's contents. These csv files include. Clinicaltrial\_2021.csv, Clinicaltrial\_2019.csv, Clinicaltrial2020.csv. Each row in the Clinicaltrial files represent an individual clinical trial, identified by an ID, listing the sponsor, the status of the study and other necessary details.

The pharma.csv file contains a small number if a publicly available list of pharmaceutical violations. For the purpose of this work, the assessment focus is on the parent company which contains the name of the pharmaceutical company in question

I am to implement three steps three times once in Spark SQL and twice in PySpark for RDD and DataFrame, For the visualisation of the results, I decided to use a tool of my choice which fulfil the assessment requirements.

## DATA CLEANING AND PREPARATION

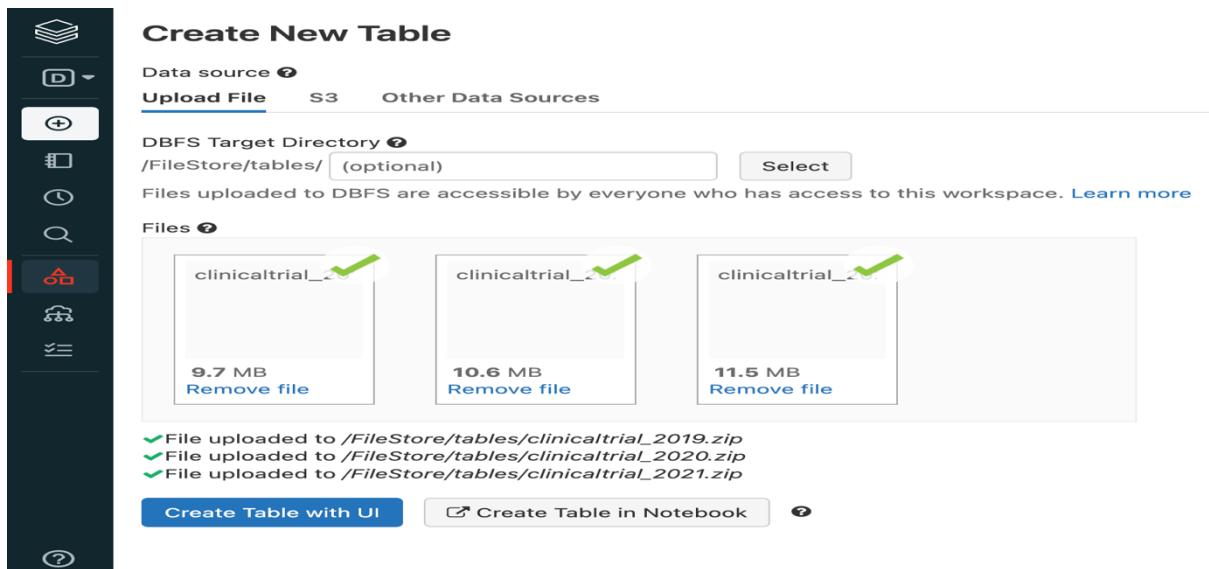
The Databricks compute tool has been created using the Create tab, it provides a powerful and scalable platform for data processing and analysis. While several computes have been created during the process of solving this evaluation due to a timeout caused by inactivity, the following example shows how compute has been created.



A screenshot of the Databricks Compute tab. The sidebar on the left has icons for clusters, jobs, notebooks, and search. The main area is titled 'Compute' with tabs for 'All-purpose compute' (selected) and 'Job compute'. A 'Create compute' button is at the top. Below it is a table with columns: Name, Runtime, Active memory, Active cores, Active DBU / h, Source, Creator, and three more columns. One row is shown with the name 'coursework-cluster', runtime 11.3, active memory 15 GB, 2 cores, 1 UI, and creator 'A.M.Kosemani@edu.salf.ac.uk'. There are also icons for edit, delete, and more options.

**Fig 1.0 Creating a compute**

Uploading of datasets was done by using graphical user's graphical interphase(GUI) tab of the Databricks File System(DBFS) and the imported files was saved in the FileStore/tables directory. A screenshot of this process is shown below.



A screenshot of the Databricks Create New Table tab. The sidebar on the left has icons for clusters, notebooks, and help. The main area is titled 'Create New Table' with tabs for 'Data source' (selected), 'Upload File' (selected), 'S3', and 'Other Data Sources'. It shows 'DBFS Target Directory' as '/FileStore/tables/'. Below that is a 'Select' button and a note about files being accessible to everyone. Under 'Files', there are three boxes for 'clinicaltrial\_2019.zip', 'clinicaltrial\_2020.zip', and 'clinicaltrial\_2021.zip', each with a green checkmark and a 'Remove file' link. Below the files is a list of successful uploads: 'File uploaded to /FileStore/tables/clinicaltrial\_2019.zip', 'File uploaded to /FileStore/tables/clinicaltrial\_2020.zip', and 'File uploaded to /FileStore/tables/clinicaltrial\_2021.zip'. At the bottom are buttons for 'Create Table with UI' (selected) and 'Create Table in Notebook'.

**Fig 1.1 uploading clinicaltrial .csv files using the Databricks create tab GUI.**

## LIBRARIES IMPORTATION AND FILE UNZIPPING

After uploading the 3 clinicaltrial .csv files into the FileStore/tables directory, some series of other operation were required for the process of unzipping the files, but first I was able to import the required packages which are needed for this project which is shown below



The screenshot shows a Jupyter Notebook interface with two code cells. The first cell, titled 'Setup', contains the command 'Load libraries'. The second cell contains the following Python code:

```
1 import csv
2 from operator import add
3 from matplotlib import pyplot as plt
4 from pyspark.sql.session import SparkSession
5
6
7 spark = SparkSession(sc)
```

Below the code cell, a message indicates the command took 0.11 seconds and was run by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:38 on course-work.

**Fig 1.2 Libraries setup and importation**

The first line import csv module and the following like imports the add operator from the operator module, both modules are responsible for reading and manipulating data.

The third line imports the pyplot module from the matplotlib while the fourth line created a SparkSession object, the pyplot is a data visualization library and will be used to create graphs and charts in this assessment. The SparkSession serves as an entry point to spark functionality. Hence, the “sc” param is a SparkContext which is purpose is to set up connection to the spark cluster.

After the library’s importation, the imported files were then moved from FileStore/tables to what is called the driver node which is the /tmp directory this is because the FileStore/tables doesn’t support files unzipping. For this reason, it has no command for file unzipping.

```

Cmd 3
1 ## list the files and directories located in the "tables" directory within the "FileStore" directory.
2
3 dbutils.fs.ls("FileStore/tables")

Out[183]: [FileInfo(path='dbfs:/FileStore/tables/Occupancy_Detection_Data.csv', name='Occupancy_Detection_Data.csv', size=50968, modificationTime=1677692453000),
 FileInfo(path='dbfs:/FileStore/tables/account-models/', name='account-models/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/accounts/', name='accounts/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/accounts-1.zip', name='accounts-1.zip', size=5297592, modificationTime=1675271374000),
 FileInfo(path='dbfs:/FileStore/tables/activations/', name='activations/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/activations.zip', name='activations.zip', size=8411369, modificationTime=1675269262000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2019.csv', name='clinicaltrial_2019.csv', size=42400056, modificationTime=168259854100),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2019.zip', name='clinicaltrial_2019.zip', size=9707871, modificationTime=1682569150000),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2020.csv', name='clinicaltrial_2020.csv', size=46318151, modificationTime=168259854300),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2020.zip', name='clinicaltrial_2020.zip', size=10599182, modificationTime=168256915800),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2021.csv', name='clinicaltrial_2021.csv', size=50359696, modificationTime=168260322400),
 FileInfo(path='dbfs:/FileStore/tables/clinicaltrial_2021.zip', name='clinicaltrial_2021.zip', size=11508457, modificationTime=168256918900),
 FileInfo(path='dbfs:/FileStore/tables/logs/', name='logs/', size=0, modificationTime=0),
 FileInfo(path='dbfs:/FileStore/tables/logs-1.zip', name='logs-1.zip', size=18168065, modificationTime=1675272151000),
 FileInfo(path='dbfs:/FileStore/tables/logs-2.zip', name='logs-2.zip', size=18168065, modificationTime=1675273001000).
Command took 0.22 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:42 on course-work

```

**Fig 1.3 Checking the for the imported files with the dbutils.fs.ls() command**

```

1 ## copy the clinicaltrail_2019.zip file to file:/tmp/ for unzipping
2
3 dbutils.fs.cp("dbfs:/FileStore/tables/clinicaltrial_2019.zip", "file:/tmp/")

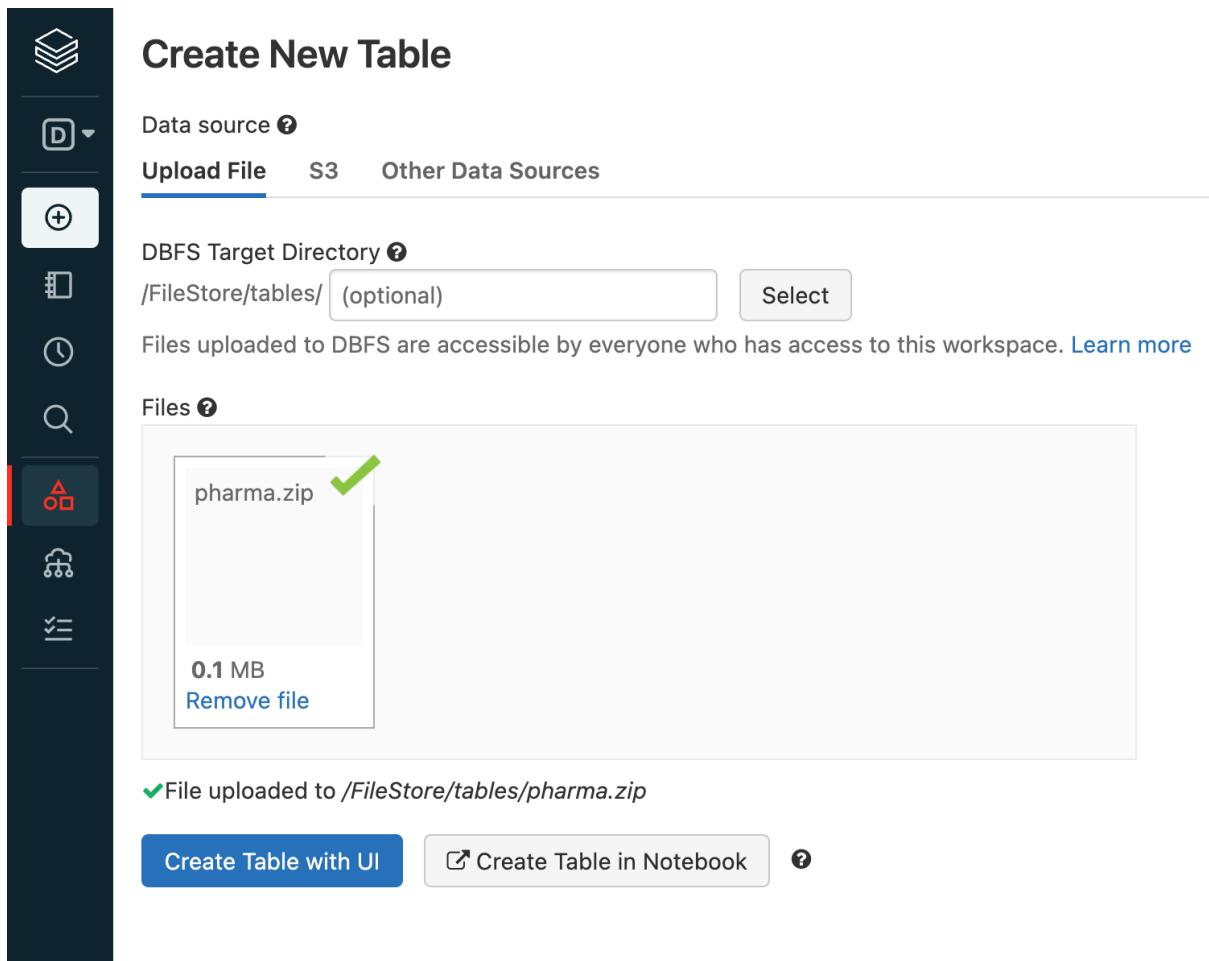
Out[184]: True
Command took 0.55 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:45 on course-work
Cmd 5
1 ## copy the clinicaltrail_2020.zip file to file:/tmp/ for unzipping
2
3 dbutils.fs.cp("dbfs:/FileStore/tables/clinicaltrial_2020.zip", "file:/tmp/")

Out[185]: True
Command took 0.47 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:48 on course-work
Cmd 6
1 ## copy the clinicaltrail_2021.zip file to file:/tmp/ for unzipping
2
3 dbutils.fs.cp("dbfs:/FileStore/tables/clinicaltrial_2021.zip", "file:/tmp/")

Out[186]: True
Command took 0.47 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:52 on course-work

```

**Fig 1.4 Coping the .zip files from FileStore/tables to /tmp dir.**



**Fig 1.5 uploading pharma.csv file using the Databricks create tab GUI.**

```
Cmd 7
1 ## list the files and directories located in the "tables" directory within the "FileStore" directory.
2
3 dbutils.fs.ls("FileStore/tables")

Cmd 8
1 ## copy the pharma.zip file to file:/tmp/ for unzipping
2
3 dbutils.fs.cp("dbfs:/FileStore/tables/pharma.zip", "file:/tmp/")

Out[188]: True
Command took 0.23 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:59 on course-work
```

**Fig 1.6**

The ls command checks the **FileStore/tables** dir to confirm that pharma.csv file has been uploaded, and the UNIX “cp” command moves this file to \tmp directory for unzipping

```
Cmd 10
1
2 %sh
3 ls /tmp/
Rserv
RtmpgWL5it
chauffeur-daemon-params
chauffeur-daemon.pid
chauffeur-env.sh
clinicaltrial_2019.csv
clinicaltrial_2019.zip
clinicaltrial_2020.csv
clinicaltrial_2020.zip
clinicaltrial_2021.zip
custom-spark.conf
driver-daemon-params
driver-daemon.pid
driver-env.sh
hsperfdata_root
pharma.csv
pharma.zip
systemd-private-9d5d5cdef0574a1d8af7fb2d9839932-apache2.service-ud3oDi
systemd-private-9d5d5cdef0574a1d8af7fb2d9839932-ntp.service-SLkpih
systemd-private-9d5d5cdef0574a1d8af7fb2d9839932-systemd-logind.service-jsYjZe
systemd-private-9d5d5cdef0574a1d8af7fb2d9839932-systemd-resolved.service-2wdY9f
Command took 0.13 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:52:09 on course-work
```

**Fig 1.7 Checking the /tmp directory for the moved .zip files.**

```
Cmd 11
1 %sh
2
3 unzip -d /tmp/ /tmp/clinicaltrial_2019.zip
4 unzip -d /tmp/ /tmp/clinicaltrial_2020.zip
5 unzip -d /tmp/ /tmp/clinicaltrial_2021.zip
6 unzip -d /tmp/ /tmp/pharma.zip
Archive: /tmp/clinicaltrial_2019.zip
  inflating: /tmp/clinicaltrial_2019.csv
Archive: /tmp/clinicaltrial_2020.zip
  inflating: /tmp/clinicaltrial_2020.csv
Archive: /tmp/clinicaltrial_2021.zip
  inflating: /tmp/clinicaltrial_2021.csv
Archive: /tmp/pharma.zip
  inflating: /tmp/pharma.csv
Command took 1.40 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:56:12 on course-work
```

**Fig 1.8 unzipping the .zip files**

The “**unzip**” command is used to unzip files which are in the zip format, the “**-d**” which is followed by the dir path is used to specify the directory the extracted files should be placed.

The above image in **Fig 1.8** extracts the contents of the .zip files into the **/tmp/** directory.



```
Cmd 12
1 ## Create tables in the FileStore directory
2
3 dbutils.fs.mkdirs("FileStore/tables")

Out[192]: True
Command took 0.15 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:53:04 on course-work
```

Cmd 12

Python

## Fig 1.9 Creating tables directory in the FileStore.

A directory named "tables" in the FileStore folder is created in **Fig 1.9**



```
Cmd 15
1 ## ## moving clinicaltrial_2021.csv file from /tmp directory to /tables directory
2
3 dbutils.fs.mv("file:/tmp/clinicaltrial_2021.csv", "FileStore/tables/clinicaltrial_2021.csv", True)
4

Out[222]: True
Command took 2.30 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:56:29 on course-work
```

```
Cmd 16
1 ## moving pharma.csv file from /tmp directory to /tables directory
2
3 dbutils.fs.mv("file:/tmp/pharma.csv", "FileStore/tables/pharma.csv", True)

Out[196]: True
Command took 0.25 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:53:25 on course-work
```

## Fig 2.0 Moving the unzipped .csv files to the FileStore/tables

In Fig 2.0, the dbutils module is used to move a file from a location file path to another. This dbutils.fs.mv method is used in the Fig above to move a file within the Databricks file system (DBFS). The ‘true’ option which is provided in the argument in the method which means to move the subdirectories and files under the source directory is a recurse statement indicating whether to move files recursively.

```

Cmd 18
1 ## extracting the filename from the given path and append it to a directory path.
2
3 id = clinicaltrial_2021.rfind('/')
4 FileStoreTableDir = 'FileStore/tables'
5 clinicalTrial = FileStoreTableDir+clinicaltrial_2021[id:]

Command took 0.09 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:51:51 Insert a new cell -work

```

**Fig 2.1 extracting the filename from the given path to append it to a directory path.**

```

Cmd 20
1 ## reading the csv file text with sc.textFile to return an RDD
2
3 clinicaltrial_2021 = sc.textFile(clinicaltrial_2021)\n
4 .map(lambda line: line.split('|'))\n
5 .filter(lambda line: len(line)>1)
6
7 header = clinicaltrial_2021.first()
8 clinicaltrial_2021 = clinicaltrial_2021.filter(lambda x: x!= header)

▶ (1) Spark Jobs
Command took 1.06 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:53:40 on course-work

```

**Fig 2.2 reading the csv file text with sc.textFile to return an RDD**

It is safe to assume that the text file contains a “|” delimiter and a header as the first row, and that the header has more than one element, if these assumptions are not met for some reasons this code will produce error or an unexpected result.

In **Fig 2.2** the clinicaltrial\_2021 file was read and splits each line by a delimiter which filters out the lines with only one element. The header is stored separately in the header variable which is filtered out of the dataset.

The result of this operation is a Resilient Distributed Dataset object which called RDD.

```

1 #return all the rows in the DataFrame
2
3 clinicaltrial_2021.collect()

▶ (1) Spark Jobs
Out[201]: [[{'NCT02758028',
  'The University of Hong Kong',
  'Recruiting',
  'Aug 2001',
  'Aug 2001',
  'Interventional',
  'Apr 2016',
  '',
  ''},
  {'NCT02751957',
  'Duke University',
  'Completed',
  'Jul 2016',
  'Jul 2020',
  'Interventional',
  'Apr 2016',
  'Autistic Disorder, Autism Spectrum Disorder',
  ''},
  {'NCT02758483',
  'Universidade Federal do Rio de Janeiro',
  'Completed',
  ''}]

Command took 5.14 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:53:44 on course-work

```

**Fig 2.3, using collect method to return all rows in the DataFrame**

## PROBLEMS ANSWERS IMPLEMENTATION(PySpark RRD).

### SOLUTION 1.



```
QUESTION 1
```

```
1 ## Numbers of studies in a dataset
2
3 clinicaltrial_2021.count()
```

▶ (1) Spark Jobs  
Out[203]: 387261  
Command took 2.17 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:00 on course-work

**Fig 2.4 Numbers of count in a dataset.**

The number of studies was obtained by using the count method. In essence, it returns the numbers of rows in an entire dataset that correspond to the number of different IDs within a given Clinicaltrial Dataset. Therefore, the number of studies can be seen in the screenshot above for the clinicaltrial\_2021 dataset with a total number of distinct counts is 387261.

### SOLUTION 2.



```
QUESTION 2
```

```
1 header
```

```
Out[204]: ['Id',
'Sponsor',
>Status',
'Start',
'Completion',
'Type',
'Submission',
'Conditions',
'Interventions']
```

Command took 0.12 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:07 on course-work

**Fig 2.5 Checking the file header's values.**

```

Cmd 27
1 # the value of type is 5
2
3 Type = 5

Command took 0.02 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:11 on course-work

Cmd 28
1 clinicaltrial_2021.map(lambda col: col[Type])\
2 .filter(lambda x: x!=None)\ 
3 .map(lambda x: (x, 1))\ 
4 .reduceByKey(add)\ 
5 .map(lambda x: (x[1], x[0]))\ 
6 .sortByKey(ascending=False)\ 
7 .map(lambda x: (x[1], x[0])).collect()

▶ (3) Spark Jobs
Out[206]: [('Interventional', 301472),
('Observational', 77540),
('Observational [Patient Registry]', 8180),
('Expanded Access', 69)]

Command took 3.16 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:14 on course-work

Cmd 29

```

**Fig 2.6**

As shown in Fig 2.5, the first step is to check for the value of Type in the header's Array, with the value of an array read from zero, the value of Type is equals 5, After which a map function is used to select the Type column in the clinicaltrial-2021 file for all the years rows and observations under the column, any missing observation in the column was filtered out. Aggregation was done which returns an output in a random order. The output was then sorted out in ascending order which a collection method invoked to display the result.

## SOLUTION 3

```

QUESTION 3
Markdown ⓘ ▾ ×

Cmd 30
1 Conditions = 7

Command took 0.11 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:21 on course-work

Cmd 31
1 clinicaltrial_2021.map(lambda col: col[Conditions])\
2 .filter(lambda x: x!=None)\ 
3 .map(lambda x: x.split(','))\ 
4 .flatMap(lambda x: x.collect()


```

**Fig 2.7**

```

CMD 34
Five top condition in the dataset

Cmd 33
1 clinicaltrial_2021.map(lambda col: col[Conditions])\
2 .filter(lambda x: x!= '')\
3 .flatMap(lambda x: x.split(','))\
4 .map(lambda x: (x, 1))\
5 .reduceByKey(add)\
6 .sortBy(lambda x: x[1], ascending=False)\
7 .take(5)

▶ (3) Spark Jobs
Out[209]: [('Carcinoma', 13389),
('Diabetes Mellitus', 11080),
('Neoplasms', 9371),
('Breast Neoplasms', 8640),
('Syndrome', 8032)]
Command took 3.54 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:32 on course-work
CMD 34

```

## Fig 2.8 The five top conditions in the dataset.

The assumption was that the column condition is not empty or null and that there are rows in a column with more than one conditions which require separation to get an output

The value of the conditions was confirmed by checking the header as show in **Fig 2.5**, A map function was then used to select the conditions column inside the clinicaltrial file which returns every observation or rows under conditions. Filtering was done to take care of the missing observations. There are rows with more than one condition and a split was needed to separate the conditions rows.

Finally, to count the number of different conditions in each row under conditions column a map function was used, the resulting count was now aggregated which was sorted sing the **sortByKey** function as seen in **Fig 2.8**

## SOLUTION 4

```

QUESTION 4

CMD 35
pharma = sc.textFile(pharma) \
.map(lambda line: line.split(',')) \
.filter(lambda line: len(line)>1)

pharma_header = pharma.first()
pharma = pharma.filter(lambda x: x!= header)

▶ (1) Spark Jobs
Command took 0.34 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:40 on course-work
CMD 36

```

## Fig 2.9 reading the .csv file text with sc.textFile

```

1 pharma_header

Out[211]: ['Company',
 'Parent_Company',
 'Penalty_Amount',
 'Subtraction_From_Penalty',
 'Penalty_Amount_Adjusted_For_Eliminating_Multiple_Counting',
 'Penalty_Year',
 'Penalty_Date',
 'Offense_Group',
 'Primary_Offense',
 'Secondary_Offense',
 'Description',
 'Level_of_Government',
 'Action_Type',
 'Agency',
 'Civil/Criminal',
 'Prosecution_Agreement',
 'Court',
 'Case_ID',
 'Private_Litigation_Case_Title',
 'Lawsuit_Resolution',
 'Facility_State']

Command took 0.13 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:43 on course-work

```

**Fig 3.0 Checking the pharma's columns.**

```

Cmd 37
1 Parent_Company = 1
2 Sponsor = 1

Command took 0.04 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:47 on course-work

```

**Fig 3.1 giving a value position to parent company and sponsor column**

```

Cmd 38
1 sponsor = clinicaltrial_2021.map(lambda col: col[Sponsor])
2 pharmaceutical = pharma.map(lambda col: col[Parent_Company])\
3 .map(lambda x: x.strip('\"'))
4
5 pharmaceutical_rdd = sc.broadcast(pharmaceutical.collect())
6
7 sponsor.filter(lambda x: x not in pharmaceutical_rdd.value) \
8 .map(lambda x: (x, 1)) \
9 .reduceByKey(add) \
10 .sortBy(lambda x: x[1], ascending = False) \
11 .take(10)

▶ (4) Spark Jobs

Out[213]: [('National Cancer Institute (NCI)', 3218),
 ('M.D. Anderson Cancer Center', 2414),
 ('Assistance Publique - Hôpitaux de Paris', 2369),
 ('Mayo Clinic', 2300),
 ('Merck Sharp & Dohme Corp.', 2243),
 ('Assiut University', 2154),
 ('Novartis Pharmaceuticals', 2088),
 ('Massachusetts General Hospital', 1971),
 ('Cairo University', 1928),
 ('Hoffmann-La Roche', 1828)]

Command took 8.81 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:54:54 on course-work

```

**Fig 3.2 ten most common sponsors that are not pharmaceutical companies, along with the number of clinical trials they have sponsored.**

Before carrying out the task, my assumption was that the parent company column should have all the pharmaceuticals, but the not sponsors and the sponsors. Based on **Fig 3.0**, the column number for both parent company and sponsor are 1 according to **Fig 3.1**. In **Fig 3.2**, the RDD created for sponsor was filtered out to get a list of 10 sponsors that are not pharmaceuticals.

## SOLUTION 5

### QUESTION 5

Cmd 40

```
1 year='2021'
2
3 ans = clinicaltrial_2021.map(lambda col: [col[x] for x in [2,4]])\
4 .filter(lambda x: x[0] == 'Completed')\
5 .map(lambda x: x[1])\
6 .filter(lambda x: x!= '')\
7 .map(lambda x: tuple(x.split()))\
8 .filter(lambda x: x[1] == year)\
9 .groupByKey().mapValues(len)\
10 .sortBy(lambda x: x[1], ascending=False)
11
12 ans.collect()
```

► (3) Spark Jobs

```
Out[214]: [('Mar', 1227),
('Jan', 1131),
('Jun', 1094),
('May', 984),
('Apr', 967),
('Feb', 934),
('Jul', 819),
('Aug', 700),
('Sep', 528),
('Oct', 187)]
```

Command took 3.44 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:55:07 on course-work

Cmd 41

```
1 def changer(month):
2     if month[0]=='Jan':
3         return (1,month)
4     elif month[0]=='Feb':
5         return (2,month)
6     elif month[0]=='Mar':
7         return (3, month)
8     elif month[0]=='Apr':
9         return (4, month)
10    elif month[0]=='May':
11        return (5, month)
12    elif month[0]=='Jun':
13        return (6, month)
14    elif month[0]=='Jul':
15        return (7, month)
16    elif month[0]=='Aug':
17        return (8, month)
18    elif month[0]=='Sep':
19        return (9, month)
20    elif month[0]=='Oct':
21        return (10, month)
22    elif month[0]=='Nov':
23        return (11, month)
24    else:
25        return (12,month)
```

Command took 0.08 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:55:15 on course-work

Fig 3.3

Cmd 42

```
1 final = ans.map(changer)\n2 .sortBy(lambda x: x[0], ascending=True)\n3 .map(lambda x: x[1])\n4 final.collect()
```

► (3) Spark Jobs

```
Out[216]: [('Jan', 1131),\n            ('Feb', 934),\n            ('Mar', 1227),\n            ('Apr', 967),\n            ('May', 984),\n            ('Jun', 1094),\n            ('Jul', 819),\n            ('Aug', 700),\n            ('Sep', 528),\n            ('Oct', 187)]
```

Command took 0.47 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:55:19 on course-work

Cmd 43

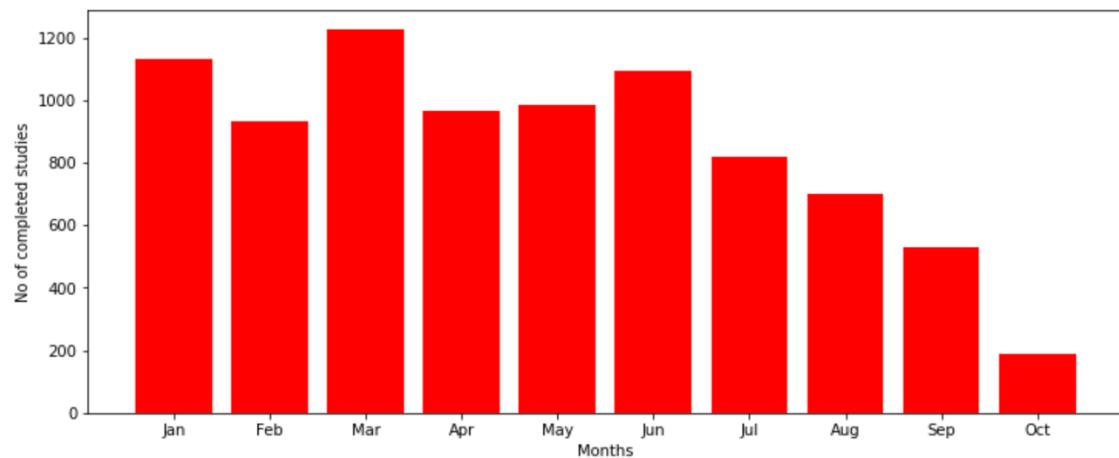
```
1 months = final.map(lambda x: x[0]).collect()\n2 counts =final.map(lambda x: x[1]).collect()
```

► (2) Spark Jobs

Command took 0.31 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:55:24 on course-work

Cmd 44

```
1 plt.figure(figsize=(12, 5))\n2 plt.bar(months, counts, color='red')\n3 plt.xlabel('Months')\n4 plt.ylabel('No of completed studies')\n5 plt.show()\n6
```



Command took 0.37 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:55:28 on course-work

Cmd 45

**Fig 3.4 A graph of no of completed studies against the months.**

In **Fig 3.3**, A RDD was created for completed studies by writing a for loop to loop through and select all studies under the status column and the corresponding date in the completion column then filtering was done in other to filter out the completed and the missing observations and splitting is then done after this to separate the month from year.

The next operation was the filtering by year and sorting the numbers of completed that was save in a variable name, with the result later viewed.

Finally, to have a visualization of the table, the month was then arranged accordingly using a switch if else statement with the visualization done using the matplotlib as shown in **Fig 3.4**.

## EXTRAS.

```
Cmd 45
EXTRA MARK
Comparison between the number of completed studies for 2019, 2020 and 2021

Cmd 46
1   ###Number of Completed Studies 2021
2
3
4 Year_2021 = clinicaltrial_2021.map(lambda col: [col[x] for x in [2,4]])\
5 .filter(lambda x: x[0] == 'Completed')\
6 .map(lambda x: x[1])\
7 .filter(lambda x: x!= '')\
8 .map(lambda x: (x.split()))\
9 .filter(lambda x: x[1] == '2021')\
10 .map(lambda x: x[1])
11
12 Year_2021.count()

> (1) Spark Jobs
Out[219]: 8571
Command took 2.98 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 14:55:33 on course-work

Cmd 47
1   ###Visualization of the comparison of the number of Completed Studies
2   ##between 2019, 2020, 2021.
3
4 Completed_Studies = ['Year_2019', 'Year_2020', 'Year_2021']
5 count = [17941, 16118, 8571]
6
7 plt.bar(x = Completed_Studies, height = count, color = 'pink')
8 plt.xlabel('YEAR')
9 plt.ylabel('Number of Completed Studies')
10 plt.show()



```

**Fig 3.5** The results show a fall in the number of studies that have been finalised during these three years. In the image below, we see a visual representation of this.

## PROBLEMS ANSWERS IMPLEMENTATION(PySpark SQL).

```
cmd _1
1 drop table if exists clinicaltrial_2021;
2
3 create table clinicaltrial_2021
4 using csv
5 options(
6 header='true',
7 delimiter='|',
8 inferSchema='true',
9 mode='FAILFAST',
10 path='/FileStore/tables/clinicaltrial_2019.csv'
11 );
12
13 cache table clinicaltrial_2021

▶ (4) Spark Jobs
OK
Command took 11.96 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:17:05 on course-work
```

**Fig 3.5 Creating clinicaltrial\_2021 table.**

This SQL query drop clinicaltrial\_2021 table if it already exists and create a new table clinicaltrial\_2021, the using csv implies csv as the source of data, the options take in arguments such as header, delimiter and whether the schema should be inferred automatically.

The cache table clinicaltrial\_2021 SQL statement caches a table in memory which can improve the performance of the query by reducing the need to read from disk every single time.

```

Cmd 2
1 drop table if exists pharma;
2
3
4 create table pharma
5 using csv
6 options(
7 header='true',
8 delimiter=',',
9 inferSchema='true',
10 mode='PERMISSIVE',
11 path='/FileStore/tables/pharma.csv'
12 );
13
14 cache table pharma

▶ (4) Spark Jobs
OK
Command took 2.44 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:21:52 on course-work

```

**Fig 3.5 Creating pharma's table.**

```

Cmd 3
1 select * from clinicaltrial_2021
SQL ▶ ↻ ⟲ ⟳ ×

▶ (1) Spark Jobs
Table + 



|   | Id          | Sponsor                                | Status                 |
|---|-------------|----------------------------------------|------------------------|
| 1 | NCT02758028 | The University of Hong Kong            | Recruiting             |
| 2 | NCT02751957 | Duke University                        | Completed              |
| 3 | NCT02758483 | Universidade Federal do Rio de Janeiro | Completed              |
| 4 | NCT02759848 | Istanbul Medeniyet University          | Completed              |
| 5 | NCT02758860 | University of Roma La Sapienza         | Active, not recruiting |
| 6 | NCT02757209 | Consorzio Futuro in Ricerca            | Completed              |


↙ ↘ 10,000 rows | Truncated data | 0.61 seconds runtime Refreshed 12 hours ago
Command took 0.61 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:24:19 on course-work

```

**Fig 3.6 Retrieving all columns and rows from clinicaltrial\_2021 table.**

1 | `select * from pharma`

▶ (1) Spark Jobs

**Table** +

	Company	Parent_Company	Penalty_
1	Abbott Laboratories	Abbott Laboratories	\$5,475,0
2	Abbott Laboratories Inc.	AbbVie	\$1,500,0
3	Abbott Laboratories Inc.	AbbVie	\$126,500
4	Abbott Laboratories Puerto Rico, Inc.	Abbott Laboratories	\$49,045
5	Acclarent Inc.	Johnson & Johnson	\$18,000,

↓ 968 rows | 0.48 seconds runtime      Refreshed 12 hours ago

**Fig 3.7 Retrieving all columns and rows from pharma's table.**

## SOLUTION 1.

Cmd 5

**QUESTION 1**

Cmd 6

1 | `select count(*) from clinicaltrial_2021;`

▶ (2) Spark Jobs

**Table** +

	count(1)
1	326348

↓ 1 row | 0.44 seconds runtime      Refreshed 12 hours ago

Command took 0.44 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:23:24 on course-work

**Fig 3.8 Fetching the total number of records in the clinicaltrial\_2021 table.**

The total number of record(rows) in clinicaltrial\_2021 is 326348.

## SOLUTION 2

The screenshot shows a Jupyter Notebook cell with the following content:

```
Cmd 7
QUESTION 2
Cmd 8
1 SELECT Type,
2   COUNT(Type) AS frequency
3   FROM clinicaltrial_2021
4   GROUP BY Type
5   HAVING Type IS NOT NULL
6   ORDER BY frequency DESC;
```

▶ (2) Spark Jobs

Table +

Type	frequency
1 Interventional	255945
2 Observational	64163
3 Observational [Patient Registry]	6171
4 Expanded Access	69

↓ 4 rows | 1.13 seconds runtime Refreshed 12 hours ago

Command took 1.13 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:26:58 on course-work

**Fig 3.9**

This query retrieves data from the clinicaltrial\_2021 table, grouping is done based on the Type column which then counts the number of occurrences of each type. Then ordering was done by the result of frequency occurrence in descending order

The not null value is only included in this result, the null values are filtered out by using where state which creates an exception.

## SOLUTION 3

The screenshot shows a Jupyter Notebook cell with the following content:

```
Cmd 10
QUESTION 3
Cmd 11
1 WITH ct AS (
2   SELECT explode(split(Conditions, ',')) AS separate_condition
3   FROM clinicaltrial_2021
4 )
5 SELECT separate_condition AS condition, COUNT(*) AS Count
6 FROM ct
7 GROUP BY separate_condition
8 ORDER BY Count DESC
9 LIMIT 5;
```

▶ (2) Spark Jobs

Table +

condition	Count
1 Carcinoma	11155
2 Diabetes Mellitus	9830
3 Neoplasms	7815
4 Breast Neoplasms	7486
5 Syndrome	6842

↓ 5 rows | 1.94 seconds runtime Refreshed 12 hours ago

Command took 1.94 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:29:40 on course-work

**Fig 4.0**

In **Fig 4.0**, the query creates a common table expression **CTE** called **ct**, which uses the **WITH** statement. Which uses the explode and split functions to split the values in the conditions column using commas which sperate them into different roles

The other part of this query selects the separate \_condition column from the ct. and uses count to count the number of times each condition appears. The result is then grouped by group by and order using ordered by statement.

This query returns the top 5 conditions which appears mostly in the condition column of the table **clinicaltrial\_2021**.

## SOLUTION 4

QUESTION 4

```
Cmd 12
1 ---create a view to show the sponsors
2
3 create view if not exists SPONSORS as
4 select Sponsor
5   from clinicaltrial_2021;
6
7 select * from SPONSORS
```

SQL ▶️ ↻ ⟲ ⟳ ×

▶ (1) Spark Jobs

Sponsor
1 The University of Hong Kong
2 Duke University
3 Universidade Federal do Rio de Janeiro
4 Istanbul Medeniyet University
5 University of Roma La Sapienza
6 Consorzio Futuro in Ricerca

Table +

10,000 rows | Truncated data | 0.73 seconds runtime

Refreshed 12 hours ago

Command took 0.73 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:42:52 on course-work

**Fig 4.1 Create sponsor's view from the sponsor table.**

```

1 ---create a view for the pharmaceutical companies
2
3 create view if not exists pharmaceutical as
4 select Parent_Company
5 from pharma;
6
7 select * from pharmaceutical
8

```

► (1) Spark Jobs

Parent_Company
1 Abbott Laboratories
2 AbbVie
3 AbbVie
4 Abbott Laboratories
5 Johnson & Johnson
6 Abbott Laboratories

↓ 968 rows | 0.45 seconds runtime Refreshed 12 hours ago

Command took 0.45 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:43:25 on course-work

```

1 ---10 most common sponsors that are not pharmaceuticals
2
3 select Sponsor, count(*) as Count
4 from Sponsors_not_pharmaceuticals
5 where Sponsor is not null
6 group by Sponsor
7 order by Count desc
8 limit 10;

```

SQL ►▼ ↻ — ×

► (3) Spark Jobs

Sponsor	Count
1 National Cancer Institute (NCI)	3003
2 Merck Sharp & Dohme Corp.	2124
3 M.D. Anderson Cancer Center	2097
4 Mayo Clinic	1930
5 Novartis Pharmaceuticals	1881
6 Assistance Publique - Hôpitaux de Paris	1764

↓ 10 rows | 2.48 seconds runtime Refreshed 12 hours ago

Command took 2.48 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:44:31 on course-work

Cmd 16

```

1 ---create a view for sponsors that are not pharmaceutical companies
2
3 create view if not exists Sponsors_not_pharmaceuticals as
4 select * from SPONSORS
5 where Sponsor not in (select Parent_Company from pharmaceutical);
6
7
8 select * from Sponsors_not_pharmaceuticals

```

SQL ►▼ ↻ — ×

► (2) Spark Jobs

Sponsor
1 The University of Hong Kong
2 Duke University
3 Universidade Federal do Rio de Janeiro
4 Istanbul Medeniyet University
5 University of Roma La Sapienza
6 Consorzio Futuro in Ricerca

↓ ▼ 10,000 rows | Truncated data | 1.23 seconds runtime Refreshed 12 hours ago

Command took 1.23 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 18:43:50 on course-work

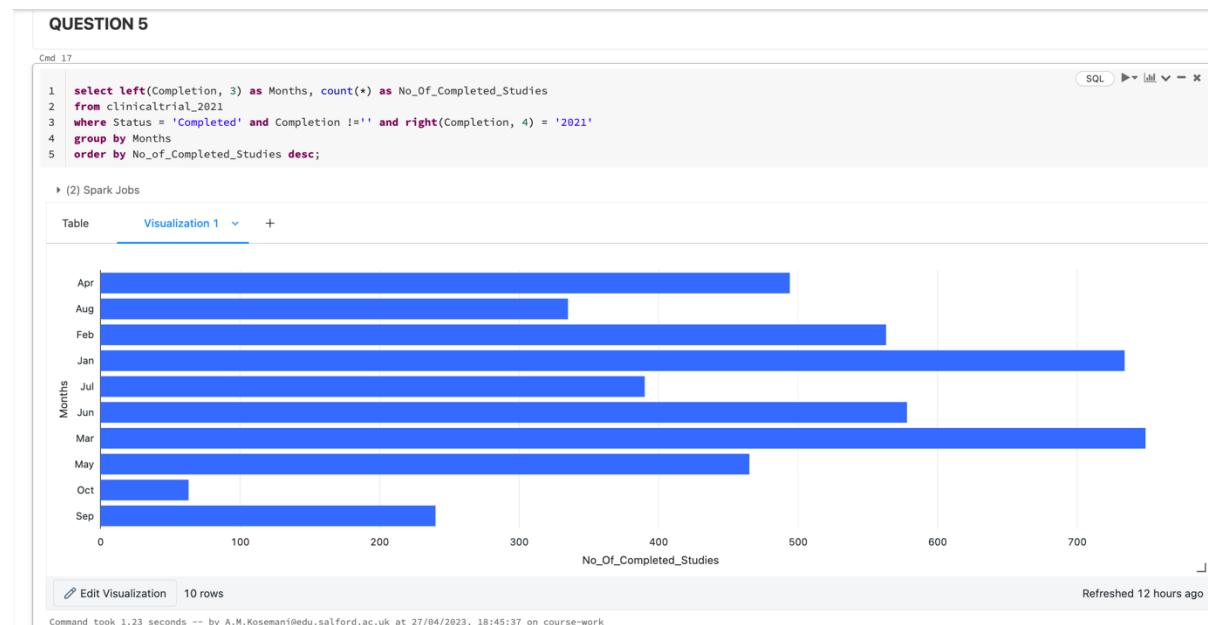
Fig 4.2

From the above result, sponsor was selected from sponsor not pharmaceutical then a count was taken to get an aggregation aliased count, which was then ordered in descending order

**Fig 4.1,** Create a table that alias sponsors consisting of sponsor from the sponsor table **clinicaltrial\_2021**, the same way is used to create pharmaceutical view which creates a view aliased pharmaceutical consisting of all companies from pharma.

A select statement was then used to select sponsors which were not same with parent companies

## SOLUTION 6



**Fig 4.3**

The substring and count functions are used to list the first three elements of a completion date column, while each month's occurrence is counted. Where filter is used to filter to where the status of the study is completed, and the year is 2021. The group by function is also employed due to the use of the count aggregate function.

Using standard Databricks visualisation tools, a graph of the results shown above can be seen below.

## PROBLEMS ANSWERS IMPLEMENTATION(PySpark DataFrame).

### DATA CLEANING AND PREPARATION

Data cleaning is the process of identifying errors or inconsistencies in a dataset and resolving them, if necessary, in order to ensure that the data being worked on is correct, consistent and usable. Although we used the correct format for columns as variables and observations in rows, data was processed dynamically so that it made things a lot easier.

The missing values were also investigated by using the pyspark data frame and it was found that there are 3 columns, namely "Completion", "Outcomes" and "Interventions," which do not have any observations.

```
1 clinicaltrial_2021.select([count(when(isnan(c) | col(c).isNull(), c)).alias(c) for c in clinicaltrial_2021.columns]).show()
```

► (2) Spark Jobs

	Id	Sponsor	Status	Start	Completion	Type	Submission	Conditions	Interventions
	0	0	0	0	13260	0	0	65131	253837

Command took 21.20 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 15:14:00 on course-work

**Fig 4.4**

Fig 4.4 checks for the number of missing observations for all the column by using the select function to loop over each column in the dataset by using the count function to count when either isnan function return as true or when the column is null

The ‘alias’ basically is a way of temporary renaming the column. As we can see above, this result is the same. However, these missed observations have not been deleted because most of the tasks dealt with are specifically for one variable and if this is to be eliminated, it would lead to significant loss of data in the dataset.

Although no general cleaning of the data base has been carried out, missing observations have been removed from certain tasks where necessary. Under the relevant task, such cleaning shall be further explained.

## SOLUTION 1

The number of studies in the dataset.

```
Cmd 15
1 clinicaltrial_2021.count()
```

▶ (2) Spark Jobs  
Out[20]: 387261  
Command took 2.62 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:04:24 on course-work

**Fig 4.5**

The numbers of studies were gotten by using count function, this function, in essence, returns a number of rows to the whole dataset as it will correspond with the total number of distinct IDs for each clinical trial data set. As a result, the following table shows the number of trials for clinicaltrial\_2021 dataset.

## SOLUTION 2

QUESTION 2

```
Cmd 17
1 clinicaltrial_2021.groupBy(clinicaltrial_2021.Type).count()\
2 .orderBy('count', ascending=False)\n3 .show(5, truncate=False)
```

▶ (2) Spark Jobs

Type	count
Interventional	301472
Observational	77540
Observational [Patient Registry]	8180
Expanded Access	69

Command took 4.97 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:05:52 on course-work

**Fig 4.6**

All rows in the ‘**Type**’ column of the clinicaltrial data was collected using the ‘**groupBy**’ function and result aggregated by ‘count’. The use of the truncate false clause with the frequency presented in descending order by using the ‘**orderBy**’ function on the count was used to avoid truncation in the results that usually occur in pyspark, and the result is shown below for the clinical trial\_2021 dataset.

## SOLUTION 3

### QUESTION 3

The top 5 conditions (from Conditions) with their frequencies

```
Cmd 19
1 from pyspark.sql.window import Window
Command took 0.06 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:08:51 on course-work

Cmd 20
1 clinicaltrial_2021.select('Conditions').na.drop()\
2 .select(split('Conditions', ',').alias('Conditions'))\
3 .withColumn('Individual_condition', explode('Conditions'))\
4 .drop('Conditions')\
5 .groupBy('Individual_condition').count()\
6 .orderBy('count', ascending=False)\
7 .limit(5)\
8 .withColumn('S/N', row_number().over(Window.orderBy(monotonically_increasing_id())))\
9 .select(col('S/N'), col('Individual_condition') , col('count'))\
10 .show(truncate=False)

▶ (2) Spark Jobs
+---+-----+-----+
|S/N|Individual_condition|count|
+---+-----+-----+
|1 |Carcinoma           |13389|
|2 |Diabetes Mellitus   |11080|
|3 |Neoplasms            |9371 |
|4 |Breast Neoplasms    |8640 |
|5 |Syndrome              |8032 |
+---+-----+-----+
Command took 7.02 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:08:54 on course-work
```

**Fig 4.7**

**Fig 4.7** started by selecting "Conditions" from the DataFrame and deleting any rows that have an empty value in them. Then it applies "split" to the condition columns, splitting each string with commas and creating a new column 'Conditions' which contains an array of individual conditions.

The next step is to create a new row for each condition on the "Conditions" array by using an explode function. At this point, the current "Conditions"

column is deleted. The resulting DataFrame shall then be divided by "Individual\_Condition", i.e., the individual conditions, and counts of each condition shall be calculated. The resulting DataFrame is sorted by the total count in a descending order, and top 5 criteria are selected using "limit" option.

Finally, using "row\_number" and window function, we can add a new column to the DataFrame named 'SDN'. For each entry in the DataFrame, this column only contains a row number. For the top five conditions, the final DataFrame is selected to display the "SN", "Individual\_Condition", and "Count" columns.

## SOLUTION 4

QUESTION 4

10 most common sponsors that are not pharmaceutical companies + the number of clinical trials sponsored.

```
Cmd 22
1 parent = list(pharma.select('Parent_Company').toPandas()['Parent_Company'])
2
3 clinicaltrial_2021.select('Sponsor')\
4 .filter(~col('Sponsor').isin(parent))\
5 .groupBy('Sponsor').count()\
6 .orderBy('count', ascending=False)\
7 .show(10, truncate=False)

▶ (3) Spark Jobs
+-----+-----+
|Sponsor|count|
+-----+-----+
|National Cancer Institute (NCI)|3218 |
|M.D. Anderson Cancer Center|2414 |
|Assistance Publique - Hôpitaux de Paris|2369 |
|Mayo Clinic|2300 |
|Merck Sharp & Dohme Corp.|2243 |
|Assiut University|2154 |
|Novartis Pharmaceuticals|2088 |
|Massachusetts General Hospital|1971 |
|Cairo University|1928 |
|Hoffmann-La Roche|1828 |
+-----+
only showing top 10 rows

Command took 14.26 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:12:12 on course-work
```

**Fig 4.8**

In Fig 4.8 A list is formed of the parent company column from the pharma's DataFrame which contains values, then the sponsor column of the clinicaltrial\_2021 was filtered to only include values that are not in the parent list. The resulting DataFrame is ordered by counting and displays the top ten results without truncation.

## QUESTION 5

Cmd 24

```

1  #####Table and Visualization for Completed Studies
2
3  DataFrame = clinicaltrial_2021.select('Status', 'Completion').na.drop()\
4  .filter(col('Status') == 'Completed')\
5  .withColumn('Month', split(col('Completion'), ' ').getItem(0))\
6  .withColumn('Year', split(col('Completion'), ' ').getItem(1))\
7  .filter(col('Year') == '2021')\
8  .drop('Status', 'Completion', 'Year')\
9  .groupBy('Month')\
10 .count()
11
12 display(Completed_Studies)

```

▶ (2) Spark Jobs

▶ DataFrame: pyspark.sql.dataframe.DataFrame = [Month: string, count: long]

Table +

	Month	count
1	Oct	187
2	Sep	528
3	Aug	700
4	May	984
5	Jun	1094
6	Feb	934

↓ 10 rows | 3.34 seconds runtime



Command took 3.34 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:47:26 on course-work

Cmd 25

```

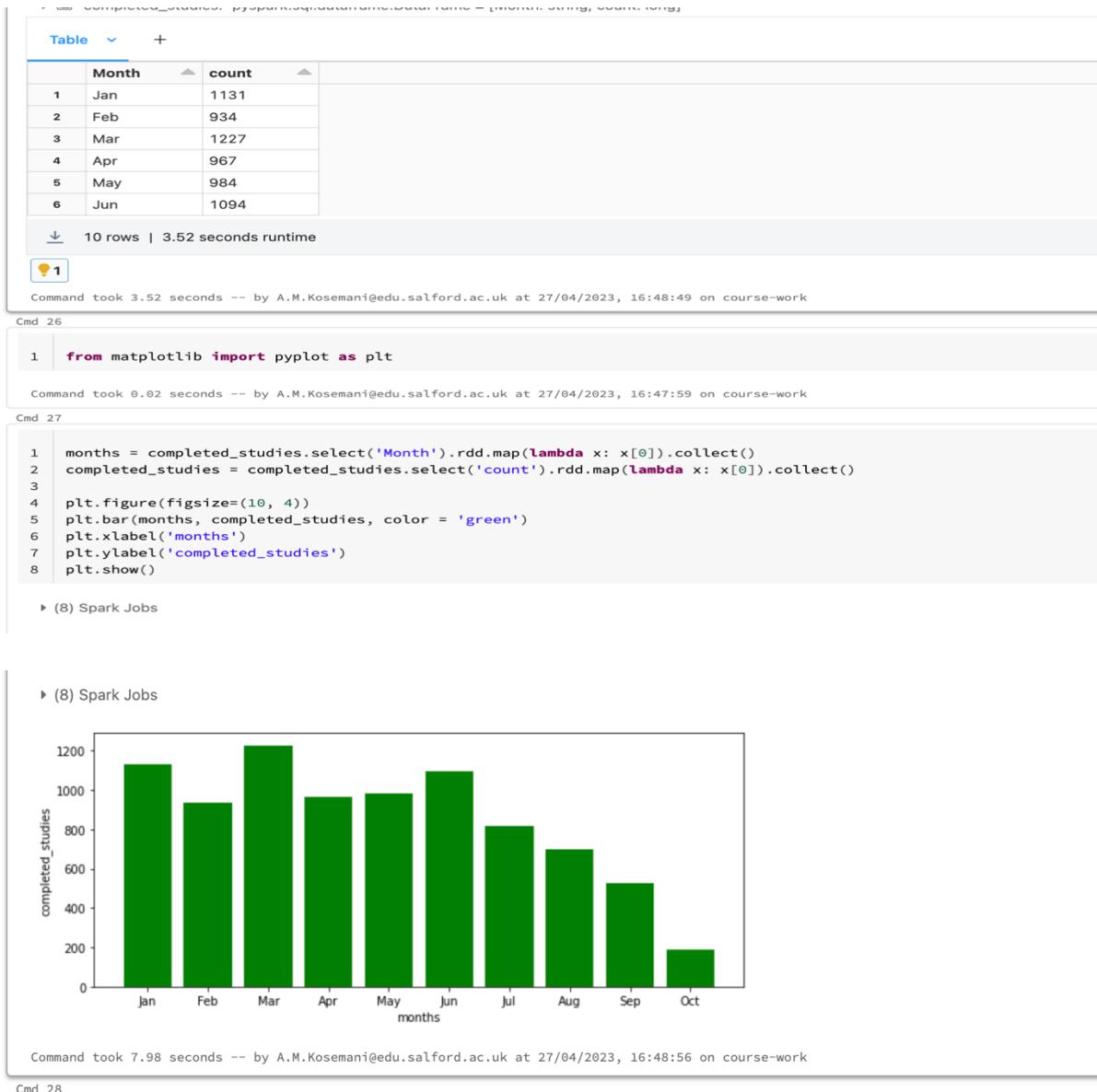
1 completed_studies = DataFrame.withColumn('S/N_Month',
2   when(col('Month')=='Jan', 1) \
3   .when(col('Month')=='Feb', 2) \
4   .when(col('Month')=='Mar', 3) \
5   .when(col('Month')=='Apr', 4) \
6   .when(col('Month')=='May', 5) \
7   .when(col('Month')=='Jun', 6) \
8   .when(col('Month')=='Jul', 7) \
9   .when(col('Month')=='Aug', 8) \
10  .when(col('Month')=='Sep', 9) \
11  .when(col('Month')=='Oct', 10) \
12  .when(col('Month')=='Nov', 11) \
13  .when(col('Month')=='Dec', 12) \
14 )
15
16
17 completed_studies = completed_studies.orderBy('S/N_Month', ascending = True) \
18 .drop('S/N_Month')
19
20 display(completed_studies)
21

```

▶ (2) Spark Jobs

▶ completed\_studies: pyspark.sql.dataframe.DataFrame = [Month: string, count: long]

**Fig 4.9**



**Fig 5.0**

In **Fig 4.9**, the completion and status columns were selected and then drops all the rows with null values, after which filtering was done for the completed studies, followed by splitting the completion column into both **Years** and **Months** columns, then after this I was able to filter through for the studies that were completed in 2021. then the data was grouped by month, which allows the counting of the number of studies which was completed for each of the month. With the result of the DataFrame operation stored in variable name **completed\_studies** .

In **Fig 4.9**, in the completed\_studies DataFrame, a new column **S/N\_Month** was created, which assigns a numerical value to each month. The result of the values is stored in **the\_completed\_studies**. A new DataFrame was created by ordering the completed\_studies DataFrame by the **S/N\_Month** column in ascending order and dropping the **S/N\_Month** column.

In **Fig 5.0**, using Matplotlib a bar graph of completed studies by month is shown.

**EXTRA MARK**

```
Cmd 29
1  #####Completed studies for 2021
2  ### For 2019 and 2020 change the Year for the codes to generate result
3
4  Year_2021 = clinicaltrial_2021.select('Status','Completion').na.drop()\
5  .filter(col('Status') == 'Completed')\
6  .withColumn('Month', split(col('Completion'), ' ').getItem(0))\
7  .withColumn('Year', split(col('Completion'), ' ').getItem(1))\
8  .filter(col('Year') == '2021')\
9  .drop('Status','Completion','Month')\
10 .groupBy('Year')\
11 .count()
12
13 display(Year_2021)

▶ (2) Spark Jobs
▶  Year_2021: pyspark.sql.dataframe.DataFrame = [Year: string, count: long]
Table  +
+---+
| Year | count |
+---+
| 2021 | 8571 |
+---+
↓ 1 row  | 4.01 seconds runtime
💡 1
Command took 4.01 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:58:41 on course-work
Cmd 29
```

(2) Spark Jobs

Year\_2021: pyspark.sql.dataframe.DataFrame = [Year: string, count: long]

Year	count
2021	8571

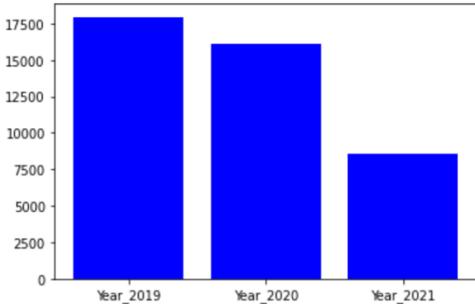
1 row | 4.01 seconds runtime

💡 1

Command took 4.01 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:58:41 on course-work

Cmd 30

```
1  ###Visualization of the comparison of the number of Completed Studies
2  ###between 2019, 2020, 2021.
3  Completed_Studies = ['Year_2019', 'Year_2020', 'Year_2021']
4  count = [17941, 16118, 8571]
5
6  plt.bar(x = Completed_Studies, height = count, color = 'blue')
7  plt.show()
```



Command took 0.30 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 16:59:17 on course-work

Shift+Enter to run  
Shift+Ctrl+Enter to run selected text

The steps of task 5 shall be repeated, except for filtering and grouping the month column by year. An analysis for 2019, 2020 and 2021 has been provided in the images above. As shown below, the results for the three years show a decrease in the number of completed studies each year with a visual representation.

# **TASK TWO**

## INTRODUCTION

The purpose of this task is to undertake a forecast analysis of FaultData that uses the readings obtained by using them, in order to determine whether machines are operating correctly. For every row containing 20 vibration sensor readings, and for each of the last columns indicating if there has been a failure when these measurements have been made, this dataset is entered in CSV format. 0 in this column indicates that there was no mechanical fault was found, while 1 indicates that a fault was found.

Data was analysed using the Pyspark data frame application, which has been performed in an environment called Databricks. As outlined in the presentation, some preliminary analysis of these figures with visualisation was performed on bokeh and matplotlib.

Pre-processing has been carried out by splitting the training and testing data set prior to the modelling of the data to make it suitable for machine learning. The training of the dataset for classification model has been performed with Spark, MLLib machine learning library and this is followed up by an experiment measuring its performance using MLflow.

## COMPUTE SETUP AND FILE IMPORT

The screenshot shows the 'Compute' section of the Databricks UI, specifically the 'coursework-cluster' configuration. Key details include:

- Databricks Runtime Version:** 13.0 ML (includes Apache Spark 3.4.0, Scala 2.12)
- Driver type:** Community Optimized (15.3 GB Memory, 2 Cores, 1 DBU)
- A note about termination: "Free 15 GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours." A link to upgrade the subscription is provided.
- Instances:** The tab is selected, showing the current configuration.
- Spark:** Configuration options for the Spark environment.
- JDBC/ODBC:** Configuration options for connecting to external databases.
- Permissions:** Configuration options for cluster access.
- Availability zone:** us-west-2a

**Fig 1.0 Compute setup.**

A 15GB cluster that included Apache Spark 3.4.0 and Scala 2.12 has been created which is used for the implementation of this task using a Databricks Community Cloud Platform to create clusters containing two attached cores. To implement the data frame, resources and settings that have been created from cluster creation were used to generate notebooks in a standard Python language.

## UPLOADING DATASET

The screenshot shows the 'Create New Table' interface in the Databricks UI. The process is for creating a new table from an uploaded CSV file:

- Data source:** The 'Upload File' tab is selected.
- DBFS Target Directory:** The target directory is set to '/FileStore/tables/' (optional).
- Files:** A single file named 'FaultDataset.csv' is listed, showing it is 1.7 MB in size and can be removed.
- Success message:** A green checkmark indicates the file was uploaded successfully to '/FileStore/tables/FaultDataset.csv'.
- Action buttons:** Buttons for 'Create Table with UI' and 'Create Table in Notebook' are available at the bottom.

**Fig 1.1 Uploading the dataset using the graphical user's interphase(GUI)**

## IMPORTING LIBRARIES

# Setup

## Load libraries

Cmd 2

```
1 # exploration libraries
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 import pyspark.sql.functions as F
5
6 # ML libraries
7 from pyspark.ml import Pipeline
8 import mlflow
9 import pyspark.ml.feature as MF
10 from pyspark.ml.classification import DecisionTreeClassifier, RandomForestClassifier
```

Command took 2.57 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 11:25:15 on coursework-cluster

Cmd 3

## Fig 1.2 Data exploration and machine learning libraries

```
| Loading the data
| Cmd 4
| FaultDataset=spark.read.csv(
|   '/FileStore/tables/FaultDataset.csv',
|   sep=',',
|   header=True,
|   inferSchema=True
| )
| FaultDataset.show(5)
|
| (3) Spark Jobs
|
| FaultDataset: pyspark.sql.dataframe.DataFrame = [1: double, 2: double ... 19 more fields]
+-----+
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
+-----+
| 0.3503125 | 0.3496875 | 0.35 | 0.3459375 | 0.3475 | 0.3459375 | 0.341875 | 0.3434375 | 0.355 | 0.3553125 | 0.3459375 | 0.3525 | 0.3575 | 0.3596625 | 0.35875 | 0.34843
75 | 0.3599625 | 0.35 | 0.3559375 | 0.3490625 | 0 |
| 0.5625 | 0.484375 | 0.464875 | 0.464875 | 0.073125 | 0.061 | 0.0634375 | 0.0575 | 0.0546875 | 0.0559375 | 0.058125 | 0.0628125 | 0.065625 | 0.0640625 | 0.0634375 | 0.0534375 | 0.0843
75 | 0.6640625 | 0.65375 | 0.646875 | 0.646875 | 0.073125 | 0.066875 |
| 0.028125 | 0.0975 | 0.109375 | 0.109375 | 0.1025 | 0.09625 | 0.1053125 | 0.09875 | 0.098125 | 0.091875 | 0.0909375 | 0.09875 | 0.103125 | 0.1 | 0.1034375 | 0.1015625 | 0.09768
25 | 0.0996625 | 0.10375 | 0.098125 | 0.1040625 | 0.1053125 | 0.096875 | 0.1053125 | 0.09875 | 0.098125 | 0.091875 | 0.0909375 | 0.09875 | 0.103125 | 0.1 | 0.1034375 | 0.1015625 | 0.09768
75 | 0.09375 | 0.089375 | 0.091875 | 0.0996875 | 0.09909375 | 0.096875 | 0.0940625 | 0.096875 | 0.096875 | 0.099375 | 0.099375 | 0.0959375 | 0.0959375 | 0.0940625 | 0.09125 | 0.09968
75 | 0.09375 | 0.0934375 | 0.091875 | 0.094375 | 0 |
| 0.036875 | 0.0440625 | 0.038125 | 0.0428125 | 0.0353125 | 0.0346625 | 0.033125 | 0.0403125 | 0.0346875 | 0.036875 | 0.035625 | 0.03625 | 0.0409375 | 0.039375 | 0.035 | 0.0406
25 | 0.0384375 | 0.036875 | 0.04 | 0.0371875 | 0 |
only showing top 5 rows

Command took 10.56 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 19:02:23 on course-work
```

### **Fig 1.3 reading the file directory**

Calling `Show()` function in order to show all five rows of DataFrame once I have read the file. This is a good way to look at the data and make sure that they are properly interpreted.

## **EXPLANATORY DATA ANALYSIS**

Cmd 6

# Exploratory Data Analysis

Cmd 7

```
1 FaultDataset.printSchema()

root
 |-- 1: double (nullable = true)
 |-- 2: double (nullable = true)
 |-- 3: double (nullable = true)
 |-- 4: double (nullable = true)
 |-- 5: double (nullable = true)
 |-- 6: double (nullable = true)
 |-- 7: double (nullable = true)
 |-- 8: double (nullable = true)
 |-- 9: double (nullable = true)
 |-- 10: double (nullable = true)
 |-- 11: double (nullable = true)
 |-- 12: double (nullable = true)
 |-- 13: double (nullable = true)
 |-- 14: double (nullable = true)
 |-- 15: double (nullable = true)
 |-- 16: double (nullable = true)
 |-- 17: double (nullable = true)
 |-- 18: double (nullable = true)
 |-- 19: double (nullable = true)
 |-- 20: double (nullable = true)
```

Command took 0.11 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 22:14:17 on course-work

## **Fig 1.4 Exploratory data analysis**

Prior to the creation of the model in machine learning, the EDA was usually used to examine the data for patterns, identify trends where they exist, check for possible assumptions using statistical descriptions and graphical representations.

The first EDA to be carried out for this task has been the checking of the type of data in respect of each variable included in the dataset. This was done using the **printSchema** syntax as seen in **Fig 1.4** below with all the 20 explanatory variables being the same data type while the response variable; **fault\_detected** column is an integer consisting of “0” and “1” with 0 meaning no fault was detected with the machine and 1 means a fault was identified.

**Fig 1.5 Fault dataset count and dimension.**

The shape of the Fault Dataset was determined in a further examination carried out under EDA. The length of rows and columns in a two-dimensional data frame structure is computed by means of the number or Len for each row and column, as well as determining the shape dimension of the data. The dimensions of the data are very common to give you an idea of how big it is. The screen shot below shows the size of a FaultDataset with an array of 21 columns.

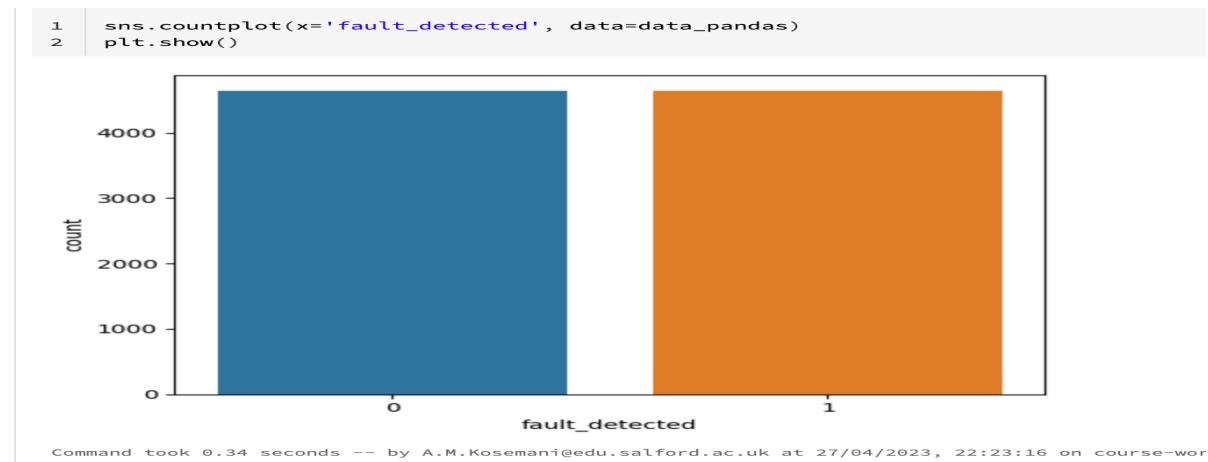
```
Cmd 10
1 ## Checking the distribution of Fault Detected (outcome/response) Column
2 (
3   FaultDataset.groupBy('fault_detected')
4   .count()
5   .show()
6 )

▶ (2) Spark Jobs
+-----+-----+
|fault_detected|count|
+-----+-----+
|          1| 4646|
|          0| 4646|
+-----+-----+

Command took 2.16 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 27/04/2023, 22:16:23 on course-work
Cmd 11
```

**Fig 1.6 Distribution of response variable**

In order to determine a pattern for distribution of response variables, a univariate analysis was carried out with the **fault\_detected** column as this factor is very important in determining suitable modelling techniques that can give accurate predictions. The response variable for this task has two classes in which observations are categorized, which are “0” and “1”. The data frame was grouped by all columns using the groupBy method, and a pivot table is made with the fault\_detected column. The resulting data file **FaultData** shows each unique value separately, including a count of the number of times it has occurred in an **error\_detected** column.



**Fig 1.7 To visualize the above result, a library called Seaborn was used.**

Cmd 12

## Descriptive Statistics

Cmd 13

```
1 display(FaultDataset.drop('fault_detected').describe())
```

► (2) Spark Jobs

Table +

	summary	1	2	3	4	5	6
1	count	9292	9292	9292	9292	9292	9292
2	mean	0.34162330499354226	0.34263116121394677	0.3421213812957383	0.34213907124407966	0.3428434405940584	0.342827936665946
3	stddev	0.28919489486260785	0.2890875372793958	0.28916422490616933	0.28916356333107296	0.2889646554403878	0.289088989972954
4	min	0.024375	0.024375	0.024375	0.024375	0.024375	0.024375
5	max	1.0809375	1.2134375	1.0809375	1.0809375	1.0809375	1.0809375

↓ 5 rows | 5.20 seconds runtime

Command took 5.20 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster

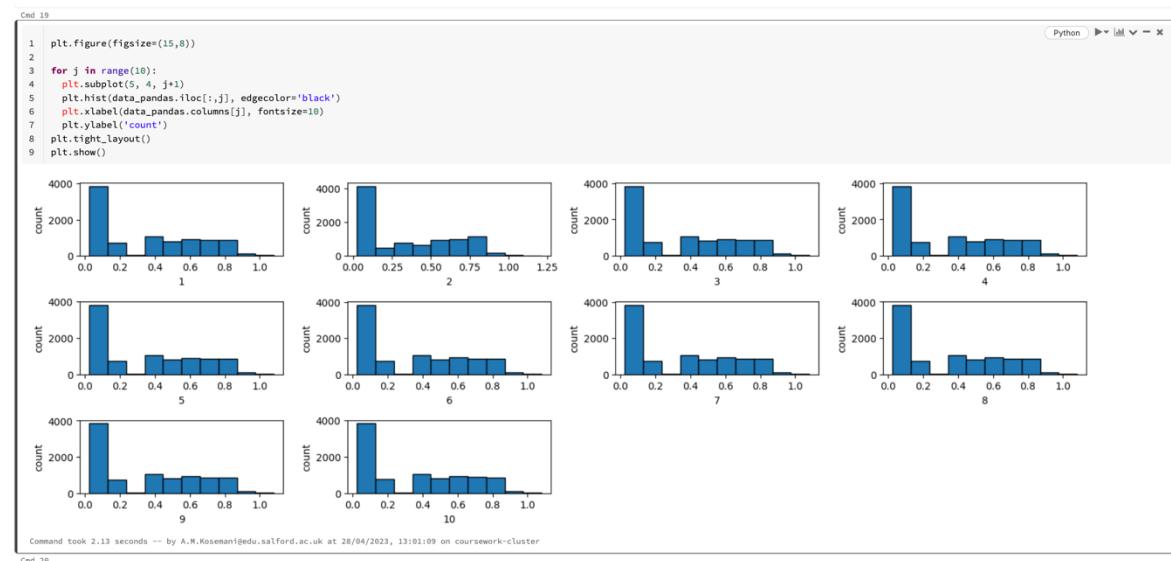
Refreshed 11 minutes ago

Cmd 14

## Fig 1.8 Descriptive statistics

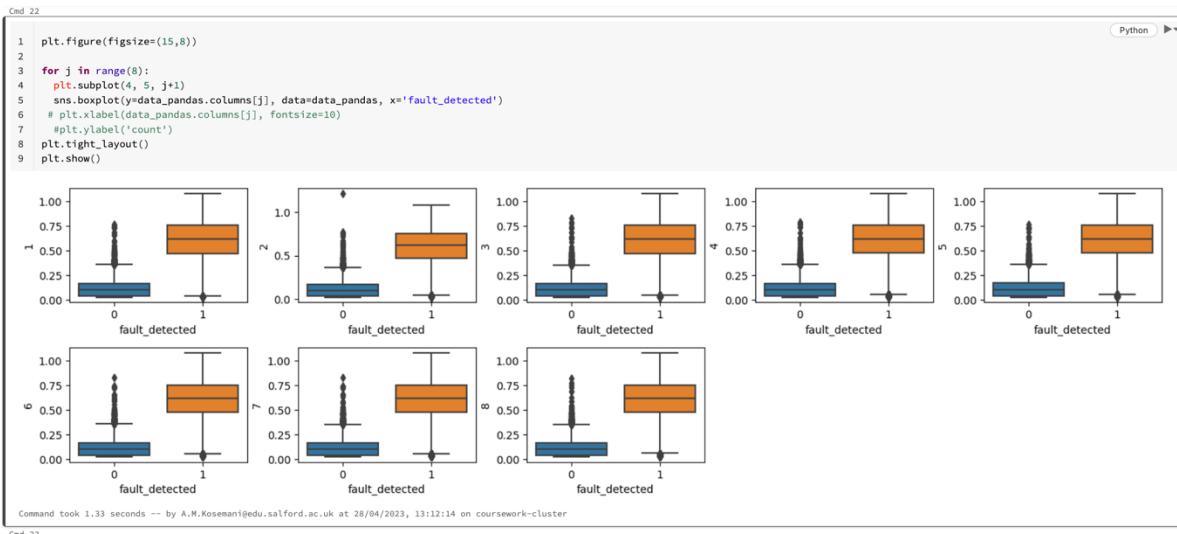
Another important analysis performed on EDA in the task is the descriptive statistics to shed light on the distribution, central tendency, and variability of the Fault Dataset. The mean and the median which reveals the data's basic or essential value were generated using a summary function.

Details regarding the distribution or dispersion of the data from a standard deviation have been revealed, which provide insight into an instrument for measuring variability. The distribution of 25th to 50th, and 75th percentiles is also shown in the data summary.



## Fig 1.9

In **Fig 1.9** A visualization of this summary is presented with histogram and box plot in the screenshot below. histogram gives a frequency distribution for the readings from vibration sensors, showing that there are 10 explanatory variables



**Fig 2.0**

A five number summary of FaultDataset, which contains minimum, lower, median, upper quartile and maximum scores, is shown in the box plot. From the diagram above, the above summary was done based on the class of the response variable of either “0” or “1”.

The middle area of the box plot is a median and seems to be positioned in the centre for each set of boxes suggesting normal distribution. Some readings below the lower quartile, which indicate that certain outliers may exist in the data set, can also be seen.



**Fig 2.1**

This measurement provides details on the relationship of two variables in terms of direction and strength. Correlation coefficients for both positive and negative correlation are found, with 0 of them indicating no correlation.

A correlation matrix was created, in which all explanatory variables are included, and a graph is displayed as a colour coded heatmap showing the very close relation of these data points. Different explanatory variables in the data set also have dark and lighter colours that indicate higher correlation coefficients. The following images show the red to be positive and the blue to be negative correlated.

```
Cmd 23
[ ] Train-test Split

Cmd 24
1 from pyspark.ml import Pipeline
2 import mlflow
3 import pyspark.ml.feature as MF
4 from pyspark.ml.classification import DecisionTreeClassifier

Command took 0.09 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster

Cmd 25
1 mlflow.pyspark.ml.autolog()

Command took 0.20 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster

Cmd 26
1 train, test=FaultDataset.randomSplit(weights=[0.5, 0.3], seed=35)
2
3 print("Training count Dataset: " + str(train.count()))
4 print("Test count Dataset: " + str(test.count()))

▶ (4) Spark Jobs
▶ [ ] train: pyspark.sql.dataframe.DataFrame = [1: double, 2: double ... 19 more fields]
▶ [ ] test: pyspark.sql.dataframe.DataFrame = [1: double, 2: double ... 19 more fields]

Training count Dataset: 5835
Test count Dataset: 3457

Command took 1.59 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 13:54:47 on coursework-cluster

Cmd 27
```

**Fig 2.2 Train-test split**

The training data and the testing data for the task has a corresponding number of rows to be 5835 and 3457.

Splitting the dataset into two parts is a first step to pre-processing data after it has been cleaned, namely training data used for machine learning model and testing data that enable models to be tested post training.

As you can see below, a training experiment using the Decision Tree Model was monitored through an ML flow.

Run Name	Created	Duration	Source	Models
honorable-wolf-245	1 hour ago	2.1s	KOSEM...	-
silent-turtle-660	2 hours ago	1.7min	KOSEM...	spark
debonair-kite-46	2 hours ago	3.5min	KOSEM...	spark, 1 more
adventurous-toad-693	2 hours ago	1.2min	KOSEM...	spark
agreeable-stoat-440	2 hours ago	2.0s	KOSEM...	-
shivering-panda-616	2 hours ago	3.3s	KOSEM...	-
judicious-loon-406	16 hours ago	1.6min	KOSEM...	spark
persistent-moose-628	16 hours ago	3.3min	KOSEM...	spark, 1 more
fortunate-grouse-691	16 hours ago	1.1min	KOSEM...	spark
rumbling-rook-499	16 hours ago	1.3s	KOSEM...	-
bright-roo-929	16 hours ago	3.5s	KOSEM...	-

Fig 2.3

### Evaluate Model on Testing Set

```
Cmd 34
1 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
Command took 0.09 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster

Cmd 35
1 evaluator=MulticlassClassificationEvaluator(labelCol='fault_detected', predictionCol='prediction', metricName='accuracy')
2 acc_DT=evaluator.evaluate(Pred_DT)
3 acc_DT
4 Out[47]: 0.970391061452514
5
6

▶ (1) Spark Jobs
Command took 1.40 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster

Cmd 36
1 f1_dt=evaluator.evaluate(Pred_DT, {evaluator.metricName: "f1"})
2 f1_dt

▶ (1) Spark Jobs
Out[32]: 0.9703912000678954
Command took 1.19 seconds -- by A.M.Kosemani@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster
```

Fig 2.3 Evaluating Model on Testing set

```

1 GS = tvs.fit(scaled_train)

▶ (56) Spark Jobs
▼ (19) MLflow runs
  Logged 19 runs to an experiment in MLflow. Learn more

2023/04/28 11:42:24 INFO mlflow.utils.autologging_utils: Created MLflow autologging run with ID '3fea5de3460347c2889f4c3c3f874bac', which will track hyperparameters, performance metrics, model artifacts, and lineage information for the current pyspark.ml workflow
2023/04/28 11:43:54 INFO mlflow.spark: Inferring pip requirements by reloading the logged model from the databricks artifact repository, which can be time-consuming. To speed up, explicitly specify the co
nda_env or pip_requirements when calling log_model().
2023/04/28 11:44:57 INFO mlflow.spark: Inferring pip requirements by reloading the logged model from the databricks artifact repository, which can be time-consuming. To speed up, explicitly specify the co
nda_env or pip_requirements when calling log_model().

Command took 3.46 minutes -- by A.M.Koseman@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster
Cmd 41

1 bestModel=GS.bestModel
2
3 print(f'Best Impurity Measure: {bestModel.getImpurity()}')
4 print(f'Best Max Bin Measure: {bestModel.getMaxBins()}')
5 print(f'Best Max Depth Measure: {bestModel.getMaxDepth()}')
6

Best Impurity Measure: entropy
Best Max Bin Measure: 32
Best Max Depth Measure: 20

Command took 0.08 seconds -- by A.M.Koseman@edu.salford.ac.uk at 28/04/2023, 12:40:08 on coursework-cluster

```

**Fig 2.4 Hyperparameter optimization**