

Assignment #2 (Team work)

Personal Ethics & Academic Integrity Statement

Student name: Rahul Raj Sharma

Student ID: 300111840

Student Name: Harman Bhutani

Student ID: 300144160

By typing in my name and student ID on this form and submitting it electronically, I am attesting to the fact that I have reviewed not only my own work, but the work of my team member, in its entirety.

I attest to the fact that my own work in this project adheres to the fraud policies as outlined in the Academic Regulations in the University's Undergraduate Studies Calendar. I further attest that I have knowledge of and have respected the "Beware of Plagiarism" brochure found on the Telfer School of Management's doc-depot site. To the best of my knowledge, I also believe that each of my group colleagues has also met the aforementioned requirements and regulations. I understand that if my group assignment is submitted without a completed copy of this Personal Work Statement from each group member, it will be interpreted by the school that the missing student(s) name is confirmation of non-participation of the aforementioned student(s) in the required work.

We, by typing in our names and student IDs on this form and submitting it electronically,

- warrant that the work submitted herein is our own group members' work and not the work of others
- acknowledge that we have read and understood the University Regulations on Academic Misconduct
- acknowledge that it is a breach of University Regulations to give or receive unauthorized and/or unacknowledged assistance on a graded piece of work

Part A - For this part of the assignment please answer questions to the best explanation possible.

Note: You may consult with resources on the Internet to research and answer these the questions. However, please do not copy and paste your answers directly from a google search! The instructor is looking for your understanding of the topic in your own words.

[5 Pts for each question]

Question 1: What is RDD? Briefly explain Transformations and Actions in the context of RDDs. Give 5 examples each for transformations and actions functions with a short description. (1 line description only for the examples)

Answer: RDD : Resilient Distributed Datasets is a fundamental data structure of Spark which is an immutable distribution collection of objects which computes on the different node of the cluster. Each dataset in RDD is divided into logical partitions, which may be computed on different nodes of the cluster.

Transformation : Transformation is a operation which will transform RDD data from one form to another and when we apply this operation to any RDD, we will get the transformed data in new RDD. Transformation are of two types first one is Narrow and another one is Wide. Below are the operations performed in transformation:

1. **map(func)** - The map function iterates over every line in RDD and split into new RDD
2. **flatMap()** - With the help of **flatMap()** function, to each input element, we have many elements in an output RDD
3. **union(dataset)** - With the **union()** function, we get the elements of both the RDD in new RDD.
4. **distinct()** - It returns a new dataset that contains the **distinct** elements of the source dataset. It is helpful to remove duplicate data.
5. **coalesce()** - To avoid full shuffling of data we use coalesce() function. In **coalesce()** we use existing partition so that less data is shuffled.

Action : Transformations create RDDs from each other, but when we want to work with the actual dataset, at that point action is performed. When the action is triggered after the result, new RDD is not formed. Actions are Spark RDD operations that give non-RDD values. Below are the operations that are performed under action:

1. **collect()** - The action **collect()** is the common and simplest operation that returns our entire RDDs content to driver program.
2. **take(n)** - The action **take(n)** returns n number of elements from RDD.
3. **countByValue()** - The **countByValue()** returns, many times each element occur in RDD.
4. **reduce()** - The **reduce()** function takes the two elements as input from the RDD and then produces the output of the same type as that of the input elements.

5. **Aggregate()** - It gives us the flexibility to get data type different from the input type.
The **aggregate()** takes two functions to get the final result.

Question 2: What is the Spark stack? Briefly explain about any 4 libraries that constitute the Spark ecosystem? (3-4 lines each, no diagrams required)

Answer: Spark is a general purpose cluster computing framework that empowers its central engine to exploit certain higher-level components. It is operationable with Apache Hadoop, in the sense that it can read and write data to and from HDFS and can also integrate the data with other storage systems that are supported by Hadoop API. Below are the libraries that constitute the spark ecosystem :

1. **Spark Core** : In a way the Spark center is close to an operating system's kernel. It is the engine of general operation, which is both quick and sensitive of fault. Upon this main engine is designed the whole Spark ecosystem. It is designed primarily to do job preparation, task delivery and job control through worker nodes. It is also responsible for memory allocation, interaction with various unidentical storage systems and various other operations.
2. **Spark SQL** : This Spark framework is designed for querying, evaluating and running organized system operations. It is a very critical component of the whole Spark stack due to the fact that most operational data are organized, while unstructured data are increasingly that. Acting as a distributed database engine, it makes requests from Hadoop Hive to run up to 100 times faster without any change
3. **Mlib** : MLlib is the built-in machine learning library in the Spark stack. This was introduced in Spark 0.8. The aim is to render machine learning fast and scalable. Developers will use Spark SQL, Spark Streaming and GraphX effortlessly in their favorite programming language, whether it's Java, Python or Scala. MLlib includes the functions required to conduct different statistical research such as associations, regression, hypothesis checking, etc.
4. **SparkR** : The SparkR project started to combine R's mathematical analytics and machine learning capabilities with Spark's scalability. This discussed R's weakness, which was their capacity to store as many data as a single machine's memory suit. R programs will now scale over SparkR in a distributed environment.

Question 3: You have been appointed as a Data Consultant by a company that would like to take advantage of Big data analytics. You have been asked to provide a brief comparative analysis of Hadoop and Spark big data processing tools and present a summary report to the executive committee. (Maximum 1 page) You are free to make assumptions about the requirements of the company for big data analytics. Remember your assumptions will set the grounds of your preferred tool. (Refer to presentation slides on Spark for requirements assumption)

Answer: Bhutani Corporation is the leading Computer, Laptop and Accessory retailer in the country. The company has recently moved from in-store to the e-commerce platform. As a business, a year of product sales has been successfully completed both in the physical and online E-commerce store.

Harman Bhutani, the company's CEO needs an analysis of last year's business revenues, because the company is facing pressure and consumers have become more competitive—because they want business processes to be quicker, quicker, and product quality to be superior and low compared to their rivals.

Company wants to learn more about consumer desires, expectations, the desire to buy for any viable business. With the rise in data sources, Harman intended to take a step away from conventional resources. Business made use of common resources like Access and Excel. As a Data Consultant, suggestion to company is to migrate to more advanced data analytics tools such as Apache Hadoop, Apache Spark and many more.

Hadoop is the oldest and most mature platform, but it would always suggest incorporating Spark in our setting as Apache Spark is the contemporary big data processing engine that provides rapid solutions for any vulnerabilities compared to Hadoop. It can be used effectively in identifying specific trends that are helpful to the company perspective

General difference between Spark and Hadoop

The Basis Of Comparison Between Hadoop vs Spark	Hadoop	Spark
Category	Basic Data processing engine	Data analytics engine

Usage	Batch processing with a huge volume of data	Process real-time data, from real-time events like Twitter, Facebook
Latency	High latency computing	Low latency computing
Data	Process data in batch mode	Can process interactively
Ease of use	Hadoop's MapReduce model is complex, need to handle low-level APIs	Easier to use, abstraction enables a user to process data using high-level operators
Scheduler	External job scheduler is required	In-memory computation, no external scheduler required
Security	Highly secure	Less secure as compare to Hadoop
Cost	Less costly since MapReduce model provide a cheaper strategy	Costlier than Hadoop since it has an in-memory solution

Productive Implication of Spark in ABC environment

Commonly, retail sales data include the latest product details as it is continually updated. In our company, for example, the mega deals will be updated both at the sites and in the retail stores. To uplift the company, cashback deals will be suggested in order to encourage buyers to shop at certain retail locations, as a result of which they will get huge profits. Thus, Apache Spark aims to deeply process data in real-time.

Because Apache Spark offers an application programming interface based on a data structure known as the robust distributed dataset (RDD), sales data objects would be spread around a computer cluster, which will be managed in a fault-tolerant manner.

Via Apache Spark, streaming data can be obtained in this way, which will make it easier for us to sort them into different categories, which will lead to improved business skills. These collected data items can also be used in predictive analytics such as assessing particular customer interest in various demographic and geographic aspects.

Spark is an open source and highly flexible to understand the success of industry and can actually meet the demands of ever-growing data scale.

Question 4: Briefly explain the following in the context of Apache Spark. (2 Sentences maximum)

Answer:

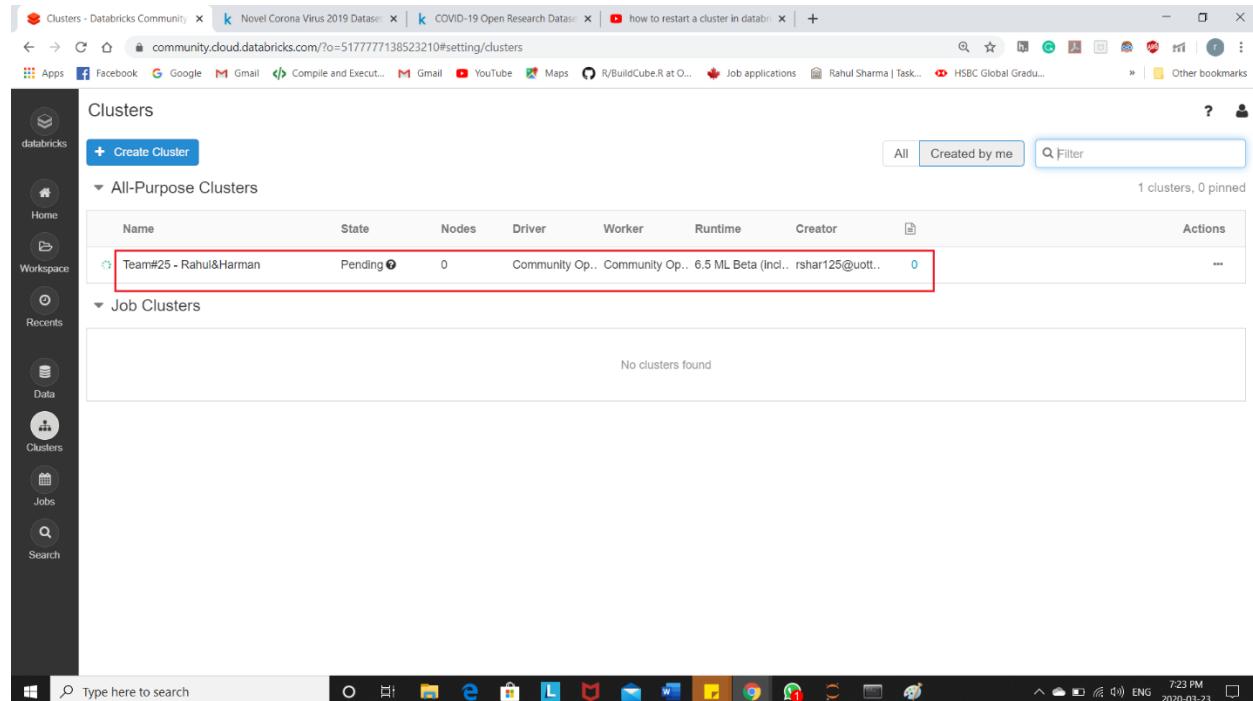
- a. **Apache Spark :** Apache Spark is a data processing system that can execute computational tasks easily on very broad data sets, and can also spread data processing tasks through many devices, either individually or in combination with other distributed computing resources. Such two attributes are essential to the environments of big data and machine learning that involve a huge computational power organising to crunch across broad data stores.
- b. **Spark Context :** SparkContext is the main entry point for spark functionality. This represents the connection to a Spark cluster, and can be used to create RDDs, accumulator and broadcast variables on that cluster. Only one SparkCluster may be active per JVM.
- c. **Master & Slave :** The machine on which the spark standalone cluster managers runs is called the master node. The resources are allocated by master. It uses the "Workers" running throughout the cluster for the creation of Executors for the "Driver". Slave node or the worker node are spawned during the life of a spark application. They execute work on the data within the node and report back to the master node.
- d. **Fault Tolerance :** Fault refers to failure, thus fault tolerance in Apache Spark is the potential to work and restore damage after a failure. If we want our device to be tolerant to error, it should be redundant, because we need a redundant portion to get the missing information. Redundant data is used to restore the defective data.

- e. **DAG** : Directed Acyclic Graph is a finite direct graph with no directed cycles. There are finite number of vertices and edges. Where each edge directed from one vertex to another. DAG operation can do better global optimization than other systems like MapReduce.

Part B – For this part of the assignment please refer to the manual posted on the course website called Spark Experiment on Databricks.doc. Create an account on Databricks® community edition and provide snapshots of your experiment as part of the answers to the questions (5) – (8). You are free to choose any two data files. One must have normal text content (up to 50 lines) and the second data file must be a data file (in excel or CSV format, with up to 1000 records/rows). Please submit a copy of your data files along with this assignment through BrightSpace.

Question 5 [2 pts]: Create a Cluster in the databricks community edition, name the cluster in the format Team# and provide the following **SCREENSHOTS** with the cluster name clearly visible:

- a) Cluster pending state – Below screenshot shows the cluster is in pending state



The screenshot shows the Databricks Community interface. On the left is a sidebar with icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area is titled 'Clusters' and shows a table of clusters. A cluster named 'Team#25 - Rahul&Harman' is listed under 'All-Purpose Clusters'. The 'Name' column shows 'Team#25 - Rahul&Harman', the 'State' column shows 'Pending', and the 'Nodes' column shows '0'. A red box highlights this row. The table has columns for Name, State, Nodes, Driver, Worker, Runtime, Creator, and Actions. The 'Actions' column for this row contains a three-dot menu icon. The top right of the interface shows '1 clusters, 0 pinned'. The bottom right corner of the screenshot shows the Windows taskbar with the date and time as '7:23 PM 2020-03-23'.

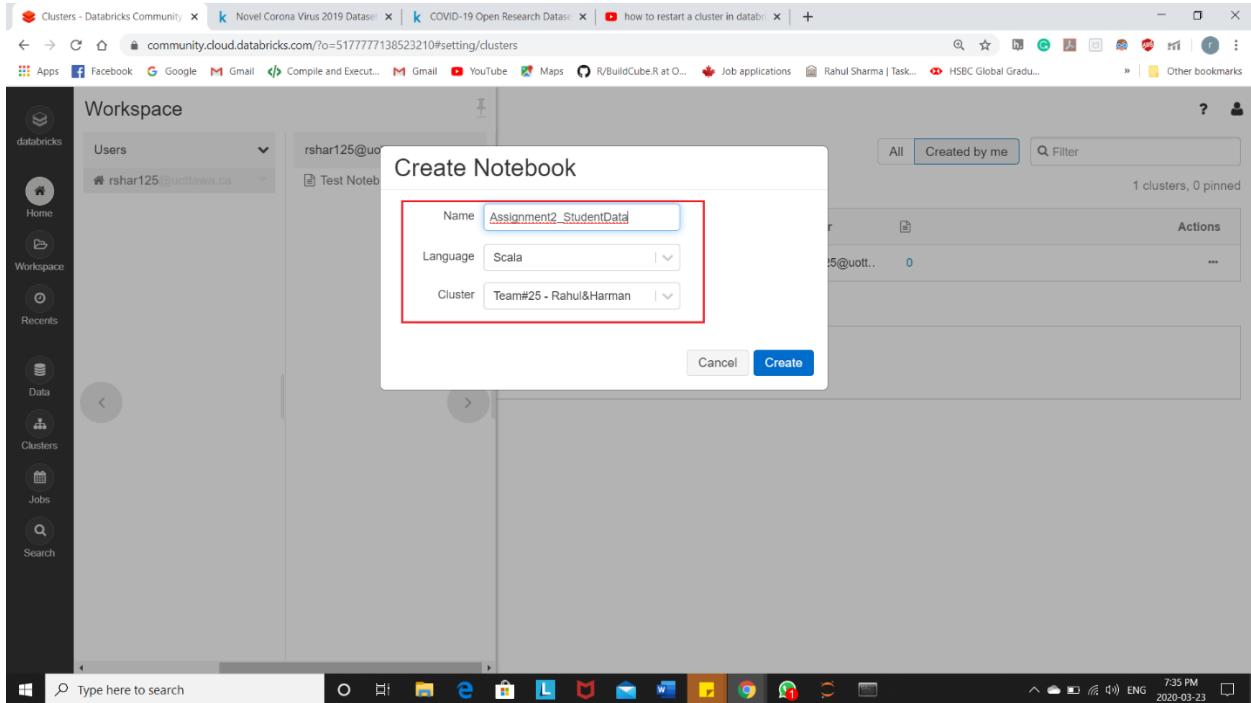
- b) Cluster running state – Below screenshot shows cluster is in running state

The screenshot shows the Databricks Clusters interface. On the left sidebar, there are links for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area is titled 'Clusters' and shows a table of clusters. A button '+ Create Cluster' is at the top left of the table. The table has columns: Name, State, Nodes, Driver, Worker, Runtime, Creator, and Actions. One row is highlighted with a red border: 'Team#25 - Rahul&Harman' (State: Running, Nodes: 1 (0 spot), Driver: Community Op., Worker: Community Op., Runtime: 6.5 ML Beta (incl. rshar125@uott..), Creator: 0). Below the table, a section for 'Job Clusters' is shown with the message 'No clusters found'. At the bottom of the screen, there is a taskbar with various icons and a system tray showing the date and time.

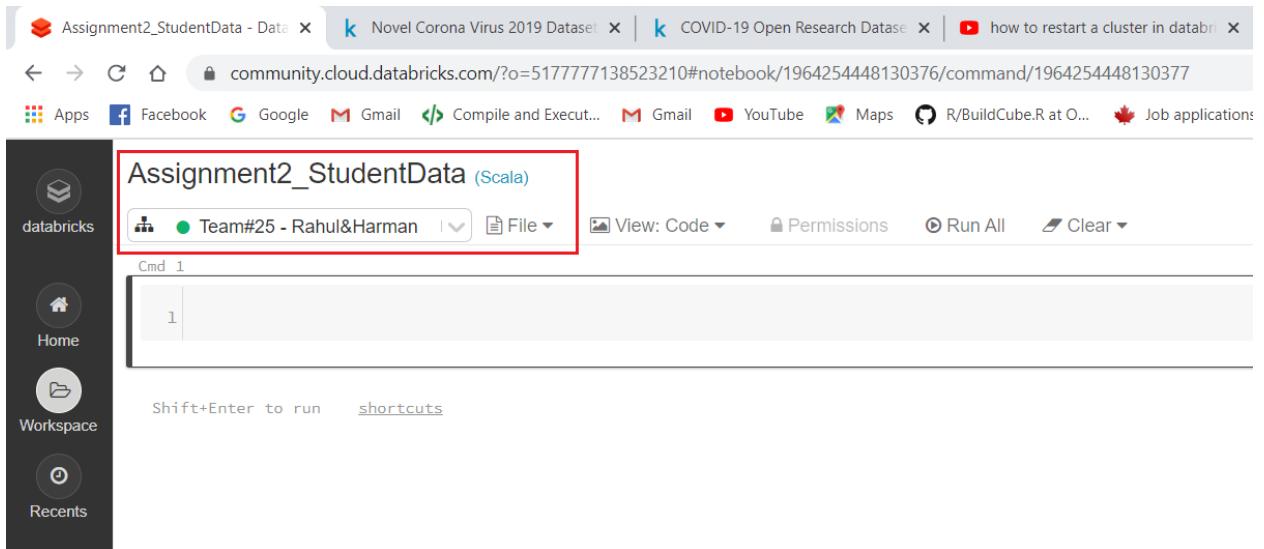
Name	State	Nodes	Driver	Worker	Runtime	Creator	Actions
Team#25 - Rahul&Harman	Running	1 (0 spot)	Community Op.	Community Op.	6.5 ML Beta (incl. rshar125@uott..)	0	...

Question 6 [6 pts]: Create a notebook in Scala using your active cluster, upload the text file (.txt format) selected for your team, and copy the data source file link. Provide **SCREENSHOTS** of the following:

- Your new Scala notebook – Below screenshot shows the creation on new notebook elements which we have set like name of the notebook, language which we use and the cluster name associated with the notebook.

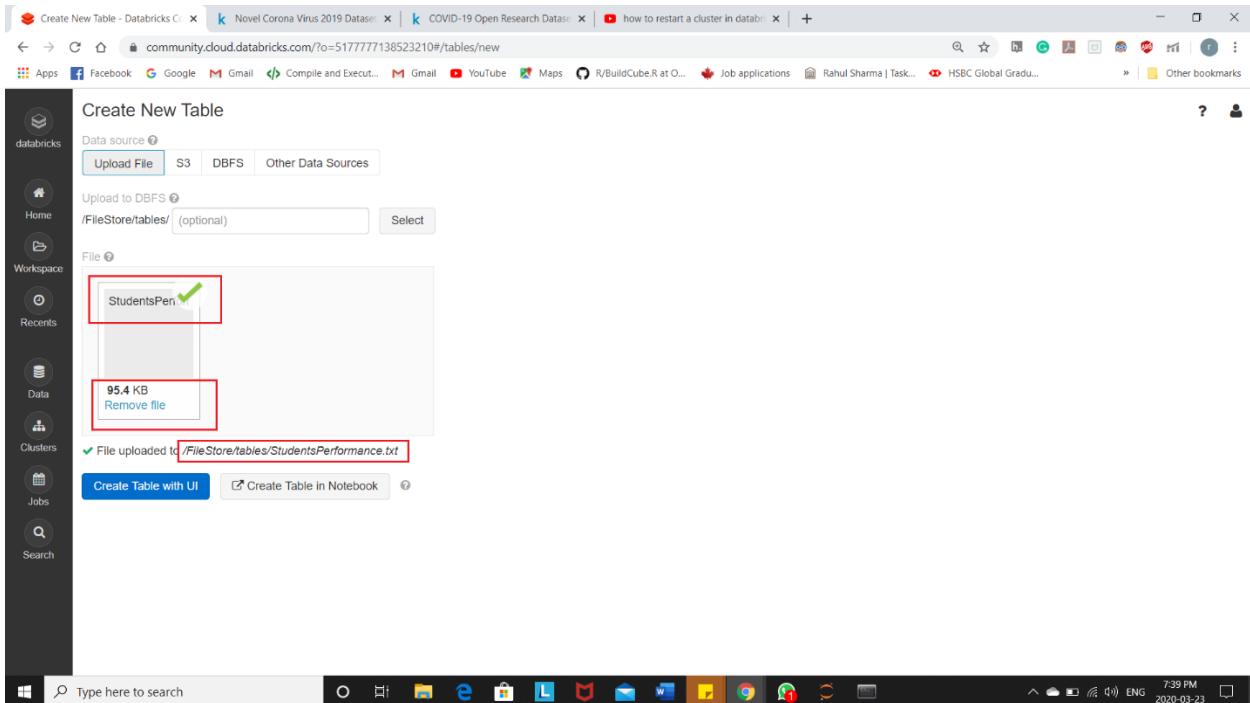


Below screenshot shows the notebook is ready for the coding



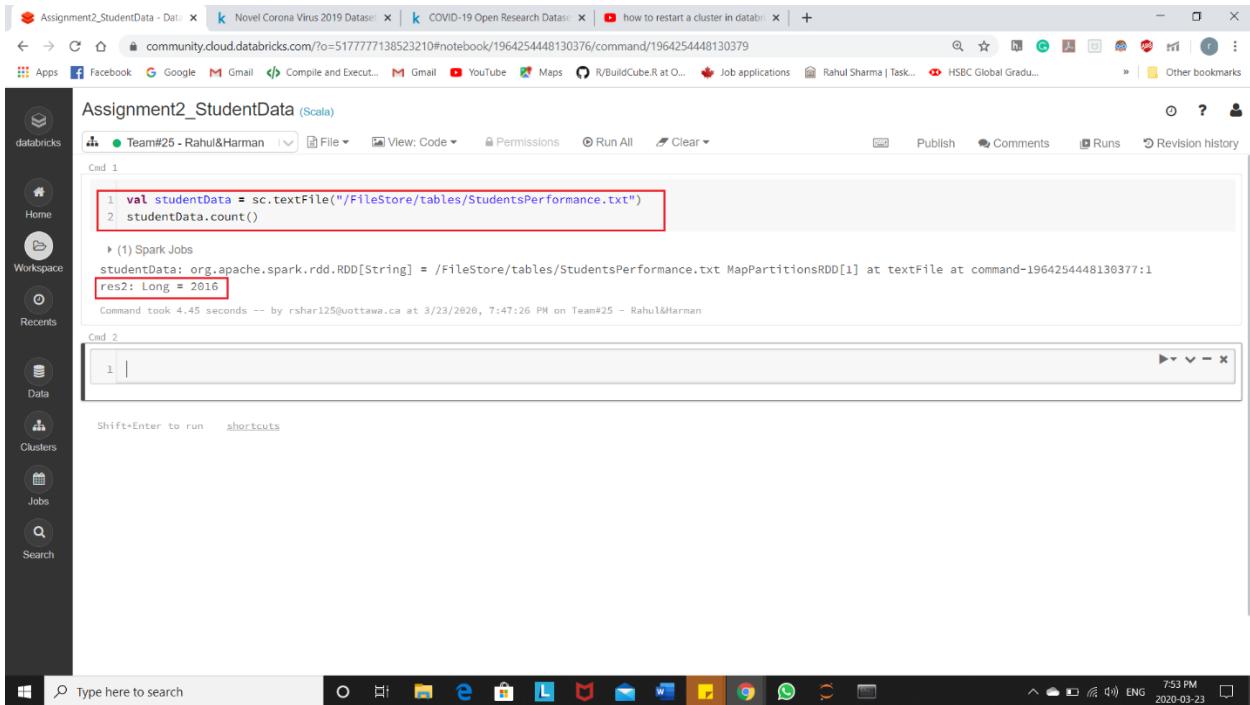
- b) Uploaded data file in databricks showing the file name, size, and the databricks file link –

Below screenshot shows the file is successfully uploaded along with the name of the file, size of the file and the link of the file where the file is uploaded

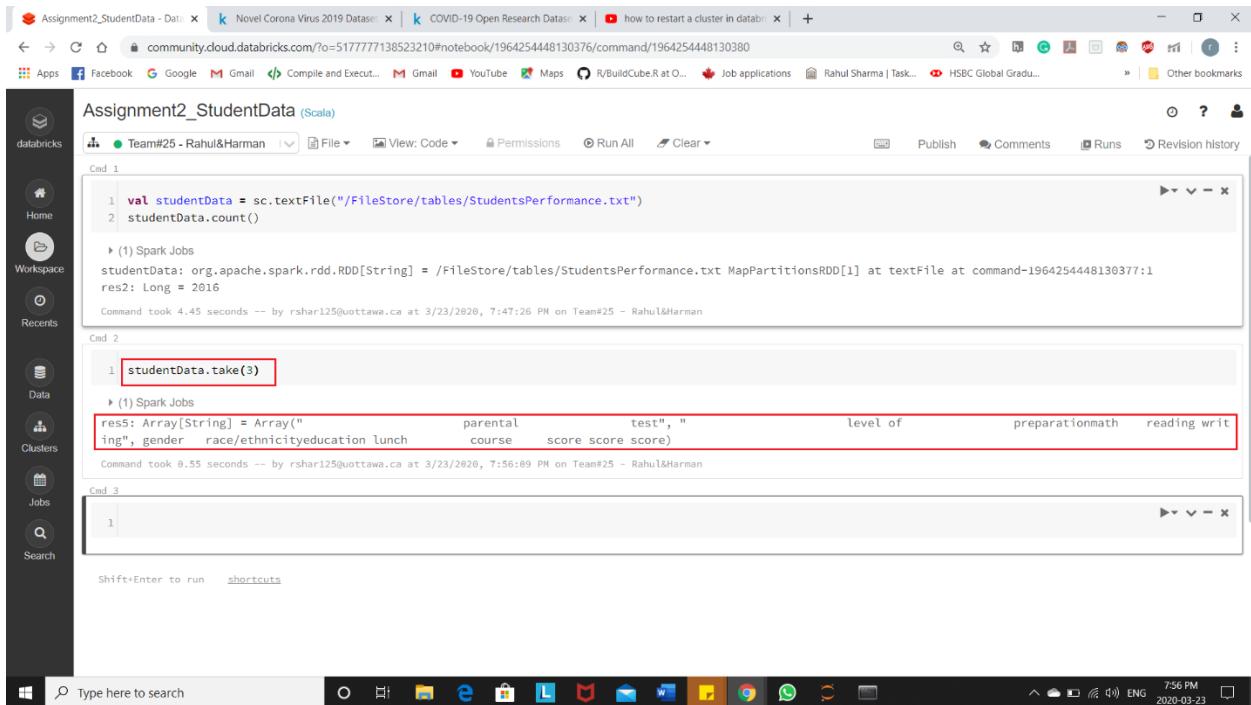


For questions (c) – (f), please provide a **SCREENSHOTS** of the code with the result from your Scala notebook:

- c) Count the number of lines of text in your data file – Below screenshot shows the file is successfully uploaded and the no. of lines is equal to 2016



- d) Retrieve the first 3 rows of the data file – Below screenshot shows the first 3 line of the file



The screenshot shows a Databricks notebook interface. The sidebar on the left includes links for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area has three command cells (Cmd 1, Cmd 2, Cmd 3). Cmd 1 contains Scala code to read a text file and count the rows. Cmd 2 contains code to take the first 3 rows of the RDD. Cmd 3 shows the resulting DataFrame output.

```

1 val studentData = sc.textFile("/FileStore/tables/StudentsPerformance.txt")
2 studentData.count()

(1) Spark Jobs
studentData: org.apache.spark.rdd.RDD[String] = /FileStore/tables/StudentsPerformance.txt MapPartitionsRDD[1] at textFile at command-1964254448130377:1
res2: Long = 2016

Command took 4.45 seconds -- by rshari125@uottawa.ca at 3/23/2020, 7:47:26 PM on Team#25 - Rahul&Harman

Cmd 2

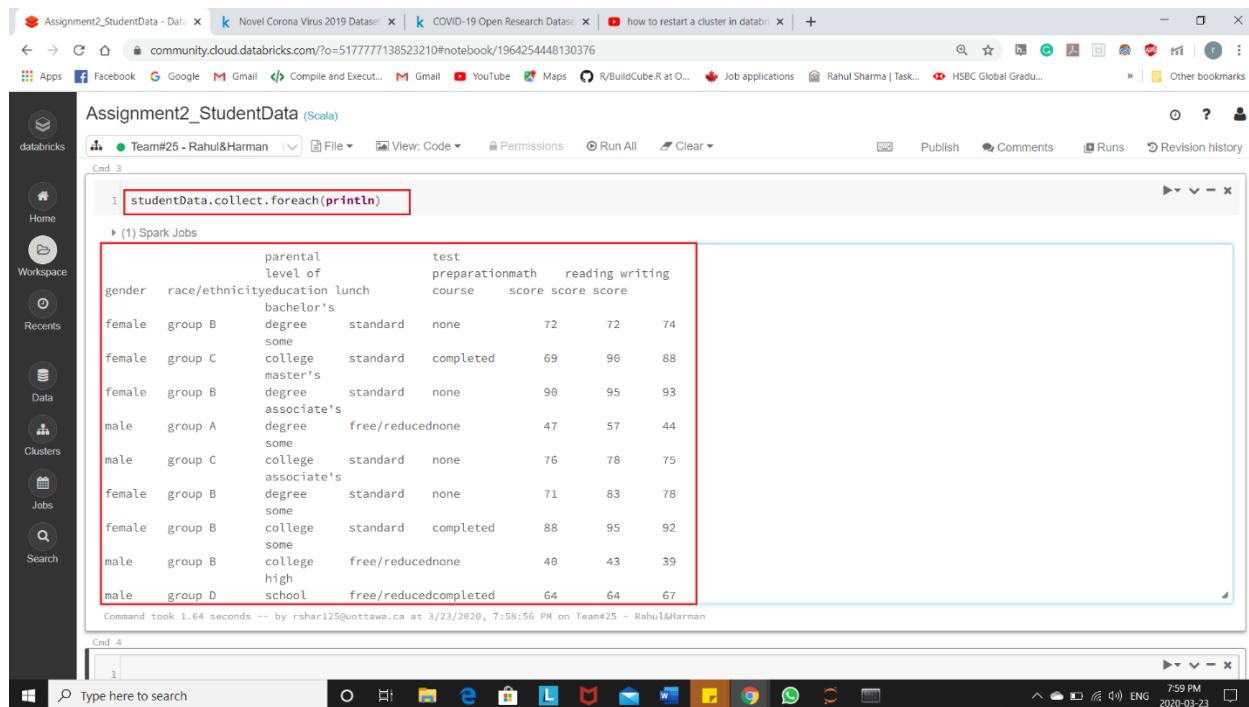
studentData.take(3)

(1) Spark Jobs
res: Array[String] = Array("ing", gender race/ethnicityeducation lunch parental course test", " level of preparationmath reading write
Command took 0.55 seconds -- by rshari125@uottawa.ca at 3/23/2020, 7:56:09 PM on Team#25 - Rahul&Harman

Cmd 3

```

- e) Print of all the lines in your data file (please provide a screenshot of only the lines that are visible on screen at one time) – Below screenshot shows the data of the file



The screenshot shows a Databricks notebook interface. The sidebar on the left includes links for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area has two command cells (Cmd 3, Cmd 4). Cmd 3 contains code to collect and foreach print the data. Cmd 4 shows the resulting DataFrame output.

```

1 studentData.collect.foreach(println)

(1) Spark Jobs
parental test
level of preparationmath reading writing
gender race/ethnicityeducation lunch bachelor's course score score score
female group B degree standard none 72 72 74
female group C college master's standard completed 69 90 88
female group B degree standard none 98 95 93
male group A degree some free/reducednone 47 57 44
male group C college standard none 76 78 75
female group B degree some standard none 71 83 78
female group B college associate's standard none 88 95 92
male group B college free/reducednone 49 43 39
male group D high school free/reducedcompleted 64 64 67

Command took 1.64 seconds -- by rshari125@uottawa.ca at 3/23/2020, 7:58:56 PM on Team#25 - Rahul&Harman

```

- f) Pick any two words that exist in your data file, filter the data file for lines with these words, and count the total occurrence of the words.

Below screenshot shows the total no. of occurrences for the word "group A" which is 89 highlighted first and second highlighted shows the no. of occurrences for word "group B" which is 190. We have chosen these 2 words to find out the category in data for these 2 groups

The screenshot shows a Databricks notebook interface with several tabs at the top: Assignment2_StudentData - Data, Novel Corona Virus 2019 Dataset, COVID-19 Open Research Dataset, how to restart a cluster in datab... The main area displays Scala code in a command cell (Cmd 4) and another cell (Cmd 5). The code filters a dataset for lines containing "group A" and "group B" respectively, then counts the results. The output shows 89 for group A and 190 for group B.

```
Assignment2_StudentData (Scala)
1 val linesWithGroupA = studentData.filter(line => line.contains("group A"))
2 linesWithGroupA.count()
1 val linesWithGroupB = studentData.filter(line => line.contains("group B"))
2 linesWithGroupB.count()
```

Cmd 4

```
1 val linesWithGroupA = studentData.filter(line => line.contains("group A"))
2 linesWithGroupA.count()
  ↗ (1) Spark Jobs
  LinesWithGroupA: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at filter at command-1964254448130381:1
  res7: Long = 89
  Command took 0.64 seconds -- by rshar125@uottawa.ca at 3/23/2020, 8:04:35 PM on Team#25 - Rahul&Harman
```

Cmd 5

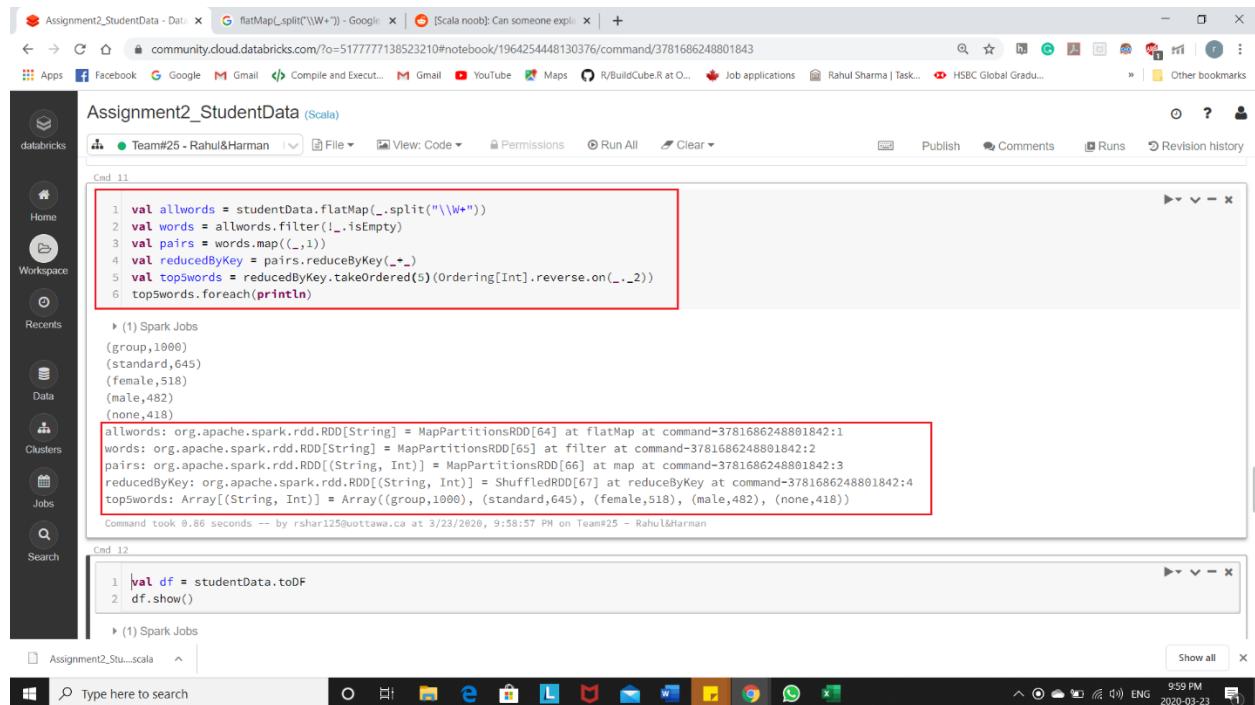
```
1 val linesWithGroupB = studentData.filter(line => line.contains("group B"))
2 linesWithGroupB.count()
  ↗ (1) Spark Jobs
  LinesWithGroupB: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[3] at filter at command-1964254448130382:1
  res8: Long = 190
  Command took 0.54 seconds -- by rshar125@uottawa.ca at 3/23/2020, 8:05:09 PM on Team#25 - Rahul&Harman
```

Shift+Enter to run shortcuts

Type here to search

Question 7 [6 pts]: Use any 4 Scala transformation functions in a pipeline format to extract some analytic information from your data file (from question 7). Please provide:

- a) A screenshot of the notebook with code showing the 4 transformation functions in a pipeline – Below screenshot shows the 4 transformation which we have used like flatMap(), filter(), map(), pairs()



The screenshot shows a Databricks notebook interface. The sidebar on the left includes icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area has two command cells. The first cell, labeled 'Cmd 11', contains Scala code to process student data and output a top 5 word count. The second cell, labeled 'Cmd 12', shows the resulting DataFrame 'df' being displayed. The bottom status bar shows the command took 0.86 seconds.

```

1 val allwords = studentData.flatMap(_.split("\\W+"))
2 val words = allwords.filter(_._isEmpty)
3 val pairs = words.map((_,1))
4 val reducedByKey = pairs.reduceByKey(_+_)
5 val topwords = reducedByKey.takeOrdered(5)(Ordering[Int].reverse.on(_._2))
6 topwords.foreach(println)

> (1) Spark Jobs
(group,1000)
(standard,645)
(female,518)
(male,482)
(None,418)

allwords: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[64] at flatMap at command-3781686248801842:1
words: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[65] at filter at command-3781686248801842:2
pairs: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[66] at map at command-3781686248801842:3
reducedByKey: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[67] at reduceByKey at command-3781686248801842:4
topwords: Array[(String, Int)] = Array((group,1000), (standard,645), (female,518), (male,482), (None,418))

Command took 0.86 seconds -- by rshar125@ottawa.ca at 3/23/2020, 9:58:57 PM on Team#25 - Rahul&Harman

1 val df = studentData.toDF
2 df.show()

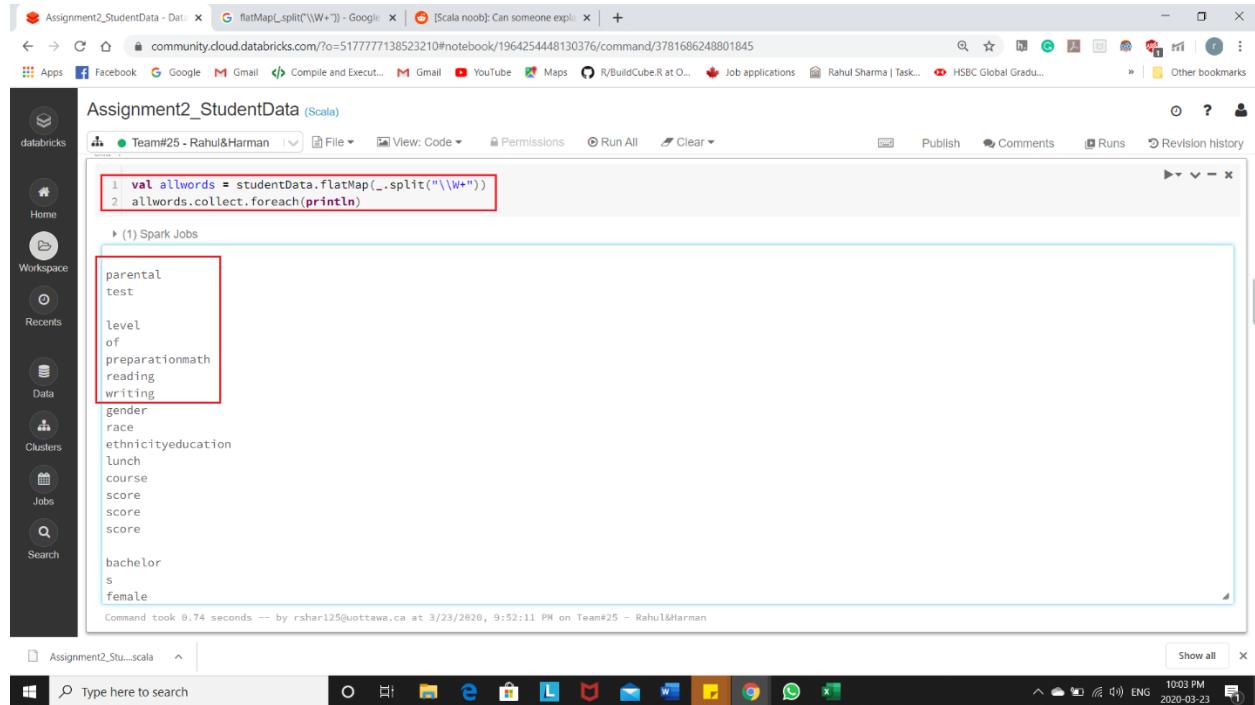
> (1) Spark Jobs

```

b) Provide a brief explanation of the transformations

Transformation 1

First transformation is done on flatMap() which will divide the whole file into single words and it will easy for us to read all the data by words



The screenshot shows a Databricks notebook interface. The sidebar on the left includes icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area shows the output of the flatMap transformation. The code `allwords.collect.foreach(println)` is highlighted. The output displays various student attributes like gender, race, ethnicity, and scores, each on a new line. The bottom status bar shows the command took 0.74 seconds.

```

1 val allwords = studentData.flatMap(_.split("\\W+"))
2 allwords.collect.foreach(println)

> (1) Spark Jobs
parental
test
level
of
preparationmath
reading
writing
gender
race
ethnicityeducation
lunch
course
score
score
score
score
bachelor
s
female

Command took 0.74 seconds -- by rshar125@ottawa.ca at 3/23/2020, 9:52:11 PM on Team#25 - Rahul&Harman

```

Transformation 2

We have used filter() in second transformation in order to filter out the empty lines from the file

The screenshot shows a Databricks notebook interface. The sidebar on the left includes icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area displays Scala code and its execution results:

```
1 val words = allwords.filter(!_._.isEmpty)
2 words.collect.foreach(println)
```

(1) Spark Jobs

```
parental
test
level
of
preparationmath
reading
writing
gender
race
ethnicityeducation
lunch
course
score
score
score
bachelor
s
female
group
B
degree
```

Command took 0.76 seconds -- by rshari125@ottawa.ca at 3/23/2020, 9:54:00 PM on Team#25 - Rahul&Harman

The output shows the words from the input file, excluding empty lines.

Transformation 3

In third transformation we have used map in order to pair all the word in key value pair like for each occurrence just write it as word and no. of occurrence as 1

The screenshot shows a Databricks notebook interface. The sidebar on the left includes icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area displays Scala code and its execution results:

```
1 val pairs = words.map((_,1))
2 pairs.collect.foreach(println)
```

(1) Spark Jobs

```
(parental,1)
(test,1)
(level,1)
(of,1)
(preparationmath,1)
(reading,1)
(writing,1)
(gender,1)
(race,1)
(ethnicityeducation,1)
(lunch,1)
(course,1)
(score,1)
(score,1)
(score,1)
(bachelor,1)
(s,1)
(female,1)
(group,1)
(B,1)
(degree,1)
```

Command took 0.73 seconds -- by rshari125@ottawa.ca at 3/23/2020, 9:54:55 PM on Team#25 - Rahul&Harman

The output shows the words from the input file, paired with their count of 1 for each occurrence.

Transformation 4

In transformation 4 we have used pair() in order to pair the same words and to find out the total no. of occurrences for each word

The screenshot shows a Jupyter Notebook interface in Databricks. The code cell contains:

```
1 val reducedByKey = pairs.reduceByKey(_+_)
2 reducedByKey.collect.foreach(println)
```

The output shows the results of the reduceByKey operation:

```
(88,35)
(score,3)
(reading,1)
(high,375)
(associate,222)
(82,62)
(89,58)
(19,2)
(42,20)
(Leve,1,1)
(53,54)

(B,190)
(62,84)
(Lunch,1)
(37,9)
(46,27)
(some,485)
(writing,1)
(86,39)
(93,18)
(28,3)
```

Below the code cell, it says "Command took 0.71 seconds -- by rshari125@uottawa.ca at 3/23/2020, 9:55:29 PM on Team#25 – Rahul&Harman".

c) A screenshot of the transformation data frame result

Below screenshot shows the transformed data after all the above 4 transformation is successfull

The screenshot shows a Jupyter Notebook interface in Databricks. The code cell contains:

```
1 val df = studentData.toDF("Student Data")
2 df.show()
```

The output shows the transformed data frame:

```
+-----+
| Student Data |
+-----+
| ... |
| gender race/eth... |
| ... |
| female group B ... |
| female group C ... |
| female group B ... |
| male group A ... |
| male group C ... |
| female group B ... |
| female group B ... |
| ... |
+-----+
```

Below the code cell, it says "Command took 1.33 seconds by rshari125@uottawa.ca at 3/23/2020, 10:00:41 PM on Team#25 – Rahul&Harman".

Question 8 [6 pts]: Please upload the selected csv file for your team into databricks and create a table. Create a new python notebook and display the data frame showing the first 30 rows using the Spark sql context. Using the same notebook, query your data for any 3 columns (including at least one summary data) from your data file and display the result data frame showing those columns. Convert the result into a dashboard of your choice. Please provide snapshots of the following activities:

- a) Data import confirmation page – Below screenshot shows the file is successfully uploaded in the cluster. Highlighted are the file name, size and link of the file

The screenshot shows the 'Create New Table' interface in Databricks. On the left is a sidebar with icons for Home, Workspace, Recents, Data, Clusters, Jobs, and Search. The main area has tabs for 'Upload File', 'S3', 'DBFS', and 'Other Data Sources'. Under 'Upload File', there's a 'File to DBFS' input field with the path '/FileStore/tables/' and an optional dropdown. Below it is a list of files. A red box highlights the file 'covid_19_data.csv' which has a green checkmark next to it. Another red box highlights the file size '0.5 MB'. A third red box highlights the success message 'File uploaded to /FileStore/tables/covid_19_data.csv'. At the bottom are buttons for 'Create Table with UI' and 'Create Table in Notebook'.

- b) Created table from imported data – Below screenshot shows the formation of the table with the uploaded data. Highlighted are the table name, file type and column delimiters

The screenshot shows the Databricks interface for creating a new table. The 'Table Name' field is set to 'covid_data'. The 'File Type' is selected as 'CSV'. The 'Column Delimiter' is set to a comma (','). Under schema settings, 'First row is header' and 'Infer schema' are checked. The 'Create Table' button is highlighted with a red box. A preview table is shown with 6 rows of data from China.

SNo	ObservationDate	ProvinceState	CountryRegion	LastUpdate	Confirmed
1	01/22/2020	Anhui	Mainland China	1/22/2020 17:00	1
2	01/22/2020	Beijing	Mainland China	1/22/2020 17:00	14
3	01/22/2020	Chongqing	Mainland China	1/22/2020 17:00	6
4	01/22/2020	Fujian	Mainland China	1/22/2020 17:00	1
5	01/22/2020	Gansu	Mainland China	1/22/2020 17:00	0
6	01/22/2020	Guangdong	Mainland China	1/22/2020 17:00	26

Below screenshot shows the table is created successfully

The screenshot shows the Databricks Table view for 'covid_data'. The table has been successfully created. The schema is displayed with columns: col_name, data_type, and comment. The sample data table shows 3 rows of data from China.

col_name	data_type	comment
SNo	int	null
ObservationDate	string	null
ProvinceState	string	null
CountryRegion	string	null
LastUpdate	string	null
Confirmed	double	null
Deaths	double	null
Recovered	double	null

SNo	ObservationDate	ProvinceState	CountryRegion	LastUpdate	Confirmed	Deaths	Recovered
1	01/22/2020	Anhui	Mainland China	1/22/2020 17:00	1	0	0
2	01/22/2020	Beijing	Mainland China	1/22/2020 17:00	14	0	0
3	01/22/2020	Chongqing	Mainland China	1/22/2020 17:00	6	0	0

- c) Your python notebook with code to create a data frame using the Spark sql context –

Below screenshot shows the python notebook is created successfully and the data frame is successfully created with the data of the table

```
Assignment2_covid-19 (Python)
1 df = sqlContext.sql("SELECT * FROM covid_data")
2
3 df: pyspark.sql.dataframe.DataFrame
4   SNo: integer
5   ObservationDate: string
6   ProvinceState: string
7   CountryRegion: string
8   LastUpdate: string
9   Confirmed: double
10  Deaths: double
11  Recovered: double
Command took 0.39 seconds -- by rshari25@uottawa.ca at 3/23/2020, 11:04:00 PM on Team#25 - Rahul&Harman
```

- d) Data frame showing the first 30 rows of your data file – Below screenshot shows the first 30 rows of the data which we have uploaded as a table

```
Assignment2_covid-19 (Python)
1 df.take(30)
2
3 Out[2]: [Row(SNo=1, ObservationDate='01/22/2020', ProvinceState='Anhui', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
4   Row(SNo=2, ObservationDate='01/22/2020', ProvinceState='Beijing', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=14.0, Deaths=0.0, Recovered=0.0),
5   Row(SNo=3, ObservationDate='01/22/2020', ProvinceState='Chongqing', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=6.0, Deaths=0.0, Recovered=0.0),
6   Row(SNo=4, ObservationDate='01/22/2020', ProvinceState='Fujian', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
7   Row(SNo=5, ObservationDate='01/22/2020', ProvinceState='Gansu', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0, Recovered=0.0),
8   Row(SNo=6, ObservationDate='01/22/2020', ProvinceState='Guangdong', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=26.0, Deaths=0.0, Recovered=0.0),
9   Row(SNo=7, ObservationDate='01/22/2020', ProvinceState='Guangxi', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=2.0, Deaths=0.0, Recovered=0.0),
10  Row(SNo=8, ObservationDate='01/22/2020', ProvinceState='Guizhou', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
11  Row(SNo=9, ObservationDate='01/22/2020', ProvinceState='Hainan', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=4.0, Deaths=0.0, Recovered=0.0),
12  Row(SNo=10, ObservationDate='01/22/2020', ProvinceState='Hebei', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
13  Row(SNo=11, ObservationDate='01/22/2020', ProvinceState='Heilongjiang', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0,
14  Row(SNo=12, ObservationDate='01/22/2020', ProvinceState='Henan', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
15  Row(SNo=13, ObservationDate='01/22/2020', ProvinceState='Hubei', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=10.0, Deaths=0.0, Recovered=0.0),
16  Row(SNo=14, ObservationDate='01/22/2020', ProvinceState='Jiangxi', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
17  Row(SNo=15, ObservationDate='01/22/2020', ProvinceState='Jiangsu', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
18  Row(SNo=16, ObservationDate='01/22/2020', ProvinceState='Liaoning', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
19  Row(SNo=17, ObservationDate='01/22/2020', ProvinceState='Ningxia', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0, Recovered=0.0),
20  Row(SNo=18, ObservationDate='01/22/2020', ProvinceState='Qinghai', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0, Recovered=0.0),
21  Row(SNo=19, ObservationDate='01/22/2020', ProvinceState='Shaanxi', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
22  Row(SNo=20, ObservationDate='01/22/2020', ProvinceState='Shandong', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
23  Row(SNo=21, ObservationDate='01/22/2020', ProvinceState='Shanghai', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
24  Row(SNo=22, ObservationDate='01/22/2020', ProvinceState='Sichuan', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
25  Row(SNo=23, ObservationDate='01/22/2020', ProvinceState='Tibet', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0, Recovered=0.0),
26  Row(SNo=24, ObservationDate='01/22/2020', ProvinceState='Xinjiang', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0, Recovered=0.0),
27  Row(SNo=25, ObservationDate='01/22/2020', ProvinceState='Yunnan', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=0.0, Deaths=0.0, Recovered=0.0),
28  Row(SNo=26, ObservationDate='01/22/2020', ProvinceState='Zhejiang', CountryRegion='Mainland China', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
29  Row(SNo=27, ObservationDate='01/22/2020', ProvinceState='Azerbaijan', CountryRegion='Azerbaijan', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
30  Row(SNo=28, ObservationDate='01/22/2020', ProvinceState='Brazil', CountryRegion='Brazil', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
31  Row(SNo=29, ObservationDate='01/22/2020', ProvinceState='Russia', CountryRegion='Russia', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0),
32  Row(SNo=30, ObservationDate='01/22/2020', ProvinceState='United States', CountryRegion='United States', LastUpdate='1/22/2020 17:00', Confirmed=1.0, Deaths=0.0, Recovered=0.0)]
Command took 0.52 seconds -- by rshari25@uottawa.ca at 3/23/2020, 11:15:33 PM on Team#25 - Rahul&Harman
```

Below screenshot shows the same 30 rows which we have retrieved using sql commands

The screenshot shows a Databricks notebook titled "Assignment2_covid-19 (Python)". In the code editor (Cmd 3), the following SQL command is run:

```
1 %sql
2 select * from covid_data
3 limit 30;
```

The resulting data frame is displayed as a table with columns: SNo, ObservationDate, ProvinceState, CountryRegion, LastUpdate, Confirmed, Deaths, and Recovered. The data is as follows:

SNo	ObservationDate	ProvinceState	CountryRegion	LastUpdate	Confirmed	Deaths	Recovered
1	01/22/2020	Anhui	Mainland China	1/22/2020 17:00	1	0	0
2	01/22/2020	Beijing	Mainland China	1/22/2020 17:00	14	0	0
3	01/22/2020	Chongqing	Mainland China	1/22/2020 17:00	6	0	0
4	01/22/2020	Fujian	Mainland China	1/22/2020 17:00	1	0	0
5	01/22/2020	Gansu	Mainland China	1/22/2020 17:00	0	0	0
6	01/22/2020	Guangdong	Mainland China	1/22/2020 17:00	26	0	0
7	01/22/2020	Guangxi	Mainland China	1/22/2020 17:00	2	0	0
8	01/22/2020	Gulzhou	Mainland China	1/22/2020 17:00	1	0	0
9	01/22/2020	Hainan	Mainland China	1/22/2020 17:00	4	0	0

Command took 0.56 seconds -- by rshar125@ottawa.ca at 3/23/2020, 11:18:48 PM on Team#25 - Rahul&Harman

- e) Code to create a data frame showing the 3 columns from your data file – Below is the screenshot shows the 3 columns selected to do the analysis of the data

The screenshot shows a Databricks notebook titled "Assignment2_covid-19 (Python)". In the code editor (Cmd 4), the following Python code is run:

```
1 df1 = sqlContext.sql("SELECT ProvinceState,CountryRegion,Confirmed FROM covid_data")
```

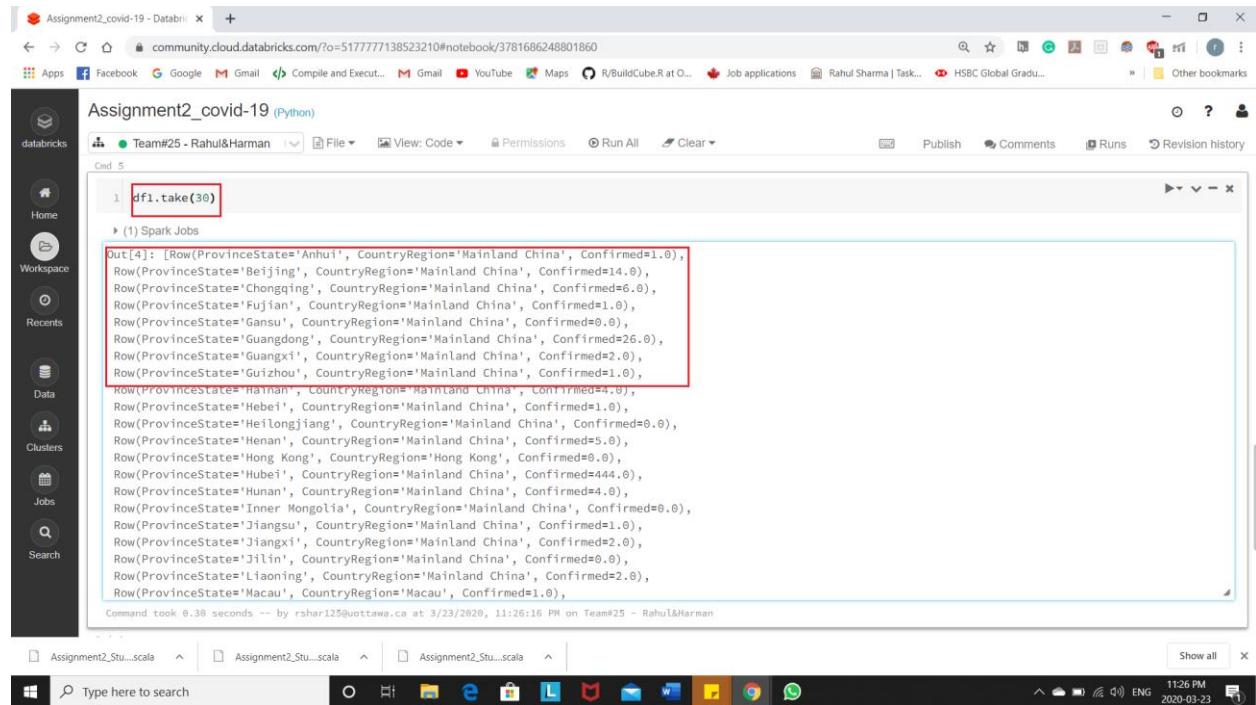
The resulting DataFrame is displayed as a table with columns: ProvinceState, CountryRegion, and Confirmed. The data is as follows:

ProvinceState	CountryRegion	Confirmed
Anhui	Mainland China	1
Beijing	Mainland China	14
Chongqing	Mainland China	6
Fujian	Mainland China	1
Gansu	Mainland China	0
Guangdong	Mainland China	26
Guangxi	Mainland China	2
Gulzhou	Mainland China	1
Hainan	Mainland China	4

df1: pyspark.sql.dataframe.DataFrame

Command took 0.25 seconds -- by rshar125@ottawa.ca at 3/23/2020, 11:21:24 PM on Team#25 - Rahul&Harman

Below screenshot shows the data associated with the 3 columns which we have selected for the analysis using data frame



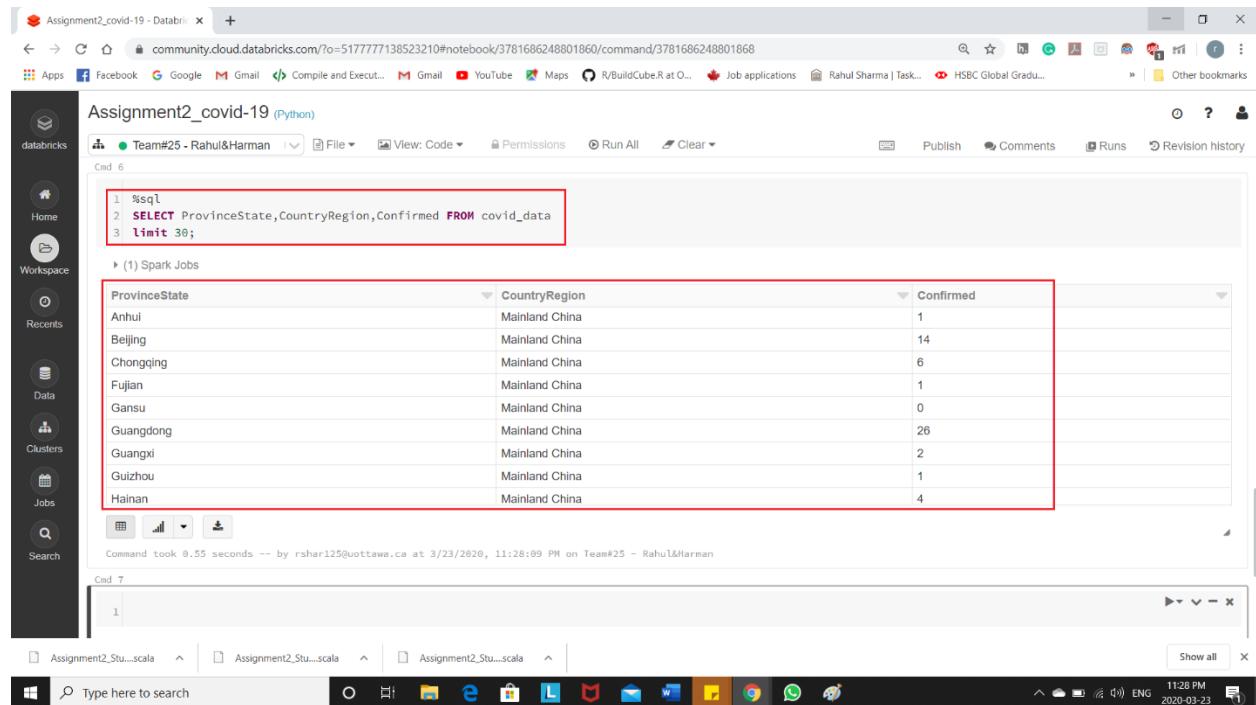
The screenshot shows a Databricks notebook titled "Assignment2_covid-19 - Python". In the command editor, the code `df1.take(30)` is run, resulting in a list of 30 rows. The output is highlighted with a red box. The rows represent data from Mainland China, with columns including ProvinceState, CountryRegion, and Confirmed cases.

```

1 df1.take(30)
▶ (1) Spark Jobs
Out[4]: [Row(ProvinceState='Anhui', CountryRegion='Mainland China', Confirmed=1.0),
Row(ProvinceState='Beijing', CountryRegion='Mainland China', Confirmed=14.0),
Row(ProvinceState='Chongqing', CountryRegion='Mainland China', Confirmed=6.0),
Row(ProvinceState='Fujian', CountryRegion='Mainland China', Confirmed=1.0),
Row(ProvinceState='Gansu', CountryRegion='Mainland China', Confirmed=0.0),
Row(ProvinceState='Guangdong', CountryRegion='Mainland China', Confirmed=26.0),
Row(ProvinceState='Guangxi', CountryRegion='Mainland China', Confirmed=2.0),
Row(ProvinceState='Guizhou', CountryRegion='Mainland China', Confirmed=1.0),
Row(ProvinceState='Hainan', CountryRegion='Mainland China', Confirmed=4.0),
Row(ProvinceState='Hebei', CountryRegion='Mainland China', Confirmed=1.0),
Row(ProvinceState='Heilongjiang', CountryRegion='Mainland China', Confirmed=0.0),
Row(ProvinceState='Henan', CountryRegion='Mainland China', Confirmed=5.0),
Row(ProvinceState='Hong Kong', CountryRegion='Hong Kong', Confirmed=0.0),
Row(ProvinceState='Hubei', CountryRegion='Mainland China', Confirmed=444.0),
Row(ProvinceState='Hunan', CountryRegion='Mainland China', Confirmed=4.0),
Row(ProvinceState='Inner Mongolia', CountryRegion='Mainland China', Confirmed=0.0),
Row(ProvinceState='Jiangsu', CountryRegion='Mainland China', Confirmed=1.0),
Row(ProvinceState='Jiangxi', CountryRegion='Mainland China', Confirmed=2.0),
Row(ProvinceState='Jilin', CountryRegion='Mainland China', Confirmed=0.0),
Row(ProvinceState='Liaoning', CountryRegion='Mainland China', Confirmed=2.0),
Row(ProvinceState='Macau', CountryRegion='Macau', Confirmed=1.0),
Command took 0.30 seconds -- by rshari25@ottawa.ca at 3/23/2020, 11:26:16 PM on Team#25 - Rahul&Harman

```

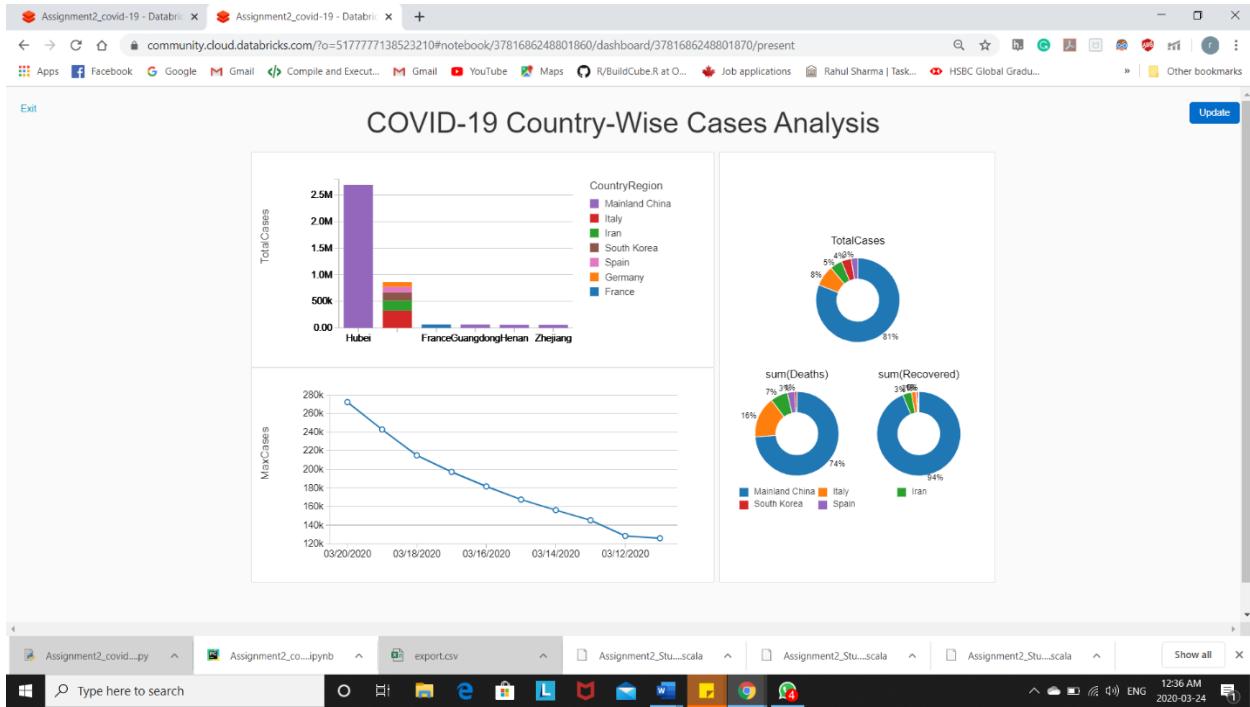
Below is the screenshot of the data associated with the 3 columns using sql commands



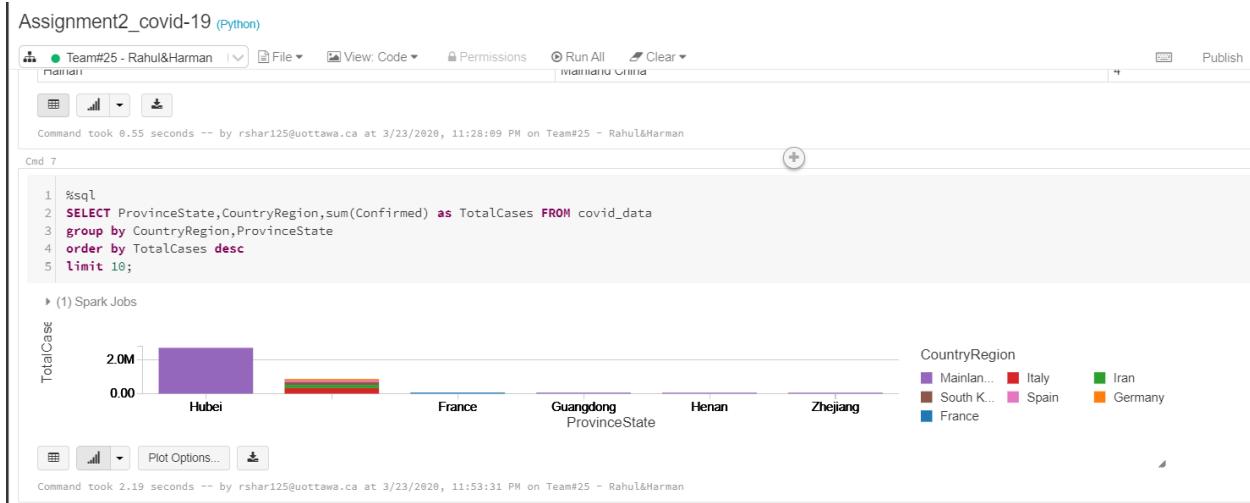
The screenshot shows a Databricks notebook titled "Assignment2_covid-19 - Python". In the command editor, the SQL query `SELECT ProvinceState, CountryRegion, Confirmed FROM covid_data limit 30;` is run, resulting in a table of 30 rows. The table is highlighted with a red box. The columns are ProvinceState, CountryRegion, and Confirmed.

ProvinceState	CountryRegion	Confirmed
Anhui	Mainland China	1
Beijing	Mainland China	14
Chongqing	Mainland China	6
Fujian	Mainland China	1
Gansu	Mainland China	0
Guangdong	Mainland China	26
Guangxi	Mainland China	2
Guizhou	Mainland China	1
Hainan	Mainland China	4

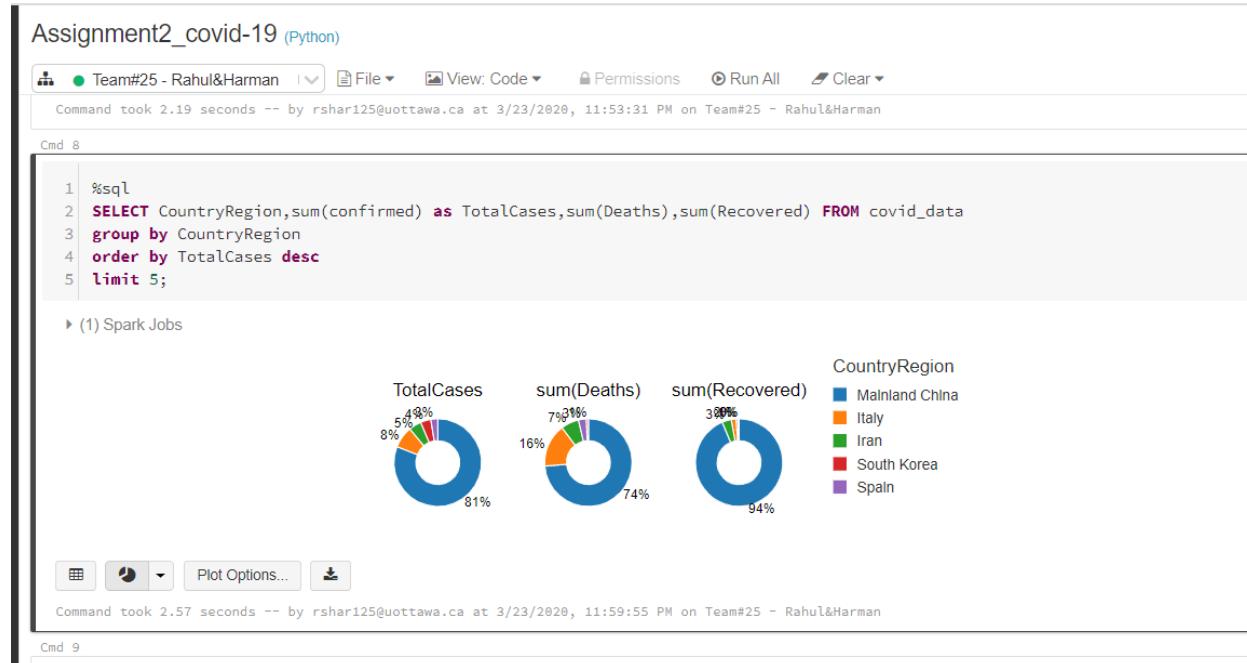
- f) Screenshot of a dashboard using the results from (e) – Below is the screenshot of the dashboard created with the 3 graphs



Below screenshot shows the first graph of the dashboard which basically shows the top-10 country and province which has total maximum cases of the covid-19



Below screenshot shows the top-5 countries with total no. of cases along with the cases which they have recovered and the total death of there particular region due to covid-19



Below screenshot shows the affect of covid -19 day by day

