



# ELG5166

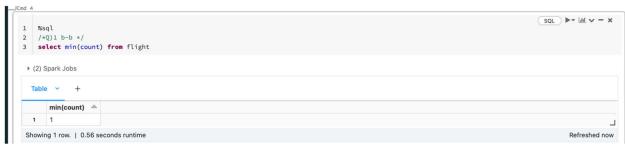
Spark APIs



NOVEMBER 3, 2022
ALI AMIN EL-SAYED MAHMOUD EL-SHERIF
300327246

#### Question 1:

```
1 /* Q) 1 */
  val sales = spark
     .option("inferSchema", "true")
4
     .option("header", "true")
      .csv("/FileStore/tables/Data/Sales.csv")
7 sales.count()
8 sales.createOrReplaceTempView("Sales")
10 sales.show(3)
 ▶ (5) Spark Jobs
 ▶ ■ sales: org.apache.spark.sql.DataFrame = [Transaction_date: string, Product: string ... 10 more fields]
|Transaction_date| Product|Price|Payment_Type|
                                                    Name
                                                                       City| State| Country|Account_Created| Last_Login|Latitude| Longitu
de
     1/2/09 6:17|Product1| 1200| Mastercard| carolina|
                                                                  Basildon|England|United Kingdom| 1/2/09 6:00| 1/2/09 6:08| 51.5|-1.11666
67|
     1/2/09 4:53|Product1| 1200| Visa|
                                                 Betina|Parkville
                                                                        ...| MO| United States| 1/2/09 4:42| 1/2/09 7:49| 39.195| -94.681
94|
    1/2/09 13:08|Product1| 1200| Mastercard|Federica e Andrea|Astoria
                                                                        ...| OR| United States| 1/1/09 16:21|1/3/09 12:32|46.18806| -123.
83|
+---
only showing top 3 rows
```



#### Question 2:

```
Cmd 5
     %python
 3 flight1 = spark.read.option("inferSchema", "true").option("header", "true").csv("/FileStore/tables/flight_data.csv")
 4 flight1.createOrReplaceTempView("flight_data_2015")
   12) Snark Johs
Cmd 6
 1 %sal
 2 /*Q)2 a-a */
 select sum(count) as total_number_flights,ORIGIN_COUNTRY_NAME from flight_data_2015 group by ORIGIN_COUNTRY_NAME order by total_number_flights limit 3
   ▶ (2) Spark Jobs
   Table v +
         total_number_flights  
ORIGIN_COUNTRY_NAME  

                                Singapore
    2
                                Croatia
  Showing all 3 rows. | 2.31 seconds runtime
                                                                                                                                                       Refreshed now
 1 /*Q)2 a-b */
 2 var query = spark.sql("select sum(count) as total_number_flights,ORIGIN_COUNTRY_NAME from flight_data_2015 group by ORIGIN_COUNTRY_NAME order by
      total_number_flights limit 3")
 3 query.show()
  \bullet \  \  \, \blacksquare \  \  \, \text{query: org.apache.spark.sql.DataFrame} = [\text{total\_number\_flights: long, ORIGIN\_COUNTRY\_NAME: string}]
  |total_number_flights|ORIGIN_COUNTRY_NAME|
                                   Singapore
                                   Lithuania
                      1|
                                    Croatia
 query: org.apache.spark.sql.DataFrame = [total_number_flights: bigint, ORIGIN_COUNTRY_NAME: string]
```

```
Scala FY V - X
2 var query = spark.sql("select sum(count) as total_number_flights,ORIGIN_COUNTRY_NAME from flight_data_2015 group by ORIGIN_COUNTRY_NAME order by
    total_number_flights limit 3")
3 query.show()
5 flightData
     .groupBy("ORIGIN_COUNTRY_NAME")
      .sort("total_number_flights")
      .limit(3)
11
12
 (2) Spark Jobs
 |total_number_flights|ORIGIN_COUNTRY_NAME|
                                Singapore|
                    11
                                Lithuania
                    1
                    1|
                                 Croatia
 ⊕ AnalysisException: Column 'total_number_flights' does not exist. Did you mean one of the following? [sum(count), ORIGIN_COUNTRY_NAME];
   'Sort ['total_number_flights ASC NULLS FIRST], true
   +- Aggregate [ORIGIN_COUNTRY_NAME#192], [ORIGIN_COUNTRY_NAME#192, sum(count#193) AS sum(count)#503L]
     +- Repartition 5, false
        +- Relation [DEST_COUNTRY_NAME#191,ORIGIN_COUNTRY_NAME#192,count#193] csv
```

```
Scala DT V - X
    /*Q)2 b */
2 var query = spark.sql("select sum(count) as total_number_flights,ORIGIN_COUNTRY_NAME from flight_data_2015 group by ORIGIN_COUNTRY_NAME order by
    total_number_flights limit 3")
    query.show()
4 query.explain()
 (2) Spark Jobs
 \bullet \  \  \, \blacksquare \  \  \, \text{query: org.apache.spark.sql.DataFrame} = [\text{total\_number\_flights: long, ORIGIN\_COUNTRY\_NAME: string}]
|total_number_flights|ORIGIN_COUNTRY_NAME|
                                        Singapore|
                         1|
                        1|
                                          Croatia
== Physical Plan ==
AdaptiveSparkPlan isFinalPlan=false
+- TakeOrderedAndProject(limit=3, orderBy=[total_number_flights#507L ASC NULLS FIRST], output=[total_number_flights#507L,ORIGIN_COUNTRY_NAME#421])
   +- HashAggregate(keys=[ORIGIN_COUNTRY_NAME#421], functions=[finalmerge_sum(merge sum#518L) AS sum(count#422)#508L])
+- Exchange hashpartitioning(ORIGIN_COUNTRY_NAME#421, 200), ENSURE_REQUIREMENTS, [id=#792]
           +- HashAggregate(keys=[ORIGIN_COUNTRY_NAME#421], functions=[partial_sum(count#422) AS sum#518L])
+- FileScan csv [ORIGIN_COUNTRY_NAME#421,count#422] Batched: false, DataFilters: [], Format: CSV, Location: InMemoryFileIndex(1 paths)[dbfs:/FileStore/tables/flight_data.csv], PartitionFilters: [], PushedFilters: [], ReadSchema: struct<ORIGIN_COUNTRY_NAME:string,count:int>
query: org.apache.spark.sql.DataFrame = [total_number_flights: bigint, ORIGIN_COUNTRY_NAME: string]
```

### Question 3:

- ▶ (2) Spark Jobs
- ▶ flightsDataFile: org.apache.spark.sql.DataFrame = [DEST\_COUNTRY\_NAME: string, ORIGIN\_COUNTRY\_NAME: string ... 1 more field]
- Fights: org.apache.spark.sql.Dataset[Flight] = [DEST\_COUNTRY\_NAME: string, ORIGIN\_COUNTRY\_NAME: string ... 1 more field]
- ▶ flights\_Metadata: org.apache.spark.sql.Dataset[FlightMetadata] = [count: long, random: long]

```
|count|
    0|-5790796081456756646|
    1 | -7406709754318658195 |
    2 | 3723754083696420047
    31-3678622977270750094
    4 | -8820342020588039680 |
    5 | 554201728505525082 |
    6|-4210117342112615418
    7 | 8640157511789465638 |
    8 | 5967035757421660660 |
    9 | 3968210518744556750
   10 | 3652462764861634911 |
   11 4179064758276548709
    12 | 3928348356421848608 |
   13 | -4441871604060770695 |
   14|-5151248049971105193|
   15 | -1522132841203857217 |
    16 4844185957588460211
   171-53120101167592655531
```

```
1 /*Q)3 a*/
    case class Flight(DEST_COUNTRY_NAME: String,
                        ORIGIN_COUNTRY_NAME: String, count: BigInt)
    /* Reads the Parquet file */
     val flightsDataFile = spark.read
      . \verb|parquet("/FileStore/tables/Data/flight-data/parquet/2010-summary.parquet/part_r_00000_1a9822ba_b8fb_4d8e_844a_ea30d0801b9e_gz.parquet")|
10 /* map the file to a typed dataset using Flight class*/
val flights = flightsDataFile.as[Flight]
12
13 /* create a structured dataset of the Type FlightMetadata */
14 case class FlightMetadata(count: BigInt, random: BigInt)
15
16 /* 500 entries with an index column and another with random numbers */
   val flights_Metadata = spark.range(500)
18
      .map(x => (x, scala.util.Random.nextLong))
      .withColumnRenamed("_1", "count")
.withColumnRenamed("_2", "random")
19
20
21
       .as[FlightMetadata]
22 flights_Metadata.show()
23
                                                                                                                                                    Scala ▶▼ ∨ - ×
   /*Q)3 b*/
 2
    val flights_2 = flights
       .joinWith(flights_Metadata, flights.col("count") === flights_Metadata.col("count"))
.withColumnRenamed("_1", "count")
.withColumnRenamed("_2", "random")
 6 flights_2.show()
```

```
|\{ \texttt{United States, U} ... | \{ \texttt{1, -208177811221} ... |
|{United States, F...|{1, -208177811221...
|{Bulgaria, United...|{1, -208177811221...|
 |{United States, S...|{1, -208177811221...
|{United States, C...|{1, -208177811221...
|{United States, B...|{1, -208177811221...
|{United States, I...|{1, -208177811221...
|{Vietnam, United ...|{1, -208177811221...
|{United States, P...|{1, -208177811221...
|{United States, B...|{1, -208177811221...
 |{United States, G...|{1, -208177811221...
 |{United States, L...|{1, -208177811221...
|{Malaysia, United...|{1, -208177811221...
|{United States, S...|{1, -208177811221...|
 |{United States, A...|{1, -208177811221...
|{United States, E...|{1, -208177811221...
|{United States, V...|{1, -208177811221...
```

## Question 4:

```
| SQL | SQL
```

```
1 /* Q)4 b*/
2 flights.rdd.getNumPartitions
resl8: Int = 1
```

# Question 5:

```
1 /* Q)5 a*/
     3 val textFile = spark.sparkContext.textFile("/FileStore/tables/Data/Adult.csv")
    4 textFile.count()
        ▶ (1) Spark Jobs
      textFile: org.apache.spark.rdd.RDD[String] = /FileStore/tables/Data/Adult.csv \ MapPartitionsRDD[108] \ at \ textFile \ at \ command-3522973057872627:3 \ and \ textFile \ at \ command-3522973057872627:3 \ at \ command-3522973057872627:3 \ at \ command-3522973057872627:3 \ at \ command-3522973057872627:3 \ at \ command-35229730572627:3 \ at \ command-352297305726
     res19: Long = 32562
                                                                                                                                                                                                                                                                                                                                                      Scala > V - X
     1 /* Q)5 b*/
            val count = textFile.flatMap(line => line.split(","))
                                                               .map(word => (word, 1))
                                                                 . {\tt reduceByKey(\_+\_).take0rdered(5)(0rdering[Int].on(\_.\_2))}\\
     6 count.foreach(println)
       ▶ (1) Spark Jobs
      (87,1)
      (marital_status,1)
      (Holand-Nether,1)
      (occupation.1)
      count: Array[(String, Int)] = Array((87,1), (marital_status,1), (age,1), (Holand-Nether,1), (occupation,1))
Question 6:
                                                                                                                                                                                                                                                                                                                                                       Scala > - - x
     1 /* Q)6 */
     3 val dataframe = spark.read.format("csv")
                .option("header", "true")
.option("inferSchema", "true")
.load("/FileStore/tables/2010_12_10.csv")
     8 dataframe.printSchema()
    dataframe.createOrReplaceTempView("dfTable")
        ▶ ■ dataframe: org.apache.spark.sql.DataFrame = [InvoiceNo: string, StockCode: string ... 6 more fields]
         |-- InvoiceNo: string (nullable = true)
         |-- StockCode: string (nullable = true)
         |-- Description: string (nullable = true)
         |-- Quantity: integer (nullable = true)
         |-- InvoiceDate: timestamp (nullable = true)
         |-- UnitPrice: double (nullable = true)
          |-- CustomerID: double (nullable = true)
         |-- Country: string (nullable = true)
```

dataframe: org.apache.spark.sql.DataFrame = [InvoiceNo: string, StockCode: string ... 6 more fields]

```
Cmd 17
                                                                                                                                                               Scala > - - x
    /*0)6 a*/
 3 import org.apache.spark.sql.functions.col
 d dataframe.where(col("Quantity").equalTo(12) && col("UnitPrice").gt(2))
5    .select("*")
        .show(5, false)
 8
  ▶ (1) Spark Jobs
 |InvoiceNo|StockCode|Description
                                                                  |Quantity|InvoiceDate
                                                                                                  |UnitPrice|CustomerID|Country
                                                                             |2010-12-10 09:33:00|2.55
  |538172
             84558A | 3D DOG PICTURE PLAYING CARDS
                                                                  12
                                                                             |2010-12-10 09:33:00|2.95
                                                                                                                15805.0
                                                                                                                             |United Kingdom|
                        ROUND SNACK BOXES SET OF4 WOODLAND |12
  1538174
             122326
                                                                             |2010-12-10 09:35:00|2.95
                                                                                                                112471.0
                                                                                                                             | Germany
  |538174 |22472
                        TV DINNER TRAY DOLLY GIRL
                                                                             |2010-12-10 09:35:00|4.95
                                                                                                                12471.0
                                                                  12
                                                                                                                             Germany
                       PENCIL CASE LIFE IS BEAUTIFUL
                                                                             |2010-12-10 09:35:00|2.95
 only showing top 5 rows
 import org.apache.spark.sql.functions.col
Cmd 18
                                                                                                                                                               Scala > V - X
 1 /*Q)6 b*/
import org.apache.spark.sql.functions.split
dataframe.select(split(col("Description"), " ").alias("Detailed description")).show(false)
  ▶ (1) Spark Jobs
  |Detailed description
  |[HAWAIIAN, GRASS, SKIRT, ]
  |[CHILLI, LIGHTS]
  [["RECORD, FRAME, 7"", SINGLE, SIZE, "]
  [3D, DOG, PICTURE, PLAYING, CARDS]
  [60, CAKE, CASES, VINTAGE, CHRISTMAS]
  |[PAPER, CHAIN, KIT, VINTAGE, CHRISTMAS]
|[CHRISTMAS, TOILET, ROLL]
  ["ASSORTED, FLOWER, COLOUR, ""LEIS"""]
  [I, CAN, ONLY, PLEASE, ONE, PERSON, MUG]
[HAND, WARMER, BABUSHKA, DESIGN]
  [BLUE, HARMONICA, IN, BOX, ]
  null
  [SET, OF, 72, RETROSPOT, PAPER, , DOILIES]
  |[WOODLAND, DESIGN, , COTTON, TOTE, BAG]
  |[BIG, DOUGHNUT, FRIDGE, MAGNETS]
|[RECYCLED, PENCIL, WITH, RABBIT, ERASER]
Cmd 19
 1 /*Q)6 2-a*/
 2 ls /FileStore/tables
   Table v +
         path
                                              name
                                                                 size
                                                                                 modificationTime

        1
        dbfs:/FileStore/tables/2010_12_01.csv
        2010_12_01.csv
        275001

        2
        dbfs:/FileStore/tables/2010_12_02.csv
        2010_12_02.csv
        191826

        3
        dbfs:/FileStore/tables/2010_12_10.csv
        2010_12_10.csv
        241468

                                                                                     1667513508000
                                                                                     1667513508000
                                                                                     1667513718000
                                            Data/ 0
   4 dbfs:/FileStore/tables/Data/
                                                                                    0
                                                                   7080
    5 dbfs:/FileStore/tables/flight_data.csv
                                                 flight_data.csv
                                                                                    1666034682000
  Showing all 5 rows. | 1.13 seconds runtime
                                                                                                                                                              Refreshed 2 minutes ago
```

## Question 7:

```
Scala ▶▼ ∨ - ×
    /* Q)7 */
3 val dataset = spark.read.format("csv")
4 .option("header", "true")
       .option("inferSchema", "true")
       .load("/FileStore/tables/2010_12_01.csv")
 8 dataset.createOrReplaceTempView("RetailView")
 9 dataset.count()
  ▶ (4) Spark Jobs
  → 🔳 dataset: org.apache.spark.sql.DataFrame = [InvoiceNo: string, StockCode: string ... 6 more fields]
  dataset: org.apache.spark.sql.DataFrame = [InvoiceNo: string, StockCode: string ... 6 more fields]
 res29: Long = 3108
Cmd 23
                                                                                                                                                          Scala > - - x
    /*Q)7 a*/
      val datadescription= dataset.filter(x =>
                            val allStr = x.getString(2).split(" ")
                            var check = true
 6
                             for(x <- allStr ){</pre>
                               if (!x(0).isUpper)
                                   check=false
 10
 11
 12
 13
                           (check)
 14
 15
 16 datadescription.show(5,false)
  ▶ (1) Spark Jobs
  ▶ ■ datadescription: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [InvoiceNo: string, StockCode: string ... 6 more fields]
  |InvoiceNo|StockCode|Description
                                                                |Quantity|InvoiceDate
                                                                                                 |UnitPrice|CustomerID|Country
  |536365
            |85123A |WHITE HANGING HEART T-LIGHT HOLDER |6
                                                                           |2010-12-01 08:26:00|2.55
                                                                                                             17850.0
                                                                                                                         |United Kingdom|
            | WHITE METAL LANTERN | 6 |
| 84406B | CREAM CUPID HEARTS COAT HANGER | 8 |
| 84029G | KNITTED UNION FLAG HOT WATER BOTTLE| 6
                                                                           |2010-12-01 08:26:00|3.39
                                                                                                             |17850.0
                                                                                                                         |United Kingdom|
                                                                        |2010-12-01 08:26:00|3.39

|2010-12-01 08:26:00|2.75

|2010-12-01 08:26:00|3.39
  |536365
                                                                                                             |17850.0
                                                                                                                         |United Kingdom|
                                                                                                                         |United Kingdom|
  1536365
                                                                                                             117850.0
  |536365 |84029E |RED WOOLLY HOTTIE WHITE HEART.
                                                                          |2010-12-01 08:26:00|3.39
                                                                                                             |17850.0
                                                                                                                         |United Kingdom|
  only showing top 5 rows
```