# Name/ Hosam Mahmoud Ibrahim Mahmoud.
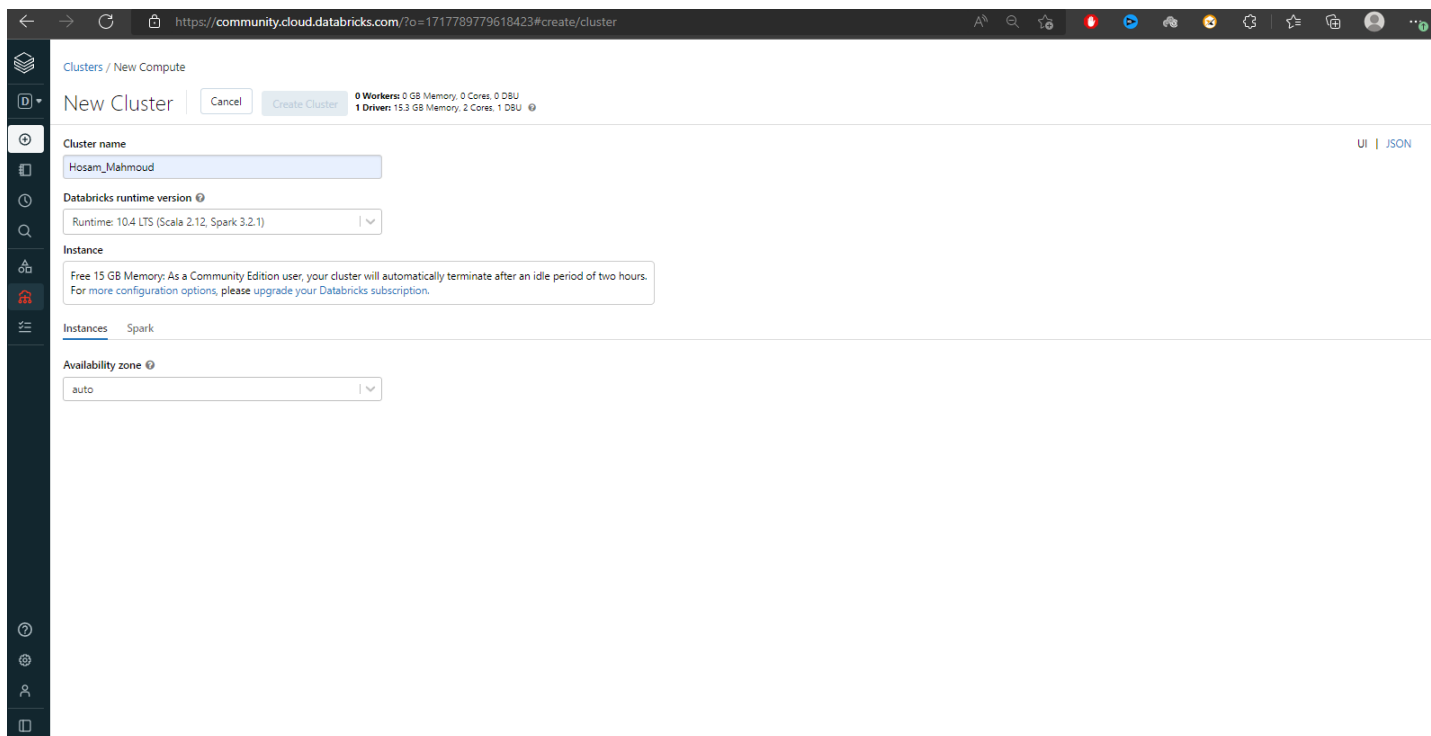
Please use the provided *HTML* Files to see the code and the output, these images are just to show the cluster name for each notebook.

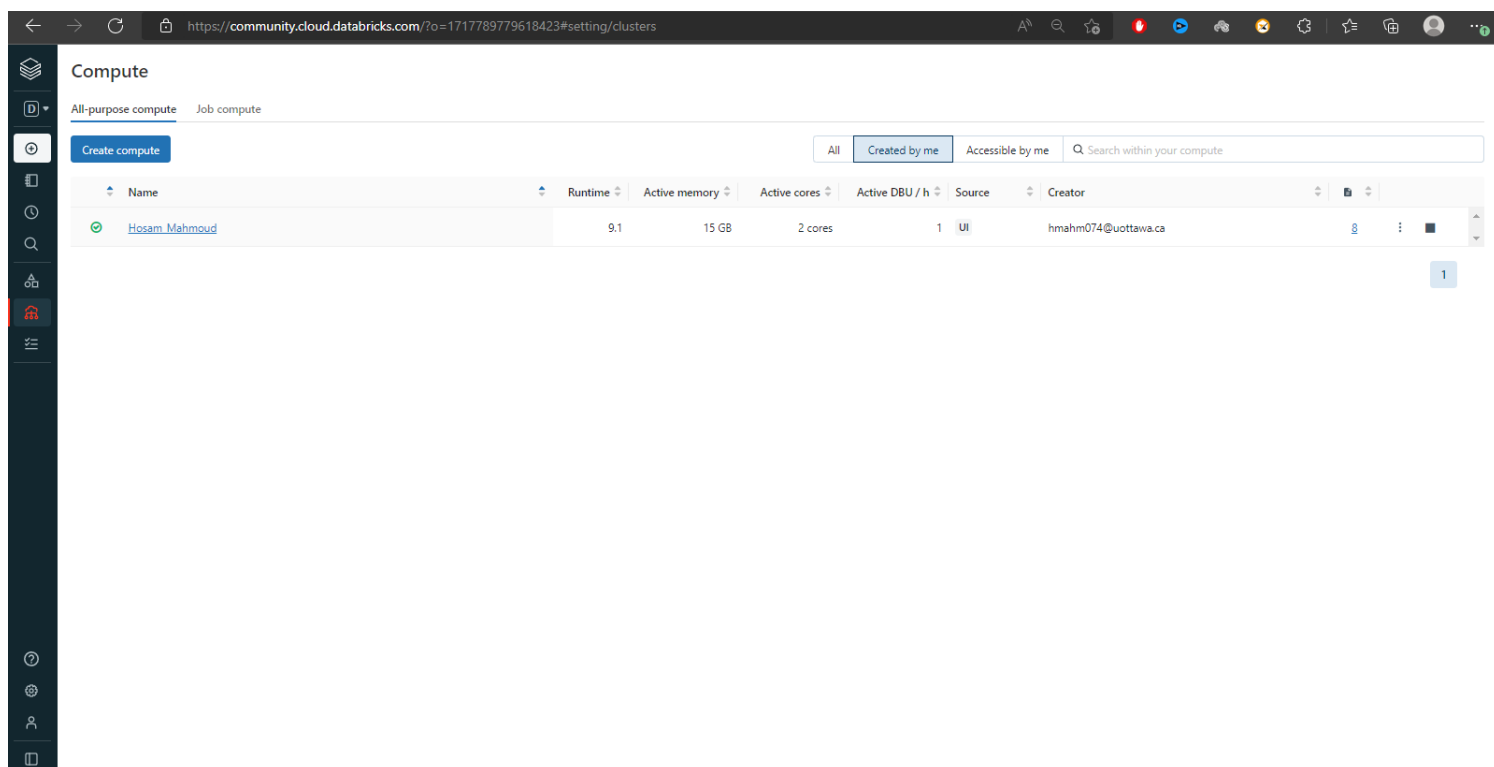*Note please don't forget to Scroll to see the full output.*

1- Cluster creation.



2- Cluster Created successfully

## 3- Q1

**Q1 (1)** Scala

File  Edit  View  Run  Help    Last edit was 2 hours ago    Give feedback

Run all   • Hosam_Mahmoud   Publish

Cmd 1

```scala
1  // Print the first 3 records of sales data. - Hosam Mahmoud
2  val salesData = spark
3    .read
4    .option("inferSchema", "true")
5    .option("header", "true")
6    .csv("/FileStore/tables/Sales.csv")
7
8  salesData.show(3)
```

▸ (3) Spark Jobs

▸ ▦ salesData: 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| Longitude|
+---------------+--------+-----+------------+----------------+--------+-----+--------------+---------------+--------------+--------+----------+
|    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|-1.1166667|
|    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| -94.68194|
|   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

salesData: org.apache.spark.sql.DataFrame = [Transaction_date: string, Product: string ... 10 more fields]
```

Command took 7.39 seconds -- by hmahm074@uottawa.ca at 11/3/2022, 3:02:27 AM on Hosam_Mahmoud

Cmd 2

Scala

```scala
1  // Print the minimum count of flights in 2 ways from flight data. - Hosam Mahmoud
2  val flightData = spark
3    .read
4    .option("inferSchema", "true")
5    .option("header", "true")
6    .csv("/FileStore/tables/flight_data.csv")
7
8  // first way
9  import org.apache.spark.sql.functions.{min,max}
10 flightData.agg(max($"count"), min($"count")).show()
11
12 // second way
13 flightData.agg(max(flightData(flightData.columns(2))),min(flightData(flightData.columns(2)))).show()
```

## 4- Q2

**Q2 (1)** Python

File  Edit  View  Run  Help    Last edit was 2 minutes ago    Give feedback

Run all   • Hosam_Mahmoud   Publish

Cmd 1

Python

```python
1  # In three ways, do the following:
2  # Print the total number of the flights counts grouped by the origin country, show only the three minimum counts of these flights. - Hosam Mahmoud
3
4  flightData2015 = spark\
5    .read\
6    .option("inferSchema", "true")\
7    .option("header", "true")\
8    .csv("/FileStore/tables/flight_data.csv")
9  flightData2015.createOrReplaceTempView("flight_data_2015")
10
11 print(flightData2015.head(5))
12
13 # first way
14 flightData2015\
15   .groupBy("ORIGIN_COUNTRY_NAME")\
16   .sum("count")\
17   .withColumnRenamed("sum(count)", "Origin_total")\
18   .sort("Origin_total")\
19   .limit(3)\
20   .show()
21
22 # second way
23 maxSql = spark.sql("""
24 SELECT ORIGIN_COUNTRY_NAME, sum(count) as Origin_total
25 FROM flight_data_2015
26 GROUP BY ORIGIN_COUNTRY_NAME
27 ORDER BY sum(count)
28 LIMIT 3
29 """)
30
31 maxSql.show()
32
33 # Print the execution plan for one of these ways.
34 print(maxSql.explain())
35 print(flightData2015.explain())
```

▸ (7) Spark Jobs

[Row(DEST_COUNTRY_NAME='United States', ORIGIN_COUNTRY_NAME='Romania', count=15), Row(DEST_COUNTRY_NAME='United States', ORIGIN_COUNTRY_NAME='Croatia', count=1), Row(DEST_COUNTRY_NAME='United States', ORIGIN_CO
UNTRY_NAME='Ireland', count=344), Row(DEST_COUNTRY_NAME='Egypt', ORIGIN_COUNTRY_NAME='United States', count=15), Row(DEST_COUNTRY_NAME='United States', ORIGIN_COUNTRY_NAME='India', count=62)]

## 5- Q3

Q3 (1)  Scala ∨

File  Edit  View  Run  Help   Last edit was 2 hours ago   Give feedback

Run all   • Hosam_Mahmoud ∨   Publish

Cmd 1

```scala
1   // Create a structured data set with the name "Flights_info" using the map function, this data set includes two columns, "random" column that holds a random variable with the type integer, and "count"
    column
2   // with range of (500). ------- Hosam Mahmoud
3
4   case class Flight(DEST_COUNTRY_NAME: String,
5                     ORIGIN_COUNTRY_NAME: String, count: BigInt)
6   val flightsDF = spark.read
7     .parquet("/FileStore/tables/part_r_00000_1a9822ba_b8fb_4d8e_844a_ea30d0801b9e_gz.parquet")
8   val flights = flightsDF.as[Flight]
9
10  flights.show(5)
11
12  case class FlightMetadata(count: BigInt, random: BigInt)
13
14  val flightsMeta = spark.range(500)
15    .map(x => (x, scala.util.Random.nextInt(1000)))
16    .withColumnRenamed("_1", "count")
17    .withColumnRenamed("_2", "random")
18    .as[FlightMetadata]
19
20  flightsMeta.show(5)
21
22  // Using "count" column, join the "Flights_info" data set to the Flight data coming from this path "/FileStore/tables/flightdata/parquet/part_r_00000_1a9822ba_b8fb_4d8e_844a_ea30d0801b9e_gz.parquet"
23  // ------- Hosam Mahmoud
24
25  val flights2 = flights
26    .joinWith(flightsMeta, flights.col("count") === flightsMeta.col("count"))
27    .withColumnRenamed("_1", "count")
28    .withColumnRenamed("_2", "random")
29
30  flights2.show(10, false)
```

▸ (4) Spark Jobs

▸ ▦ flightsDF: org.apache.spark.sql.DataFrame = [DEST_COUNTRY_NAME: string, ORIGIN_COUNTRY_NAME: string ... 1 more field]
▸ ▦ flights: org.apache.spark.sql.Dataset[Flight] = [DEST_COUNTRY_NAME: string, ORIGIN_COUNTRY_NAME: string ... 1 more field]
▸ ▦ flightsMeta: org.apache.spark.sql.Dataset[FlightMetadata] = [count: long, random: integer]
▸ ▦ flights2: org.apache.spark.sql.DataFrame = [count: struct, random: struct]

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
```

## 6- Q4

Q4 (1)  SQL ∨

File  Edit  View  Run  Help   Last edit was 2 hours ago   Give feedback

Run all   • Hosam_Mahmoud ∨   Publish

Cmd 1

```sql
1   -- Create a table that contains all the information in the flights data. (1 pt)
2   -- Note that, flights data contains all the json files in this path '/FileStore/tables/json/*.json
3   -- ########################### Hosam Mahmoud
4
5   DROP TABLE IF EXISTS flights;
6   CREATE TABLE flights (
7     DEST_COUNTRY_NAME STRING, ORIGIN_COUNTRY_NAME STRING, count LONG)
8   USING JSON OPTIONS (path '/FileStore/tables/*.json')
```

OK

Command took 2.24 seconds -- by hmahm074@uottawa.ca at 11/3/2022, 4:49:22 AM on Hosam_Mahmoud

Cmd 2

```sql
1   %sql
2   select * from flights
```

▸ (2) Spark Jobs

Table ∨  +

| | DEST_COUNTRY_NAME | ORIGIN_COUNTRY_NAME | count |
|---|---|---|---|
| 1 | United States | Romania | 15 |
| 2 | United States | Croatia | 1 |
| 3 | United States | Ireland | 344 |
| 4 | Egypt | United States | 15 |
| 5 | United States | India | 62 |
| 6 | United States | Singapore | 1 |
| 7 | United States | Grenada | 62 |

Truncated results, showing first 1,000 rows. ∨  |  1.58 seconds runtime                                          Refreshed now

Cmd 3

```
1   -- i've read the data from 6 json files (partitions)
```

Command took 0.00 seconds -- by hmahm074@uottawa.ca at 11/3/2022, 4:49:23 AM on Hosam_Mahmoud

## 7- Q5

**Q5 (1)** Scala ⌄

File  Edit  View  Run  Help    Last edit was 1 minute ago    Give feedback

Run all    • Hosam_Mahmoud ⌄    Publish

Cmd 1

```scala
1  // Show the five words that have the least number of occurences in Adult data. ---- Hosam Mahmoud
2
3  val textFile = spark.sparkContext.textFile("/FileStore/tables/Adult.csv")
4
5  // print the text lines
6  textFile.collect.foreach(println)
```

▸ (1) Spark Jobs

```
age,education,marital_status,occupation,race,sex,native_country
37,HS-grad,Divorced,?,Asian-Pac-Islander,Female,?
39,Masters,Married-civ-spouse,?,Asian-Pac-Islander,Female,?
24,HS-grad,Never-married,?,Asian-Pac-Islander,Female,?
28,Some-college,Never-married,Adm-clerical,Asian-Pac-Islander,Female,?
26,Assoc-acdm,Married-spouse-absent,Craft-repair,Asian-Pac-Islander,Female,?
50,Some-college,Divorced,Exec-managerial,Asian-Pac-Islander,Female,?
45,Bachelors,Married-civ-spouse,Exec-managerial,Asian-Pac-Islander,Female,?
45,Bachelors,Married-civ-spouse,Exec-managerial,Asian-Pac-Islander,Female,?
44,HS-grad,Married-civ-spouse,Exec-managerial,Asian-Pac-Islander,Female,?
43,5th-6th,Married-civ-spouse,Machine-op-inspct,Asian-Pac-Islander,Female,?
20,Some-college,Never-married,Machine-op-inspct,Asian-Pac-Islander,Female,?
42,7th-8th,Married-civ-spouse,Other-service,Asian-Pac-Islander,Female,?
43,Bachelors,Married-civ-spouse,Other-service,Asian-Pac-Islander,Female,?
25,Assoc-voc,Never-married,Other-service,Asian-Pac-Islander,Female,?
20,HS-grad,Never-married,Other-service,Asian-Pac-Islander,Female,?
37,HS-grad,Divorced,Prof-specialty,Asian-Pac-Islander,Female,?
30,Prof-school,Never-married,Prof-specialty,Asian-Pac-Islander,Female,?
48,HS-grad,Divorced,Sales,Asian-Pac-Islander,Female,?
37,12th,Married-civ-spouse,Sales,Asian-Pac-Islander,Female,?
31,12th,Married-civ-spouse,Sales,Asian-Pac-Islander,Female,?
```

Command took 5.10 seconds -- by hmahm074@uottawa.ca at 11/3/2022, 4:50:09 AM on Hosam_Mahmoud

Cmd 2

```scala
1  // print the occurences of each word  ---- Hosam Mahmoud
2
3  val counts = textFile.flatMap(line => line.split(","))
4                  .map(word => (word, 1))
5                  .reduceByKey(_+_)
6  counts.collect.foreach(println)
```

## 8- Q6

**Q6 (1)** Scala ⌄

File  Edit  View  Run  Help    Last edit was 2 hours ago    Give feedback

Run all    • Hosam_Mahmoud ⌄    Publish

Cmd 1

```scala
1   // Print the first 5 records that have "Quantity" of 12 and their "unitPrice" > 2 ------ Hosam Mahmoud
2
3   val df = spark.read.format("csv")
4     .option("header", "true")
5     .option("inferSchema", "true")
6     .load("/FileStore/tables/2010_12_01.csv")
7
8   // Show the inferred schema
9   df.printSchema()
10  // Create a view from the data frame.
11  df.createOrReplaceTempView("dfTable")
12
13  import org.apache.spark.sql.functions.col
14  df.where(col("Quantity").equalTo(12) && col("UnitPrice").gt(2))
15    .select("*")
16    .show(5, false)
17
18  // Split the "Description" column and re-name it as "Detailed description" ----- Hosam Mahmoud
19  import org.apache.spark.sql.functions.split
20  df.select(split(col("Description"), " ").alias("Detailed description")).show(false)
```

▸ (4) Spark Jobs

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

```
root
 |-- InvoiceNo: string (nullable = true)
 |-- StockCode: string (nullable = true)
 |-- Description: string (nullable = true)
 |-- Quantity: integer (nullable = true)
 |-- InvoiceDate: string (nullable = true)
 |-- UnitPrice: double (nullable = true)
 |-- CustomerID: double (nullable = true)
 |-- Country: string (nullable = true)
```

| InvoiceNo | StockCode | Description | Quantity | InvoiceDate | UnitPrice | CustomerID | Country |
|---|---|---|---|---|---|---|---|
| 536370 | 22726 | ALARM CLOCK BAKELIKE GREEN | 12 | 2010-12-01 08:45:00 | 3.75 | 12583.0 | France |
| 536370 | 21913 | VINTAGE SEASIDE JIGSAW PUZZLES | 12 | 2010-12-01 08:45:00 | 3.75 | 12583.0 | France |
| 536378 | 84997B | RED 3 PIECE RETROSPOT CUTLERY SET | 12 | 2010-12-01 09:37:00 | 3.75 | 14688.0 | United Kingdom |

## 9- Q7

**Q7 (1)** Scala ∨
File  Edit  View  Run  Help  Last edit was now  Give feedback
Run all  ● Hosam_Mahmoud ∨  Publish

Cmd 1

```scala
1   // Print all the partitions in this path "/FileStore/tables/Retail" through a file system command. ---------- Hosam Mahmoud
```

Command took 3.40 seconds -- by hmahm074@uottawa.ca at 11/3/2022, 4:54:08 AM on Hosam_Mahmoud

Cmd 2

```
1   %fs
2   ls /FileStore/tables/Retail
```

Table ∨  +

| | path | name | size |
|---|---|---|---|
| 1 | dbfs:/FileStore/tables/Retail/2010_12_01.csv | 2010_12_01.csv | 275001 |
| 2 | dbfs:/FileStore/tables/Retail/2010_12_02.csv | 2010_12_02.csv | 191826 |
| 3 | dbfs:/FileStore/tables/Retail/2010_12_03.csv | 2010_12_03.csv | 190700 |
| 4 | dbfs:/FileStore/tables/Retail/2010_12_05.csv | 2010_12_05.csv | 246056 |
| 5 | dbfs:/FileStore/tables/Retail/2010_12_06.csv | 2010_12_06.csv | 339039 |
| 6 | dbfs:/FileStore/tables/Retail/2010_12_12.csv | 2010_12_12.csv | 132120 |
| 7 | dbfs:/FileStore/tables/Retail/2010_12_13.csv | 2010_12_13.csv | 202021 |

Showing all 301 rows. | 1.10 seconds runtime          Refreshed now

Cmd 3

```scala
1   //For the second partition, calculate the sum, min, max and avg of the "UnitPrice" column. ----------- Hosam
2
3   val ds = spark.read.format("csv")
4     .option("header", "true")
5     .option("inferSchema", "true")
6     .load("/FileStore/tables/Retail/2010_12_02.csv")
7     .coalesce(5)
8   ds.cache()
9   ds.createOrReplaceTempView("dfTable")
10
11  import org.apache.spark.sql.functions.{min, max, sum, avg}
12  ds.select(sum("UnitPrice"),min("UnitPrice"), max("UnitPrice"), avg("UnitPrice")).show(false)
```

▶ (4) Spark Jobs

## 10- Q8

**Q8** Scala ∨
File  Edit  View  Run  Help  Last edit was 26 minutes ago  Give feedback
Run all  ● Hosam_Mahmoud ∨  Publish

Cmd 1

```scala
1   // For the data in this path "/FileStore/tables/retail-data/by-day/2010_12_01.csv", do the following: ------ Hosam Mahmoud
2
3   val RetailDF = spark.read.format("csv")
4     .option("header", "true")
5     .option("inferSchema", "true")
6     .load("/FileStore/tables/Retail/2010_12_01.csv")
7
8   RetailDF.printSchema()
9   RetailDF.createOrReplaceTempView("dfTable")
10
```

▶ (2) Spark Jobs
▶ ▦ RetailDF: org.apache.spark.sql.DataFrame = [InvoiceNo: string, StockCode: string ... 6 more fields]
```
root
 |-- InvoiceNo: string (nullable = true)
 |-- StockCode: string (nullable = true)
 |-- Description: string (nullable = true)
 |-- Quantity: integer (nullable = true)
 |-- InvoiceDate: string (nullable = true)
 |-- UnitPrice: double (nullable = true)
 |-- CustomerID: double (nullable = true)
 |-- Country: string (nullable = true)

RetailDF: org.apache.spark.sql.DataFrame = [InvoiceNo: string, StockCode: string ... 6 more fields]
```
Command took 2.90 seconds -- by hmahm074@uottawa.ca at 11/3/2022, 4:11:53 AM on Hosam_Mahmoud

Cmd 2

```scala
1   RetailDF.show(5)
```

▶ (1) Spark Jobs

```
+---------+---------+--------------------+--------+-------------------+---------+----------+--------------+
|InvoiceNo|StockCode|         Description|Quantity|        InvoiceDate|UnitPrice|CustomerID|       Country|
+---------+---------+--------------------+--------+-------------------+---------+----------+--------------+
|   536365|   85123A|WHITE HANGING HEA...|       6|2010-12-01 08:26:00|     2.55|   17850.0|United Kingdom|
|   536365|    71053| WHITE METAL LANTERN|       6|2010-12-01 08:26:00|     3.39|   17850.0|United Kingdom|
|   536365|   84406B|CREAM CUPID HEART...|       8|2010-12-01 08:26:00|     2.75|   17850.0|United Kingdom|
|   536365|   84029G|KNITTED UNION FLA...|       6|2010-12-01 08:26:00|     3.39|   17850.0|United Kingdom|
```