# $as signment \hbox{-} 7 \hbox{-} py spark \hbox{-} data frame \hbox{-} 1$

## September 18, 2024

[1]:	sc
[1]:	<pre><sparkcontext appname="PySparkShell" master="local[*]"></sparkcontext></pre>
[2]:	spark
[2]:	<pre><pyspark.sql.session.sparksession 0x7fd6540dc3c8="" at=""></pyspark.sql.session.sparksession></pre>
[3]:	sc.stop()
	a) Create a new Spark Session with new SparkConfig
[4]:	<pre>from pyspark import SparkConf, SparkContext # setMaster() - set spark context manager which is local[cpu_cores] config = SparkConf().setMaster("local[4]").setAppName("PySparkSession") sc = SparkContext(conf=config)</pre>
[5]:	sc
[5]:	<pre><sparkcontext appname="PySparkSession" master="local[4]"></sparkcontext></pre>
	b) Create new instance of Spark SQL session and define new DataFrame using sales_data_sample.csv dataset.
[6]:	<pre>from pyspark.sql import SparkSession spark = SparkSession.builder.appName("SparkSQLSession").getOrCreate()</pre>
[7]:	spark
[7]:	<pre><pyspark.sql.session.sparksession 0x7fd638c0d7b8="" at=""></pyspark.sql.session.sparksession></pre>
[3]:	sales_df = spark.read.csv("file:///home/hadoop/Downloads/sales_data_sample.  csv",header=True,inferSchema=True)
[4]:	sales_df.show()
	+++++
	+

```
|ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES|
                                                              ORDERDATE
STATUS | QTR_ID | MONTH_ID | YEAR_ID | PRODUCTLINE | MSRP | PRODUCTCODE |
CUSTOMERNAME |
                     PHONE |
                                  ADDRESSLINE1 | ADDRESSLINE2 |
                                                                  CITY
STATE | POSTAL CODE | COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTFIRSTNAME | DEALSIZE |
+-----
_____
+-----
_____
                       301
                              95.71
                                                2 | 2871.0 | 2/24/2003
0:00|Shipped|
                             2003|Motorcycles| 95|
                        2|
                                                    S10_1678|
                1|
                                                               Land of
Toys Inc.
              2125557818|897 Long Airport ...|
                                                 null|
                                                               NYCl
NY
       10022|
                  USAI
                            NA
                                          Yul
                                                        Kwail
                                                               Small|
                       34 l
                             81.35
                                                5 | 2765.9 | 5/7/2003
      10121
0:00|Shipped|
                2|
                        51
                             2003|Motorcycles| 95|
                                                    S10_1678| Reims
Collectables
                26.47.1555| 59 rue de l'Abbaye|
                                                     nulll
null
         51100|
                 France
                            EMEA I
                                        Henriot
                                                          Paul
                                                                 Small|
      10134|
1
                       41|
                             94.74
                                                2|3884.34| 7/1/2003
0:00|Shipped|
                        7|
                             2003|Motorcycles| 95|
                                                   S10_1678|
                3|
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... |
                                                 null
                                                              Parisl
                 Francel
                                       Da Cunhal
         75508 l
                            EMEA I
                                                        Daniel | Medium |
                       45 l
                              83.26
                                                6| 3746.7| 8/25/2003
10145
0:00|Shipped|
                3|
                        81
                             2003|Motorcycles| 95|
                                                    S10 1678|
Toys4GrownUps.com
                     6265557265 | 78934 Hillside Dr.
                                                          null
Pasadenal
             CAI
                    900031
                               USAL
                                         NAI
                                                     Young |
Julie| Medium|
                                               14|5205.27|10/10/2003
      10159|
                       49|
                              100.0
                             2003|Motorcycles| 95|
                                                    S10_1678 | Corporate Gift
0:00|Shipped|
                4|
                       10|
Id...l
                       7734 Strong St.|
                                             null|San Francisco|
         6505551386
                                                                    CAL
null
         USAI
                   NA
                               Brown
                                               Julie | Medium |
      10168|
                       36 L
                              96.66
                                                1|3479.76|10/28/2003
0:00|Shipped|
                4|
                       10|
                             2003|Motorcycles| 95|
                                                    S10_1678|Technics
Stores Inc.
                65055568091
                             9408 Furth Circle
                                                    null
                                                            Burlingame |
CAI
       942171
                  USAI
                            NAI
                                       Hiranol
                                                        Juri | Medium |
      10180|
                       29|
                                                9|2497.77|11/11/2003
                             86.13
0:00|Shipped|
                4|
                       11|
                             2003|Motorcycles| 95|
                                                    S10_1678|Daedalus
              20.16.1555|184, chausse de T...|
Designs ...
                                                 null
null
         59000 l
                 Francel
                           EMEA
                                         Rancel
                                                       Martine
      101881
                       48 l
                              100.0
                                                1|5512.32|11/18/2003
                       11|
                             2003|Motorcycles| 95|
0:00|Shipped|
                4|
                                                    S10 1678|
                                                                   Herkku
      +47 2267 3215|Drammen 121, PR 7...|
                                             null
                                                                  null
                                                        Bergen
N 5804|
         Norway|
                   EMEA |
                                Oeztan
                                                Veysel| Medium|
                       22|
                              98.57
                                                2|2168.54| 12/1/2003
      10201
0:00|Shipped|
                       121
                             2003|Motorcycles| 95|
                                                    S10 1678|
                4|
               6505555787|5557 North Pendal...|
Wheels Co.
                                                 null|San Francisco|
CAI
        null
                  USAI
                            NA
                                       Murphy|
                                                       Juliel
                                                               Small|
      10211
                       41|
                              100.0
                                               14|4708.44| 1/15/2004
0:00|Shipped|
                1|
                       1|
                             2004|Motorcycles| 95| S10_1678| Auto Canal
```

```
Petit| (1) 47.55.6555|
                       25, rue Lauriston
                                              null
                                                          Paris|
                                                                   nulll
75016| France|
                  EMEA I
                                            Dominique | Medium |
                              Perrier|
                             100.01
                                               1|3965.66| 2/20/2004
      102231
                       37 l
0:00|Shipped|
                        2|
                            2004|Motorcycles| 95|
                                                   S10_1678|Australian
                1|
Collec...| 03 9520 4555|
                         636 St Kilda Road|
                                              Level 31
Melbourne|Victoria|
                      3004|Australia|
                                        APAC|
                                                   Ferguson|
Peter | Medium |
      10237 l
1
                       231
                             100.01
                                               7 | 2333.12 | 4/5/2004
0:00|Shipped|
                        4|
                            2004|Motorcycles| 95|
                                                   S10 1678|
                2|
                                                   Suite 101|
Vitachrome Inc.
                                2678 Kingston Rd.
                   2125551500
NYCl
        NYI
                10022|
                          USAI
                                    NA
                                                Frick|
                                                              Michael|
Small
                       281
                             100.0|
                                               2|3188.64| 5/18/2004
      10251
0:00|Shipped|
                        5|
                            2004|Motorcycles| 95|
                                                   S10_1678|Tekni
                2|
                                 7476 Moss Rd.|
Collectable...
                 2015559350
                                                     null
                                                                Newark
NJI
       940191
                 USAI
                           NAI
                                       Brownl
                                                     Williaml
                                                             Medium
      10263|
                       34 I
                             100.0|
                                               2|3676.76| 6/28/2004
0:00|Shipped|
                        61
                            2004|Motorcycles| 95|
                                                   S10_1678|
                2|
Depot Inc.
               2035552570 | 25593 South Bay Ln. |
                                                   null | Bridgewater |
CTI
       975621
                 USAI
                           NAI
                                        Kingl
                                                       Juliel Mediuml
                                               1|4177.35| 7/23/2004
10275
                       45 l
                             92.831
0:00|Shipped|
                        7|
                            2004|Motorcycles| 95|
                                                   S10 1678|
                                                              La Rochelle
                3|
          40.67.8555|67, rue des Cinqu...|
                                            null
                                                       Nantes
        Francel
                              Labrune
44000 l
                  EMEA I
                                              Janine | Medium |
10285 l
                       36 L
                             100.01
                                               6|4099.68| 8/27/2004
0:00|Shipped|
                3|
                        8|
                            2004|Motorcycles| 95|
                                                   S10_1678|Marta's
                 6175558555 | 39323 Spinnaker Dr.
                                                             Cambridge|
Replicas Co.|
                                                     null
MA
                 USAI
       51247
                           NA I
                                    Hernandez|
                                                       Martal
                                                             Medium
                                               9|2597.39| 9/30/2004
1
      10299
                       23|
                             100.0
0:00|Shipped|
                3|
                        91
                            2004|Motorcycles| 95|
                                                   S10_1678|Toys of
Finland, Co.
                90-224 8555
                                 Keskuskatu 451
                                                     nulll
                                                              Helsinkil
                           EMEA I
null
        21240 | Finland
                                     Karttunen|
                                                        Mattil
                                                                Small
      103091
                       41|
                             100.0|
                                               5|4394.38|10/15/2004
1
0:00|Shipped|
                       10|
                            2004|Motorcycles| 95|
                                                   S10_1678| Baane Mini
                4|
Imports |
            07-98 9555 Erling Skakkes ga...
                                              null|
                                                        Stavern
null
                           EMEA |
                                                        Jonas | Medium |
         4110|
                 Norway|
                                    Bergulfsen|
      10318
                             94.74
                                               1|4358.04| 11/2/2004
                       46|
0:00|Shipped|
                41
                       11|
                            2004|Motorcycles| 95|
                                                   S10 1678|Diecast
Classics ...
               2155551555
                            7586 Pompton St.|
                                                   null
                                                           Allentown
PAI
       702671
                 USAL
                           NAI
                                          Yul
                                                       Kyung | Medium |
1
      10329 l
                       421
                             100.0
                                               1|4396.14|11/15/2004
0:00|Shipped|
                            2004|Motorcycles| 95|
                                                   S10_1678|
                4|
                       11|
                                                              Land of
Toys Inc. |
              2125557818|897 Long Airport ...|
                                                null
                                                              NYCl
NYI
                 USAI
                            NAI
                                                       Kwai| Medium|
       100221
                                          Yul
+-----
_____
```

only showing top 20 rows

1|

31

68|

33|

25%|

482.13|

2203.11

10180

c) Find the shape of DataFrame. [5]: # Number of rows sales\_df.count() [5]: 2823 [6]: # Number of columns len(sales df.columns) [6]: 25 d) Find the Summary of DataFrame for all numerical data columns. [7]: from pyspark.sql.types import IntegerType, StringType from pyspark.sql.functions import \* [8]: numerical\_cols = [field.name for field in sales\_df.schema.fields if not\_ ⇔isinstance(field.dataType, StringType)] [9]: sales\_df[numerical\_cols].summary().show() |summary| ORDERNUMBER | QUANTITYORDERED | PRICEEACH | SALES ORDERLINENUMBER | QTR\_ID| MONTH\_ID| YEAR\_ID| MSRP -----| count| 28231 28231 2823 I 2823 l 28231 28231 28231 2823 l mean | 10258.725115125753 | 35.09280906836698 | 83.65854410201929|6.466170740347148| 3553.88907190932|2.7176762309599716|7.0924 548352816155 | 2003.8150903294368 | 100.71555083244775 | stddev| 92.0854775957196| 9.74144273706958|20.174276527840536| 4.22584096469094 | 1841.8651057401842 | 1.203878088001756 | 3.656633307661765 | 0.6996701541300869 | 40.18791167720266 | 6 I ı minl 10100 26.88

27|

1|

2|

20031

2003

1|

41

68.8

```
ı
   50%|
               10262
                              35|
                                          95.7
61
                          31
                                        81
                                                   20041
         3184.8|
991
75%|
                              43|
                                         100.0|
               10334
91
                                                   20041
         4508.01
                          4|
                                       111
124 l
   max |
               10425
                              97|
                                         100.0|
18 l
          14082.81
                           41
                                        121
                                                    2005 I
214
----+
```

e) Identify and handle missing or null values in the columns.

```
[10]: from pyspark.sql.functions import *
  sales_df.select([sum(col(column).isNull().cast('int')).alias(column) for column_u
   →in sales_df.columns]).show()
  _____
  |ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER|SALES|ORDERDATE|STATUS|QT
  R ID | MONTH ID | YEAR ID | PRODUCTLINE | MSRP | PRODUCTCODE | CUSTOMERNAME | PHONE | ADDRESSLIN
  E1 | ADDRESSLINE2 | CITY | STATE | POSTALCODE | COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTF
  IRSTNAME | DEALSIZE |
  ----+
  I
        01
               0|
                    0|
                            01
                               01
                                    01
                                       0|
       01
                                  0|
  01
          01
                0|
                   01
                         0|
                               01
                                        01
  2521
      0 | 1486 |
              76|
                  01
                       01
                                       01
  01
  ____+___
```

```
sales_df1 = sales_df.fillna("Nill")
```

----+

\_\_+\_\_\_\_

[12]: from pyspark.sql.functions import \*

```
sales_df1.select([sum(col(column).isNull().cast('int')).alias(column) for_

¬column in sales_df1.columns]).show()
__+____
|ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER|SALES|ORDERDATE|STATUS|QT
R ID | MONTH ID | YEAR ID | PRODUCTLINE | MSRP | PRODUCTCODE | CUSTOMERNAME | PHONE | ADDRESSLIN
E1 | ADDRESSLINE2 | CITY | STATE | POSTALCODE | COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTF
IRSTNAME | DEALSIZE |
____+___
__+____
----+
     01
            01
                01
                        01
                          01
                                  01
01
01
    01
       01
             01
               01
                     01
                           01
                             01
                                   01
01
  01
     01
          0|
             01
                  01
                         01
                                 01
01
____+___
__+____
```

f) Calculate the total revenue generated per country by combining the columns QUANTITY-ORDERED and PRICEEACH using Spark DataFrame operations?

```
[13]: sales_df2 = sales_df1.withColumn('Total_

→Revenue',col('QUANTITYORDERED')*col('PRICEEACH'))

sales_df2.groupBy('COUNTRY').agg(sum('Total Revenue').alias('Total Revenue')).

→show()
```

```
+----+
               Total Revenue
    COUNTRY
 ----+
     Sweden | 174264.10000000006 |
|Philippines| 80291.1699999998|
  Singapore | 227985.5000000001|
    Germany
                   178689.08
     France | 919257.8499999997
    Belgium|
                    94528.88
    Finland 268714.70000000007
      Italy| 309402.8699999999|
     Norway | 246115.8000000001|
      Spain|1021705.9700000002|
    Denmark
                   192747.63
    Ireland|
                    43237.24
```

----+

```
USA | 2986425.2099999995 |
UK | 413203.33999999997 |
Switzerland | 93344.909999999999 |
Canada | 193504.34000000003 |
Japan | 153076.68999999994 |
Australia | 521598.45999999985 |
Austria | 172793.05000000002 |
```

g) Determine the top 5 products with the highest total sales revenue using Spark DataFrame?

h) Find the average order quantity for each product using groupBy and agg operations?

```
[15]: sales_df1.groupBy('PRODUCTLINE').agg(mean('QUANTITYORDERED').alias('Avg order_uquantity')).show()
```

i) Using Spark DataFrame, filter orders where the SALES value exceeds \$10,000 and sort the results by the ORDERDATE column?

```
[16]: sales_df1 = sales_df1.withColumn("ORDERDATE", to_timestamp(col("ORDERDATE"), "M/

d/yyyy H:mm"))
    sales_df1.filter(sales_df1.SALES > 10000).orderBy(sales_df1.ORDERDATE).show()
    +-----
    __+____
    ______
    ______
    |ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES|
               STATUS|QTR_ID|MONTH_ID|YEAR_ID| PRODUCTLINE|MSRP|PRODUCTCODE|
    ORDERDATE
                                ADDRESSLINE1 | ADDRESSLINE2 |
    CUSTOMERNAME |
                      PHONE |
    CITY|STATE|POSTALCODE|
    COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTFIRSTNAME | DEALSIZE |
    +-----
    __+____
    -----+
    10127
                        46|
                              100.0|
                                              2|11279.2|2003-06-03
    00:00:00
                              6|
                                   2003 | Classic Cars | 207 |
             Shipped|
                       2|
    Muscle Machine Incl
                       2125557413|
                                  4092 Furth Circle
                                                  Suite 400
    NYCl
                100221
                         USAl
                                  NAI
                                            Young
                                                          Jeffl
    Large
                              100.01
         10150
                        45 l
                                              8|10993.5|2003-09-19
    00:00:00
             Shipped
                       3|
                              9|
                                  2003 | Classic Cars | 214 |
                                                       S10_1949|Dragon
    Souveniers... | +65 221 7555 | Bronz Sok., Bronz... |
                                               Nill
                                                      Singapore |
    Nill
            79903|Singapore|
                                     Natividad|
                                                      Eric|
                           Japan|
                                                             Large
                        441
    1
         10247
                              100.01
                                              2|10606.2|2004-05-05
    00:00:00
             Shipped |
                       2|
                              5|
                                  2004 Classic Cars | 207
                                                       S12 1108
    Suominen Souveniers | +358 9 8045 555 | Software Engineer... |
                                                     Nill
    Espoo | Nill | FIN-02271 | Finland |
                                  EMEA |
                                           Suominen
                                                           Kalle|
    Large
    47|
                              100.0
                                              6|10172.7|2004-10-11
         10304
    00:00:00
                                  2004 | Classic Cars | 214 |
             Shipped
                       4|
                              10|
                                                       S10_1949| Auto
                   30.59.8555|67, avenue de l'E...|
    Assoc. & Cie.
                                                 Nill
                                                       Versailles|
    Nill
            780001
                   France
                            EMEA |
                                       Tonini|
                                                     Daniel|
                                                             Large
         103121
                        481
                              100.01
                                              3 | 11623.7 | 2004-10-21
                                                       S10_1949|Mini
    00:00:00
                                  2004 | Classic Cars | 214 |
             Shipped
                       4|
                              10|
    Gifts Distri...
                   4155551450
                               5677 Strong St.|
                                                  Nill|
                                                         San Rafaell
                            NAI
    CAL
          975621
                   USAI
                                      Nelsonl
                                                  Valarie|
                                                           Large
                        50 l
                              100.01
                                              6|12536.5|2004-11-04
         103221
    00:00:00
                                  2004|Vintage Cars| 127|
                                                       S18_2325|Online
             Shipped
                       4|
                              11|
    Diecast Cr...
                 6035558647|2304 Long Airport...|
                                               Nill
                                                         Nashua|
    NHI
          62005 l
                   USA|
                            NA |
                                      Young
                                                  Valarie|
                                                           Large
         10333|
                        461
                              100.0
                                              2|11336.7|2004-11-18
```

11|

00:00:00

Shipped|

4|

2004|Vintage Cars| 99| S18\_3320|

```
Mini Wheels Co.| 6505555787|5557 North Pendal...| Nill|San Francisco|
             USA | NA | Murphy |
CAI
     Nill|
                                          Julie| Large|
                 55| 100.0|
                              13|10758.0|2004-11-23
     10339 l
00:00:00
        Shipped 4
                       11| 2004|Vintage Cars| 88| S24_3151|Tokyo
Collectable...|+81 3 3584 0555| 2-2-8 Roppongi| Nill| Minato-
ku|Tokyo| 106-0032| Japan| Japan| Shimamura|
                                               Akikol
Large
                 431
                       100.0|
10375
                                     2|10039.6|2005-02-03
00:00:00| Shipped| 1|
                      2| 2005| Planes| 72|
                                              S24 4278| La
Rochelle Gifts | 40.67.8555|67, rue des Cinqu...| Nill|
                                                 Nantesl
                     EMEA| Labrune|
Nill|
       44000| France|
                                            Janine|
                                                  Large
    10388|
                 46|
                       100.0|
                                     2|10066.6|2005-03-03
                      3| 2005| Planes| 91| S700_1691|
00:00:00| Shipped| 1|
FunGiftIdeas.com | 5085552555 | 1785 First Street | Nill | New
                             NA |
Bedford | MA|
             50553|
                     USA
                                     Benitez
Large
    10403|
                 66| 100.0|
                                     9|11886.6|2005-04-08
00:00:00| Shipped| 2|
                     4| 2005| Motorcycles| 193|
Collectables, ... (171) 555-2282 | Berkeley Gardens ... | Nill |
Liverpool | Nill | WX1 6LT | UK | EMEA |
                                    Devonl
                                                Elizabeth
Large
                 76| 100.0|
                                     3|11739.7|2005-04-14
     10405
00:00:00| Shipped| 2|
                      4| 2005|Classic Cars| 140| S24_3856|
Mini Caravy | 88.60.1555 | 24, place Kluber | Nill |
                                              Strasbourg
       67000| France| EMEA| Citeaux| Frederique| Large|
Nill
                 65| 100.0|
    10406|
                                     1|10468.9|2005-04-15
00:00:00| Disputed| 2|
                      4| 2005|Classic Cars| 141| S18_3685|Danish
Wholesale ... | 31 12 3555 | Vinb'ltet 34 | Nill |
                                             Kobenhavn
Nill| 1734| Denmark| EMEA| Petersen|
                                            Jvttel
                 76| 100.0|
                              2|14082.8|2005-04-22
00:00:00| On Hold| 2|
                      4| 2005|Vintage Cars| 170| S18_1749|The
Sharp Gifts W...| 4085553659|
                        3086 Ingle Ln.| Nill|
                        NA| Frick|
Jose| CA| 94217| USA|
                                                  Suel
Large
    10412
                 60| 100.0|
                                     9|11887.8|2005-05-03
1
00:00:00| Shipped | 2 | 5 | 2005 | Classic Cars | 169 | S18 3232 | Euro
Shopping Cha... | (91) 555 94 44 | C/ Moralzarzal, 86 | Nill |
       28034 | Spain | EMEA | Freyre
                                            Diego| Large|
                      100.0|
104241
                 50 l
                                     6|12001.0|2005-05-31
00:00:00|In Process|
                2| 5| 2005|Classic Cars| 214| S10_1949|Euro
Shopping Cha... | (91) 555 94 44 | C/ Moralzarzal, 86 | Nill |
                               Freyre
       28034 l
                     EMEA |
Nill
              Spain
                                            Diegol
__+____
  ____+_____
______
```

j) Filter out rows where the STATUS is 'Cancelled' and calculate the total sales from the remaining orders?

k) Use Spark Data Frame transformations to derive the yearly sales for each customer (CUS-TOMERNAME) based on the ORDERDATE column?

```
CUSTOMERNAME | YEAR |
                                 TOTAL SALES!
+----+
| Baane Mini Imports | 2003 | 56176.659999999999 |
|Stylish Desk Deco...|2004|13739.90000000001|
|Marseille Mini Autos|2003|52481.840000000004|
|Danish Wholesale ...|2004|
                                   60157.62
|Toms Spezialitten...|2003|
                                   31363.181
|Australian Collec...|2004|140859.56999999998|
|Dragon Souveniers...|2004|
                                    3127.88
    Super Scale Inc. |2003|
                                     42498.76
|Collectables For ...|2004|15110.80000000001|
|Royal Canadian Co...|2004| 74634.84999999999|
|Online Diecast Cr...|2003|
                                   76114.7
   Gifts4AllAges.com | 2005 |
                                     48316.89
        Herkku Gifts 2004
                                      16363.1
|Diecast Classics ...|2004|115971.3400000001|
|Motor Mint Distri...|2003|
                                   27398.82
|Daedalus Designs ...|2003|48874.280000000006|
|Stylish Desk Deco...|2003|
                                   75064.6
```

----+

l) Add a new column to the DataFrame that categorizes orders as "High", "Medium", or "Low" sales based on the SALES value?

```
[71]: | quantiles = sales_df1.approxQuantile('SALES', [0.33,0.67],0.0001)
    val1 = quantiles[0]
    val2 = quantiles[1]
    sales_df1 = sales_df1.withColumn('CATEGORY', when(col('SALES') <= val1, 'Low').</pre>
     ⇔when(col('SALES') <= val2,'Medium')</pre>
                           .otherwise('High'))
    sales_df1.show()
   +-----
   __+____
   |ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES|
   ORDERDATE | STATUS | QTR_ID | MONTH_ID | YEAR_ID | PRODUCTLINE | MSRP | PRODUCTCODE |
   CUSTOMERNAME
                    PHONE
                              ADDRESSLINE1 | ADDRESSLINE2 |
                                                       CITY
   STATE | POSTALCODE |
   COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTFIRSTNAME | DEALSIZE | YEAR | CATEGORY |
   +-----
   __+____
   ----+
        10107
                      301
                           95.71
                                         2 | 2871.0 | 2003-02-24
   00:00:00|Shipped|
                   1|
                         21
                             2003|Motorcycles| 95|
                                              S10 1678
                                                       Land of
   Toys Inc.
               2125557818|897 Long Airport ...|
                                         Nill
                                                    NYCl
         100221
                 USAI
                         NAI
   NYI
                                    Yul
                                               Kwail
   Small|2003| Medium|
        10121
                      341
                          81.35 l
                                         5| 2765.9|2003-05-07
   00:00:00|Shipped|
                             2003|Motorcycles| 95|
                                              S10 1678|
                   2|
                         5|
                                                      Reims
   Collectables
                 26.47.1555| 59 rue de l'Abbaye|
                                             Nill
                                                      Reims
   Nill
          51100 l
                 Francel
                         EMEA I
                                  Henriot|
                                                Paull
   Small|2003| Medium|
        10134
                      41 l
                          94.74
                                         2|3884.34|2003-07-01
   00:00:00|Shipped|
                             2003|Motorcycles| 95|
                   3|
                         7|
                                              S10_1678|
                                                        Lyon
   Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... |
                                          Nill
                                                   Paris
                                  Da Cunhal
   Nill
          755081
                 France
                         EMEA |
                                               Daniel
```

```
Medium | 2003 | Medium |
                     45| 83.26| 6| 3746.7|2003-08-25
     10145|
00:00:00|Shipped| 3| 8| 2003|Motorcycles| 95| S10_1678|
Toys4GrownUps.com | 6265557265 | 78934 Hillside Dr. | Nill |
Pasadena | CA | 90003 | USA | NA | Young |
Julie | Medium | 2003 | Medium |
                                     14|5205.27|2003-10-10
     10159 l
                  49| 100.0|
00:00:00|Shipped| 4| 10| 2003|Motorcycles| 95| S10_1678|Corporate
Gift Id...| 6505551386| 7734 Strong St.| Nill|San Francisco|
CAL Mill IISAL Brown| Julie|
                                Brown|
    Nill
                 USAl
                          NAI
Medium | 2003 | High |
      10168|
                36| 96.66|
                                              1|3479.76|2003-10-28
00:00:00|Shipped| 4| 10| 2003|Motorcycles| 95| S10_1678|Technics
               6505556809| 9408 Furth Circle| Nill| Burlingame|
Stores Inc. |
                                  Hirano|
CAI
       94217|
                 USAI
                      NAI
                                                      Juri
Medium | 2003 | Medium |
      10180|
                     29| 86.13|
                                     9|2497.77|2003-11-11
00:00:00|Shipped| 4| 11| 2003|Motorcycles| 95| S10_1678|Daedalus
Designs ... | 20.16.1555 | 184, chausse de T... | Nill | Lille |
Nill 59000 France EMEA Rance
                                               Martinel
Small|2003|
           Lowl
      10188
                  48 100.0
                                              1|5512.32|2003-11-18
00:00:00|Shipped| 4| 11| 2003|Motorcycles| 95| S10_1678|
Herkku Gifts | +47 2267 3215 | Drammen 121, PR 7... | Nill | Bergen |
Nill N 5804 | Norway | EMEA | Oeztan | Veysel |
Medium | 2003 | High |
                                     2|2168.54|2003-12-01
      10201
                     22| 98.57|
00:00:00|Shipped| 4| 12| 2003|Motorcycles| 95| S10_1678|
                                                                Mini
Wheels Co.| 6505555787|5557 North Pendal...| Nill|San Francisco|
CAI
     Nill|
              USAl
                          NA I
                                Murphy|
                                                     Julie
             Low
Small|2003|
                  41 100.0
      10211
                                            14 | 4708 . 44 | 2004 - 01 - 15
00:00:00|Shipped| 1| 1| 2004|Motorcycles| 95| S10_1678| Auto
Canal Petit | (1) 47.55.6555 | 25, rue Lauriston | Nill |
                                                            Parisl
        75016 | France | EMEA | Perrier | Dominique |
Medium | 2004 | High |
                     37| 100.0|
      10223|
                                             1|3965.66|2004-02-20
00:00:00|Shipped| 1| 2| 2004|Motorcycles| 95| S10_1678|Australian
Collec... | 03 9520 4555 | 636 St Kilda Road | Level 3 |
Melbourne|Victoria|
                     3004|Australia| APAC| Ferguson|
Peter | Medium | 2004 | Medium |
                  23| 100.0|
                                    7 | 2333 . 12 | 2004 - 04 - 05
     10237
00:00:00|Shipped| 2| 4| 2004|Motorcycles| 95| S10_1678| Vitachrome Inc.| 2125551500| 2678 Kingston Rd.| Suite 101|
NYCl
     NYI
              10022| USA| NA|
                                             Frick | Michael |
Small|2004| Low|
                     28 | 100.0
      10251
                                             2|3188.64|2004-05-18
00:00:00|Shipped| 2| 5| 2004|Motorcycles| 95| S10_1678|Tekni
```

```
Collectable...
               2015559350|
                              7476 Moss Rd.
                                                 Nill
                                                           Newarkl
NJI
                USA|
                                    Brown
                                                 Williaml
      940191
                         NA
Medium | 2004 | Medium |
     10263
                     34|
                           100.01
                                            2|3676.76|2004-06-28
                                                   S10 1678|
                  2|
                              2004|Motorcycles| 95|
00:00:00|Shipped|
                         6|
                                                              Gift
Depot Inc.
              2035552570 | 25593 South Bay Ln. |
                                               Nill | Bridgewater |
CTI
      97562
                USA|
                         NA|
                                     King
                                                   Juliel
Medium 2004 | Medium |
     10275
                     45 l
                           92.83
                                            1 | 4177.35 | 2004-07-23
                              2004|Motorcycles| 95|
00:00:00|Shipped|
                  3|
                         7|
                                                   S10 1678
                                                             La
                 40.67.8555|67, rue des Cinqu...|
Rochelle Gifts
                                                 Nill|
                                                           Nantes
Nill
        440001
                         EMEA
                                    Labrune
               France
                                                    Janine
Medium | 2004 |
             High|
                     36 L
                           100.0
                                            6|4099.68|2004-08-27
     10285
00:00:00|Shipped|
                  3|
                         81
                              2004|Motorcycles| 95|
                                                   S10_1678|Marta's
Replicas Co.
               6175558555 | 39323 Spinnaker Dr.
                                                 Nill | Cambridge |
MAI
      51247
                USAI
                         NA
                                 Hernandez
                                                   Martal
Medium 2004
             High|
     102991
                     23|
                           100.0
                                            9|2597.39|2004-09-30
00:00:00|Shipped|
                  3|
                         9|
                              2004|Motorcycles| 95|
                                                   S10 1678|Toys of
Finland, Co.
                               Keskuskatu 45|
                                                         Helsinkil
               90-224 8555|
                                                 Nill
Nill
        21240 | Finland |
                         EMEA |
                                   Karttunen
                                                    Mattil
Small|2004| Medium|
                                            5|4394.38|2004-10-15
     10309 l
                     41 l
                           100.01
00:00:00|Shipped|
                4|
                         10|
                              2004|Motorcycles| 95| S10_1678| Baane
               07-98 9555|Erling Skakkes ga...|
                                               Nill|
                                                         Stavern
Mini Imports
                                                    Jonasl
Nill
        4110
               Norway|
                         EMEA |
                                 Bergulfsen|
Medium | 2004 |
             High|
                           94.74|
     10318
                     461
                                            1 | 4358.04 | 2004-11-02
00:00:00|Shipped|
                4|
                         11|
                              2004|Motorcycles| 95|
                                                   S10 1678|Diecast
                          7586 Pompton St.|
Classics ...|
              21555515551
                                               Nilll
                                                       Allentown
PAI
      70267
                USAI
                         NA
                                       Yul
                                                   Kyung|
Medium 2004
             High|
     10329
                     42|
                           100.0|
                                            1 | 4396.14 | 2004-11-15
                  4|
                              2004|Motorcycles| 95|
                                                   S10 1678|
00:00:00|Shipped|
                         11|
                                                             Land of
             2125557818|897 Long Airport ...|
Toys Inc.
                                             Nill
                                                         NYC
NYI
      10022
                USA|
                         NA|
                                       Yu|
                                                    Kwai|
Medium | 2004 |
             High
__+____
only showing top 20 rows
```

m) Assume, If you have another DataFrame with customer demographic data, how would you perform a join to compute the total sales per demographic group?

```
[101]: customer_demographics_data = [
           ("Land of Toys Inc.", "25-34", "High", "VIP"),
           ("Reims Collectables", "35-44", "Medium", "Regular"),
           ("Lyon Souveniers", "45-54", "Medium", "Regular"),
           ("Toys4GrownUps.com", "25-34", "High", "VIP"),
           ("Corporate Gift Ideas", "35-44", "High", "VIP"),
           ("Technics Stores Inc.", "45-54", "Medium", "Regular"),
           ("Daedalus Designs Inc.", "55-64", "Low", "New"),
           ("Herkku Gifts", "65+", "Norway", "New"),
           ("Mini Wheels Co.", "25-34", "Medium", "Regular"),
           ("Auto Canal Petit", "35-44", "Medium", "Regular"),
       customer_demographics_columns = ["CUSTOMERNAME", "AGE_GROUP", "INCOME_LEVEL", __

¬"LOYALTY STATUS"]

       customer demographics df = spark.createDataFrame(customer demographics data,,,
        →schema=customer_demographics_columns)
       customer_demographics_df.show()
```

```
CUSTOMERNAME | AGE GROUP | INCOME LEVEL | LOYALTY STATUS |
  Land of Tovs Inc.
                         25-34|
                                                       VIP
                                       High|
 Reims Collectables
                         35-44|
                                     Medium
                                                   Regular|
     Lyon Souveniers
                         45-54|
                                     Medium
                                                   Regular
   Toys4GrownUps.com
                         25-34
                                       High|
                                                       VIP
|Corporate Gift Ideas|
                         35-44|
                                       High|
                                                       VIP
|Technics Stores Inc.|
                         45-54|
                                     Medium
                                                   Regular|
|Daedalus Designs ...|
                                     Low
                                                     Newl
                       55-64|
        Herkku Gifts|
                           65+l
                                     Norway
                                                       Newl
                         25-34|
     Mini Wheels Co.
                                     Medium|
                                                   Regular|
    Auto Canal Petit
                         35-441
                                     Medium
                                                   Regular |
```

```
total_sales_per_demographic_df.show()
```

```
______
                  CUSTOMERNAME | INCOME LEVEL | LOYALTY STATUS |
|AGE GROUP|
    35-44|
              Auto Canal Petit
                                    Medium
                                                  Regular | 93170.66000000002|
    45-54 l
              Lyon Souveniers
                                    Medium
                                                  Regular
                                                                   78570.34
    35-44 | Reims Collectables |
                                    Medium|
                                                  Regular
                                                                 135042.94
    25-34 Toys4GrownUps.com
                                      High|
                                                      VIP | 104561.95999999998 |
              Mini Wheels Co.
    25-34|
                                    Medium
                                                  Regular
                                                                   74476.18
                  Herkku Gifts
      65+l
                                    Norway
                                                      Newl
                                                                   111640.28
    45-54|Technics Stores Inc.|
                                    Medium
                                                  Regular | 120783.06999999999 |
           Land of Toys Inc.
                                      High|
                                                      VIP | 164069.43999999994 |
    25-34 l
```

n) Can you implement a cumulative distribution function (CDF) over the SALES value for each CUSTOMERNAME? What insights can you gather from analyzing the CDF distribution for each customer?

```
CUSTOMERNAME | SALES |
                                      CDF |
+----+
|Suominen Souveniers| 1086.6| 0.0666666666666667|
|Suominen Souveniers|1103.76|
|Suominen Souveniers| 1988.4| 0.166666666666666666
|Suominen Souveniers|2140.11|
|Suominen Souveniers|2447.76| 0.23333333333333334|
|Suominen Souveniers|2632.89| 0.266666666666666666
|Suominen Souveniers| 2773.8|
|Suominen Souveniers|2775.08|
                         0.333333333333333333
|Suominen Souveniers|2817.87| 0.3666666666666664|
|Suominen Souveniers|2851.84|
|Suominen Souveniers | 2931.98 | 0.433333333333333335 |
|Suominen Souveniers|3128.65| 0.46666666666667|
|Suominen Souveniers|3288.82|
                                      0.51
```

```
|Suominen Souveniers|3595.62|
                                0.53333333333333333
|Suominen Souveniers|3686.54|
                                0.5666666666666671
|Suominen Souveniers| 3784.8|
                                                0.61
|Suominen Souveniers| 4068.7|
                                0.63333333333333333
|Suominen Souveniers | 4142.64 |
                                0.6666666666666666
|Suominen Souveniers|4157.73|
                                                0.71
|Suominen Souveniers | 4381.25 |
                                0.73333333333333333
|Suominen Souveniers| 4836.5|
                                0.7666666666666671
|Suominen Souveniers|5154.41|
                                                0.8
|Suominen Souveniers|5500.44|
                                0.833333333333334
|Suominen Souveniers|5938.53|
                                0.866666666666667
|Suominen Souveniers|6287.66|
                                                0.91
|Suominen Souveniers| 6576.5|
                                0.9333333333333333
|Suominen Souveniers| 6756.0|
                                0.9666666666666671
|Suominen Souveniers|10606.2|
| Amica Models & Co.|
                       577.6 | 0.038461538461538464 |
| Amica Models & Co.|1381.05| 0.07692307692307693|
| Amica Models & Co. | 1557.36 | 0.11538461538461539 |
| Amica Models & Co. | 1574.0 | 0.15384615384615385 |
| Amica Models & Co. | 1656.69 | 0.19230769230769232 |
| Amica Models & Co. | 1921.92 | 0.23076923076923078 |
| Amica Models & Co. | 2084.81 |
                                0.2692307692307692
| Amica Models & Co. | 2137.05 |
                                0.3076923076923077
| Amica Models & Co. | 2418.24 | 0.34615384615384615 |
| Amica Models & Co. 2800.08 | 0.38461538461538464 |
| Amica Models & Co. | 2819.28 |
                                0.4230769230769231
| Amica Models & Co. | 2941.89 |
                               0.46153846153846156
| Amica Models & Co. | 2954.53 |
                                                0.51
| Amica Models & Co. | 3006.43 |
                                0.5384615384615384
| Amica Models & Co. | 3474.46 |
                                0.5769230769230769
| Amica Models & Co. | 3668.6|
                                0.6153846153846154
| Amica Models & Co.|3704.05|
                                0.6538461538461539
| Amica Models & Co. | 4242.24 |
                                0.69230769230769231
| Amica Models & Co. | 4455.0|
                                0.7307692307692307
| Amica Models & Co. | 4750.8|
                                0.7692307692307693
+----
```

only showing top 50 rows

#### Insights

- CDF values show how sales are distributed within each customer group.
- Sales Concentration for a customer like "Suominen Souveniers," we can see that a significant portion of the total sales is concentrated at higher values.
- For example, 0.7 CDF corresponds to a sales value of 4381.25, meaning that 70% of the sales data falls below this amount.
- o) Write spark dataframe code to rank products by total revenue within each country (COUN-

## TRY)?

+	+			+
PR	ODUCTCODE	COUNTRY	Total revenue	RANK
+	+		·	+
1	S18_4600	Sweden	4900.0	1
1	S18_4600	Sweden	4800.0	2
1	S24_2300	Sweden	4800.0	2
1	S12_4675	Sweden	4700.0	4
	S18_2949	Sweden	4700.0	4
	S18_2319	Sweden	4600.0	6
1	S24_1578	Sweden	4500.0	7
1	S10_4757	Sweden	4400.0	8
1	S18_4522	Sweden	4300.5	9
1	S18_1662	Sweden	4300.0	10
1	S24_2300	Sweden	4200.0	11
1	S12_1666	Sweden	4100.0	12
1	S18_1097	Sweden	4100.0	12
1	S24_2011	Sweden	4100.0	12
1	S24_2000	Sweden	3988.6000000000004	15
1	S18_2625	Sweden	3900.0	16
1	S18_1889	Sweden	3881.77999999999997	17
1	S10_1949	Sweden	3700.0	18
1	S24_3151	Sweden	3519.85	19
1	S12_3380	Sweden	3500.0	20
+	+		·	+
onl	y showing	top 20 1	rows	

p) Calculate a running total of SALES for each customer and show the top 5 customers by this cumulative total?

```
[93]: window_spec = Window.partitionBy("CUSTOMERNAME").orderBy("ORDERDATE").

→rowsBetween(Window.unboundedPreceding, Window.currentRow)
```

q) Identify and handle Outliers in DataFrame.

```
--+----
|ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES|
ORDERDATE
          STATUS|QTR_ID|MONTH_ID|YEAR_ID| PRODUCTLINE|MSRP|PRODUCTCODE|
                           ADDRESSLINE1 | ADDRESSLINE2 |
CUSTOMERNAME!
                 PHONE |
STATE | POSTALCODE |
COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTFIRSTNAME | DEALSIZE | YEAR |
__+____
______
___+____
10150|
                    45|
                         100.0
                                        8|10993.5|2003-09-19
00:00:00
         Shipped
                   3|
                          91
                             2003 | Classic Cars | 214 |
                                                  S10_1949|Dragon
Souveniers... | +65 221 7555 | Bronz Sok., Bronz... |
                                         Nill
                                                 Singapore |
       79903|Singapore|
                       Japan|
                                                 Eric
                                Natividad|
Large|2003|
                    34 l
     10174
                         100.01
                                         4|8014.82|2003-11-06
00:00:001
       Shipped
                   4|
                         11| 2003|Classic Cars| 214|
S10_1949|Australian Gift N...|61-7-3844-6555|31 Duncan St. Wes...|
Nill|South Brisbane|Queensland| 4101|Australia|
                                           APAC
                                                     Calaghan|
     Large | 2003 |
Tony|
102061
                    47 l
                         100.01
                                        6|9064.89|2003-12-05
                   4|
00:00:00
         Shipped
                         12|
                             2003 | Classic Cars | 214 |
S10_1949|Canadian Gift Exc...|(604) 555-3392|
                                     1900 Oak St.|
                                                       Nill
                  V3F 2K1 | Canada|
Vancouver
             BCI
                                     NA I
                                            Tannamuri |
Yoshi
      Large | 2003 |
     10280
                    34|
                         100.0
                                         2|8014.82|2004-08-17
00:00:00
         Shipped
                   3|
                         8| 2004|Classic Cars| 214|
                                                  S10 1949|
Amica Models & Co.|
                011-4988555| Via Monte Bianco 34|
                                                Nill
                                           Accortil
Torino|
         Nill
                 10100|
                         Italy|
                                 EMEA |
Paolo| Large|2004|
                                        6 | 10172.7 | 2004-10-11
     10304
                    47|
                         100.0
00:00:00
         Shipped
                   4|
                         10|
                              2004 | Classic Cars | 214 |
                                                  S10 1949 | Auto
             30.59.8555|67, avenue de 1'E...|
Assoc. & Cie.
                                          Nill|
                                                  Versailles|
Nill
       78000 l
              Francel
                       EMEA I
                                  Tonini
                                               Daniell
Large | 2004 |
                                        3|11623.7|2004-10-21
     10312
                    48 l
                         100.01
4|
                         10|
                             2004|Classic Cars| 214|
                                                  S10_1949|Mini
00:00:00
         Shipped
             4155551450
                        5677 Strong St.
Gifts Distri...
                                            Nill|
                                                   San Rafael
CAI
      97562|
               USAI
                                 Nelson
                                             Valarie
                       NA|
Large | 2004 |
                   36|
                                        3 | 8254.8 | 2005-02-17
     10381
                         100.0
00:00:00
       Shipped| 1|
                         21
                            2005|Classic Cars| 214|
S10_1949|Corporate Gift Id...| 6505551386|
                                    7734 Strong St.|
                                                       Nill
San Francisco
                CAI
                      Nill| USA|
                                        NAI
                                                 Brownl
```

```
Julie | Large | 2005 |
                     50| 100.0|
     10424|
                                              6|12001.0|2005-05-31
1
00:00:00|In Process| 2| 5| 2005|Classic Cars| 214| S10_1949|Euro
Shopping Cha...|(91) 555 94 44| C/ Moralzarzal, 86| Nill|
Nill| 28034| Spain| EMEA| Freyre| Diego|
Large | 2005 |
                     46| 100.0|
      10120|
                                               2|9264.86|2003-04-29
00:00:00| Shipped| 2| 4| 2003| Motorcycles| 193|
S10_4698|Australian Collec...| 03 9520 4555| 636 St Kilda Road|
Melbourne | Victoria | 3004 | Australia | APAC |
                                                    Ferguson
Peter | Large | 2003 |
      10180|
                     41 100.0|
                                              11 | 8892.9 | 2003-11-11
Nill| 59000| France| EMEA|
Lille
                                                   Rancel
Martine | Large | 2003 |
      10188|
                      45| 100.0|
                                       3 | 8714.7 | 2003-11-18
00:00:00| Shipped | 4| 11| 2003| Motorcycles | 193| S10_4698|
Herkku Gifts | +47 2267 3215 | Drammen 121, PR 7... | Nill | Bergen |
      N 5804 | Norway | EMEA | Oeztan | Veysel |
Nill
Large | 2003 |
| 10201| 49| 100.0| 4|8065.89|2003-12-01
| 00:00:00| Shipped| 4| 12| 2003| Motorcycles| 193| S10_4698|
Mini Wheels Co. | 6505555787 | 5557 North Pendal ... | Nill | San Francisco |
CA | Nill | USA | NA | Murphy | Julie |
Large | 2003 |
| 10223| 49| 100.0| 3|9774.03|

00:00:00| Shipped| 1| 2| 2004| Motorcycles| 193|
                                               3|9774.03|2004-02-20
S10_4698|Australian Collec...| 03 9520 4555| 636 St Kilda Road| Level 3|
Melbourne | Victoria | 3004 | Australia | APAC |
                                                    Ferguson
Peter | Large | 2004 |
     10263|
                      41 100.0
                                               4|8336.94|2004-06-28
00:00:00| Shipped| 2| 6| 2004| Motorcycles| 193| S10_4698|
Gift Depot Inc. | 2035552570 | 25593 South Bay Ln. | Nill | Bridgewater | CT | 97562 | USA | NA | King | Julie |
Large | 2004 |
                     66| 100.0|
      10403|
                                               9|11886.6|2005-04-08
00:00:00| Shipped| 2|
                              4| 2005| Motorcycles| 193| S10_4698|UK
Collectables, ...|(171) 555-2282|Berkeley Gardens ...| Nill|
Liverpool | Nill | WX1 6LT | UK | EMEA |
                                                       Devonl
Elizabeth| Large|2005|
                                     4|9218.16|2005-05-13
                      56| 100.0|
10417
00:00:00| Disputed | 2 | 5 | 2005 | Motorcycles | 193 | S10_4698 | Euro
Shopping Cha... (91) 555 94 44 | C/ Moralzarzal, 86 | Nill | Madrid |
        28034| Spain| EMEA| Freyre| Diego|
Nill
Large | 2005 |
| 10400| 64| 100.0| 9|9661.44|2005-04-01
| 00:00:00| Shipped| 2| 4| 2005|Classic Cars| 136| S10_4757|The
```

```
Sharp Gifts W...
            4085553659
                       3086 Ingle Ln.
                                      Nill
                                             San
Josel
        CAI
             94217|
                     USAI
                            NA
                                     Frick
                                                 Suel
Large | 2005 |
                42|
                     100.0
                                  7|8008.56|2003-07-02
    10135
00:00:001
       Shipped
                31
                      71
                         2003 | Classic Cars | 194 |
                                          S12 1099|Mini
Gifts Distri...
                      5677 Strong St.
                                     Nill
                                           San Rafaell
           4155551450
     975621
            USAI
                    NAI
                            Nelson
                                      Valariel
Large | 2003 |
                 48 l
                     100.01
                                  719245.7612003-09-05
    10147
                         2003 | Classic Cars | 194 |
00:00:00
       Shipped
                3|
                      91
S12_1099|Collectables For ...|
                    6175558555
                              7825 Douglas Av.|
                                              Nill|
Brickhaven
            MA
                 58339
                         USA
                                NA I
                                        Nelson
Allen
     Large | 2003 |
                41 l
                                  2|8296.35|2003-10-10
    10159
                     100.01
00:00:00
       Shipped
                4|
                     10|
                         2003 | Classic Cars | 194 |
S12_1099|Corporate Gift Id...|
                    6505551386
                               7734 Strong St.
San Franciscol
              CAL
                    Nill
                           USAI
                                  NAI
                                          Brown
Julie
     Large | 2003 |
+-----
__+____
______
___________
-+
only showing top 20 rows
__+____
______
|ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES|
ORDERDATE | STATUS | QTR ID | MONTH ID | YEAR ID | PRODUCTLINE | MSRP | PRODUCTCODE |
                        ADDRESSLINE1 | ADDRESSLINE2 |
CUSTOMERNAME
               PHONE
                                               CITY
STATE | POSTAL CODE |
COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTFIRSTNAME | DEALSIZE | YEAR |
+-----
10107|
                                  2 | 2871.0 | 2003-02-24
30 l
                      95.71
                    2|
00:00:00|Shipped|
              1|
                       2003|Motorcycles| 95|
                                       S10 1678
                                               Land of
          2125557818|897 Long Airport ...|
                                            NYC|
Toys Inc.
                                  Nill|
     10022
NY
            USA
                    NA I
                              Yul
                                        Kwai|
Small|2003|
                34 l
                     81.35 l
                                  5| 2765.9|2003-05-07
    10121
00:00:00|Shipped|
              2|
                    5|
                       2003|Motorcycles| 95|
                                       S10_1678| Reims
Collectables|
            26.47.1555| 59 rue de l'Abbaye|
                                      Nill
                                              Reims
Nill
      51100 l
            France
                    EMEA |
                            Henriot|
                                         Paul
Small|2003|
```

```
| 10134| 41| 94.74| 2|3884.34|2003-07-01
| 00:00:00|Shipped| 3| 7| 2003|Motorcycles| 95| S10_1678| Lyon
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... | Nill |
                                                                  Parisl
         75508| France|
                            EMEA I
                                        Da Cunha|
                                                            Daniell
Medium | 2003 |
       10145|
                        45| 83.26|
                                                  6| 3746.7|2003-08-25
00:00:00|Shipped| 3| 8| 2003|Motorcycles| 95| S10_1678|
Toys4GrownUps.com| 6265557265| 78934 Hillside Dr.| Nill|
Pasadena| CA| 90003| USA| NA| Young|
                     90003| USA| NA|
Juliel Medium 2003
                     49| 100.0|
      10159|
                                         14|5205.27|2003-10-10
00:00:00|Shipped| 4| 10| 2003|Motorcycles| 95| S10_1678|Corporate
Gift Id... | 6505551386 | 7734 Strong St. | Nill | San Francisco | CA | Nill | USA | NA | Brown | Julie |
Medium | 2003 |
               36| 96.66|
                                                  1|3479.76|2003-10-28
      101681
00:00:00|Shipped| 4| 10| 2003|Motorcycles| 95| S10_1678|Technics
Stores Inc. | 6505556809 | 9408 Furth Circle | Nill | Burlingame | CAL 94217 | USAL WAL Hirang
CAL
       942171
                 USAl
                         NA| Hirano|
                                                             Juri|
Medium | 2003 |
                                                   9|2497.77|2003-11-11
                     29| 86.13|
      10180|
00:00:00|Shipped| 4| 11| 2003|Motorcycles| 95| S10_1678|Daedalus
Designs ... | 20.16.1555 | 184, chausse de T... | Nill | Lille |
Nill 59000 France EMEA
                                      Rance
                                                     Martinel
Small|2003|
| 10188| 48| 100.0| 1|5512.32|2003-11-
00:00:00|Shipped| 4| 11| 2003|Motorcycles| 95| S10_1678|
                                                   1|5512.32|2003-11-18

      Herkku Gifts| +47 2267 3215|Drammen 121, PR 7...|
      Nill| N 5804| Norway| EMEA| Oeztan|
      Nill| Veysel|

Medium | 2003 |
                    22 | 98.57 | 2 | 2168.54 | 2003-12-01
      10201
00:00:00|Shipped| 4| 12| 2003|Motorcycles| 95| S10_1678|
Wheels Co.| 6505555787|5557 North Pendal...| Nill|San Francisco|
CAI
        Nill USA | NA | Murphy |
                                                          Juliel
Small|2003|
| 10211| 41| 100.0| 14|4708.44|2004-01-15
| 00:00:00|Shipped| 1| 1| 2004|Motorcycles| 95| S10_1678| Auto
Canal Petit | (1) 47.55.6555 | 25, rue Lauriston | Nill | Paris |
Nill| 75016| France| EMEA| Perrier| Dominique|
Medium | 2004 |
                        37| 100.0|
                                                  1|3965.66|2004-02-20
      102231
00:00:00|Shipped| 1| 2| 2004|Motorcycles| 95| S10_1678|Australian
          03 9520 4555| 636 St Kilda Road|
Collec...
                                                Level 3|
Melbourne|Victoria|
                       3004|Australia| APAC| Ferguson|
Peter | Medium | 2004 |
                                                   7|2333.12|2004-04-05
      10237
                        23 100.0
00:00:00|Shipped| 2| 4| 2004|Motorcycles| 95| S10_167
Vitachrome Inc.| 2125551500| 2678 Kingston Rd.| Suite 101|
                              4| 2004|Motorcycles| 95| S10_1678|
```

```
NYC| NY| 10022| USA| NA| Frick| Michael|
Small|2004|
                   28 | 100.0 |
                                      2|3188.64|2004-05-18
     10251
00:00:00|Shipped|
              2|
                       5|
                          2004|Motorcycles| 95| S10_1678|Tekni
                          7476 Moss Rd.| Nill|
Collectable...
              2015559350|
                                                    Newarkl
NJ|
     940191
              USA| NA|
                                Brownl
                                           Williaml
Medium | 2004 |
                   34 100.01
                                       2|3676.76|2004-06-28
     102631
00:00:00|Shipped|
               2|
                      6| 2004|Motorcycles| 95| S10_1678| Gift
Depot Inc.
            2035552570 | 25593 South Bay Ln. | Nill | Bridgewater |
CTI
     97562|
              USAl
                      NA I
                          King|
                                             Julie
Medium | 2004 |
                  45|
     10275
                       92.83|
                                       1|4177.35|2004-07-23
                      7| 2004|Motorcycles| 95| S10_1678|
00:00:00|Shipped|
               3|
               40.67.8555|67, rue des Cinqu...| Nill|
Rochelle Gifts
                                                    Nantes
Nilll
       44000 l
             France
                      EMEA
                           Labrune
                                              Janinel
Medium | 2004 |
                   361
                       100.0
                                       6|4099.68|2004-08-27
10285
00:00:00|Shipped| 3|
                       8|
                          2004|Motorcycles| 95| S10_1678|Marta's
             6175558555| 39323 Spinnaker Dr.| Nill| Cambridge|
Replicas Co.
MAI
     51247|
                      NAI
                             Hernandez
              USAl
                                            Martal
Medium 2004
                   231
                       100.01
                                      9|2597.39|2004-09-30
     102991
                      9| 2004|Motorcycles| 95|
00:00:00|Shipped| 3|
                                             S10 1678|Toys of
Finland, Co.
             90-224 8555 | Keskuskatu 45
                                           Nill
                                                   Helsinkil
Nill
       21240 | Finland | EMEA |
                               Karttunen|
                                              Matti
Small|2004|
                   41 100.0
                                5|4394.38|2004-10-15
    10309|
00:00:00|Shipped| 4|
                      10| 2004|Motorcycles| 95| S10_1678| Baane
Mini Imports | 07-98 9555 | Erling Skakkes ga... | Nill | Stavern |
             Norway| EMEA| Bergulfsen|
Nilll
        4110|
                                              Jonasl
Medium | 2004 |
     10318|
                   46|
                       94.74
                                       1|4358.04|2004-11-02
00:00:00|Shipped| 4|
                    11| 2004|Motorcycles| 95| S10_1678|Diecast
Classics ...|
            2155551555| 7586 Pompton St.| Nill| Allentown|
                                  Yu|
PAI
              USAl
                      NAI
     702671
                                             Kyung|
Medium | 2004 |
     103291
                  42| 100.0|
                                       1 | 4396 . 14 | 2004 - 11 - 15
00:00:00|Shipped| 4|
                          2004|Motorcycles| 95|
                      11|
                                             S10_1678| Land of
           2125557818|897 Long Airport ...|
                                        Nill
                                                   NYCl
Toys Inc. |
                      NAI
NYI
     10022
              USAl
                                  Yul
                                              Kwail
Medium 2004
__+____
______
only showing top 20 rows
```

r) How would you cache a DataFrame containing sales data from the top 10 countries by sales to avoid recomputation in subsequent transformations? What persistence level (e.g. MEMORY ONLY, MEMORY AND DISK) would you choose and why?

```
+----+
| COUNTRY| TOTAL_SALES|
+-----+
| USA| 3627982.83|
| Spain|1215686.9200000009|
| France|1110916.5199999993|
|Australia| 630623.1000000001|
| UK| 478880.4600000001|
+-----+
only showing top 5 rows
```

[106]: DataFrame[COUNTRY: string, TOTAL\_SALES: double]

**Recommended Persistence Level** For caching the DataFrame of sales data from the top 10 countries, I would recommend using MEMORY AND DISK. This choice is suitable because:

Data Size: The DataFrame might be relatively small after filtering to the top 10 countries, so Resilience: If the DataFrame size grows or if you face memory constraints, it will spill to difference Balance: Provides a good balance between performance (memory access) and reliability

s) How would you pivot the data to show PRODUCTLINE as columns and the total SALES for each ORDERDATE as the values? What are the implications of pivoting large datasets in Spark?

```
[141]: pivot_df = sales_df1.groupBy("ORDERDATE", "PRODUCTLINE").agg(sum("SALES").
    ⇔alias("TOTAL_SALES"))
   pivot_df.groupBy("ORDERDATE").pivot("PRODUCTLINE").sum("TOTAL_SALES").show()
   +----+
   -----+
         ORDERDATE| Classic Cars|
                             Motorcycles|
   Ships | Trains | Trucks and Buses | Vintage Cars |
   -----
   |2005-03-02 00:00:00|
                     null|
                                4175.6
                                            null
            null|
   null| null|
                       null
```

2004-11-09 00:00:00				null
6673.29 3807.68				
2005-05-03 00:00:00	25040.6299	99999997	null	null
null  null  :	27247.11	null		
2003-09-11 00:00:00	43593.540	00000001	null	null
null  null	null	3598.22		
2005-04-01 00:00:00   6284.0  null		9661.44	null	9036.06
6284.0  null	nul	12545.3	4	
2005-05-10 00:00:00		null	7567.8	30429.010000000002
null  null				
2004-04-26 00:00:00  null  null		null	null	null
null  null	null	7129.0		
2003-10-17 00:00:00	40321.609	99999999	null	null
null  null	nu111	nu111		
2004-09-27 00:00:00		null 530	7.9800000000005	null
null  null	null	null	7.9800000000005	
2003-10-20 00:00:00		4860.24	null	null
null  null		20424.51		
2004-04-09 00:00:00		31329.56	null	null
null  null				
2005-05-06 00:00:00				
null 122664 6000000000	27/20/0 /	o I	null 16070	വ
2003-02-17 00:00:00   6598.34  null		null	null	39205.310000000005
6598.34  null	nu.	10377.6	67	
2004-12-07 00:00:00		19489.57 1039	94.560000000001	null
null  null	null	null		null
2005-04-22 00:00:00		40207.85	null	null
null  null	null	18229.19		
2003-12-01 00:00:00			31.879999999997	7120.96
null  null		1113.6		
2003-01-06 00:00:00		null	null	null
null  null	null	12133.25		
2004-05-18 00:00:00			27987.07	null
null  null	null	null		
2005-01-05 00:00:00		null		null
null  null	null	null		
2004-11-17 00:00:00			86.9200000000021	null
		37747.89	·	·
	murr ,	01111001		
++			+	+

only showing top 20 rows

## Implications of Pivoting Large Datasets

\* Memory Usage: Pivoting reshapes data, creating many new columns for each unique pivot value, which can lead to high memory consumption.

- \* Performance: The operation can be slow with large datasets or many unique pivot values due to increased computational complexity.
- \* Column Explosion: A large number of unique pivot values can result in a DataFrame with many making it unwieldy and hard to manage.
- \* Data Skew: Uneven distribution of data (e.g., some dates having much more sales) can cause performance issues due to uneven partitioning.
- \* Shuffling: Pivoting involves shuffling data across nodes, which can be expensive in terms of time and resources.
- \* Disk I/O: If the data exceeds memory capacity, Spark will spill data to disk, potentially slowing down the process.
  - t) How would you calculate the percentage growth of total sales month over month for each PRODUCTLINE using Spark DataFrame?

```
[110]: from pyspark.sql.functions import *
       from pyspark.sql.window import Window
       salesData = sales_df1.withColumn("year", year("ORDERDATE"))
       salesData = salesData.withColumn("month", month("ORDERDATE"))
       monthlySales = salesData.groupBy("year", "month", "PRODUCTLINE").
        →agg(sum("SALES").alias("total_sales"))
       windowSpec = Window.partitionBy("PRODUCTLINE").orderBy("year", "month")
       monthlySales = monthlySales.withColumn("previous_month_sales", __
        →lag("total_sales").over(windowSpec))
       monthlySales = monthlySales.withColumn(
           "percentage_growth",
           when(
               col("previous_month_sales").isNotNull(),
               (col("total_sales") - col("previous_month_sales")) /_

¬col("previous_month_sales") * 100
           ).otherwise(lit(None))
      monthlySales.show()
```

percentage_growt	h		+
+	+-	+	+
2003  2 Moto	rcycles 2	25783.760000000002	null
null	•		
2003  3 Moto	rcycles	12639.15	25783.760000000002
-50.980190631622	•		
2003  4 Moto	rcycles 2	23475.590000000004	12639.15
85.7370946622202	-		
2003  5 Moto	rcycles	22097.32	
23475.5900000000	04 -5.871	10771486467594	
2003  6 Moto	rcycles	2642.01	22097.32
-88.043753722170	82		
2003  7 Moto	rcycles	37924.23000000001	2642.01
1335.43097868668			
2003  8 Moto	rcycles 4	44164.909999999996	37924.23000000001
16.4556538128789	•		
	•	3155.58	44164.909999999996
-92.855006384027			
	-	64235.65000000001	3155.58
1935.62102687936			
	•	109345.5	64235.65000000001
70.2255678894819			
		25431.879999999997	109345.5
-76.741722338825			
2004  1 Moto		41200.52	25431.879999999997
62.0034382043325			
2004  2 Moto	-	49066.5	41200.52
19.0919434997422		0.0000 070000000001	40000 = 1
	•	36269.07000000001	49066.5
-26.081807343095		44040 0500000041	0.0000 0.000000001
	. •	16848.950000000004	36269.07000000001
29.1705301514485		47007 441	46040 05000000004
	rcycles	47237.41	46848.950000000004
0.82917546711292		00774 01	
2004  7 Moto 47237.41 -51.788	•	22774.0	
	rcycles	62704.93	22774.0
175.335602002283	•	02704.931	22114.01
		42471.04999999999	62704.93
-32.268403776226	•	T2T11.UT333333333	02104.931
02.200400110220			

42471.0499999999|-7.1980560876173065| +----+

only showing top 20 rows

|2004| 10|Motorcycles| 39413.96|

u) How can you rebalance the data by portioning based on the COUNTRY column to ensure that large data partitions are avoided?

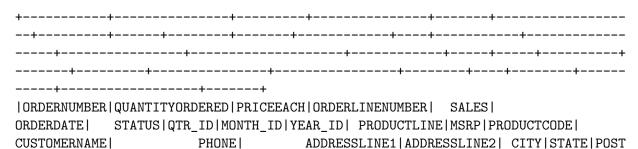
```
[126]: unique_country_count = sales_df1.select("COUNTRY").distinct().count()

repartitioned_df = sales_df1.repartition(unique_country_count, "COUNTRY")

print(f"No of partitions: {repartitioned_df.rdd.getNumPartitions()}")
```

No of partitions: 19

v) Suppose you have a smaller lookup table with customer details. How would you perform a broadcast join with the large sales\_data\_sample dataset to improve join performance? What are the key considerations when using broadcast joins?



```
ALCODE | COUNTRY | TERRITORY | CONTACTLASTNAME | CONTACTFIRSTNAME | DEALSIZE | YEAR | CATEGORY
|customer_id| name|country|
__+____
______
                  30|
                         95.7l
10107|
                                        2 | 2871.0 | 2003-02-24
00:00:00| Shipped| 1| 2| 2003| Motorcycles| 95|
                                                 S10 1678 | Land
of Toys Inc. | 2125557818 | 897 Long Airport ... | Nill | NYC | NY |
10022| USA| NA|
                         Yu| Kwai|
                                              Small|2003| Medium|
1 | Land of Toys Inc. | USA |
                  34| 81.35|
                                        5| 2765.9|2003-05-07
     10121|
00:00:00| Shipped | 2 | 5 | 2003 | Motorcycles | 95 | S10_1678 | Reims
Collectables | 26.47.1555 | 59 rue de l'Abbaye | Nill | Reims | Nill |
51100| Francel
             EMEA| Henriot| Paul| Small|2003| Medium|
2|Reims Collectables| Canada|
                                        2|3884.34|2003-07-01
     10134|
                        94.74
00:00:00| Shipped| 3| 7| 2003| Motorcycles| 95| S10_1678| Lyon
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... | Nill | Paris | Nill |
                      Da Cunha | Daniel | Medium | 2003 | Medium |
             EMEAI
75508| Francel
3 | Lyon Souveniers | Mexico |
    103291
                        100.01
                                        1|4396.14|2004-11-15
00:00:00| Shipped | 4 | 11 | 2004 | Motorcycles | 95 |
                                                 S10 1678 | Land
of Toys Inc. | 2125557818 | 897 Long Airport ... | Nill | NYC | NY |
10022| USA|
                        Yu| Kwai| Medium|2004| High|
               NA |
1 | Land of Toys Inc. | USA |
                   39|
                                        5|3896.49|2003-02-24
     10107|
                        99.91
00:00:00| Shipped| 1| 2| 2003| Motorcycles| 118|
                                                 S10_2016 | Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
       USAl
                      Yu|
               NAI
                                      Kwai| Medium|2003| Medium|
1 | Land of Toys Inc. | USA |
                   27| 100.0|
10134|
                                        5|3307.77|2003-07-01
00:00:00| Shipped| 3| 7| 2003| Motorcycles| 118| S10_2016|
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel ... | Nill | Paris | Nill |
                       Da Cunha| Daniel| Medium|2003| Medium|
75508| France|
              EMEA I
3 | Lyon Souveniers | Mexico |
                                        2 | 3176.0 | 2004-11-15
     103291
                   20 | 100.0 |
00:00:00| Shipped | 4 | 11 | 2004 | Motorcycles | 118 | S10_2016 | Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
                           Yu| Kwai| Medium|2004| Medium|
10022| USA|
             NA |
1 | Land of Toys Inc. | USA |
| 10107| 27| 100.0|
| 00:00:00| Shipped| 1| 2| 200
                                        4|6065.55|2003-02-24
                         2 | 2003 | Motorcycles | 193 | S10_4698 | Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
10022| USA|
           NAI
                     Yu| Kwai| Medium|2003| High|
1 | Land of Toys Inc. | USA |
| 10134|
                   31 100.0
                                       4|7023.98|2003-07-01
```

```
00:00:00| Shipped | 3 | 7 | 2003 | Motorcycles | 193 | S10_4698 | Lyon
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... | Nill | Paris | Nill |
                      Da Cunha| Daniel| Large|2003| High|
              EMEA I
75508| Francel
3 | Lyon Souveniers | Mexico |
              26| 100.0| 3| 5868.2|2004-11-15
     103291
00:00:00| Shipped | 4 | 11 | 2004 | Motorcycles | 193 | S10_4698 | Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
                              Yu| Kwai| Medium|2004| High|
10022| USA|
               NAI
1 | Land of Toys Inc. | USA |
     102481
                   20 | 100.0 |
                                           3 | 2910.4 | 2004-05-07
00:00:00|Cancelled| 2| 5| 2004|Classic Cars| 136| S10_4757| Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
                        Yu| Kwai|
10022| USA| NA|
                                                  Small|2004| Medium|
1 | Land of Toys Inc. | USA |
                    48| 54.68|
      10359|
                                            6|2624.64|2004-12-15
00:00:00 | Shipped | 4 | 12 | 2004 | Classic Cars | 136 | S10_4757 | Reims
Collectables | 26.47.1555 | 59 rue de l'Abbaye | Nill | Reims | Nill |
               EMEA| Henriot| Paul| Small|2004| Medium|
51100| France|
2|Reims Collectables| Canada|
                   32| 100.0|
     10395 l
                                           2|3370.56|2005-03-17
00:00:00| Shipped | 1 | 3 | 2005|Classic Cars | 136|
                                                     S10 4757| Lyon
Souveniers +33 1 46 62 7555 27 rue du Colonel... | Nill | Paris | Nill |
75508 | France | EMEA | Da Cunha | Daniel | Medium | 2005 | Medium |
3 | Lyon Souveniers | Mexico |
| 10329| 41| 71.47|
                                           5|2930.27|2004-11-15
00:00:00| Shipped| 4| 11| 2004|Classic Cars| 194| S12_1099| Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
10022| USA|
                        Yu| Kwai| Small|2004| Medium|
                NA
1 | Land of Toys Inc. | USA |
             42| 100.0|
     10359|
                                           8|4764.48|2004-12-15
00:00:00| Shipped | 4 | 12 | 2004 | Classic Cars | 207 | S12_1108 | Reims
Collectables | 26.47.1555 | 59 rue de l'Abbaye | Nill | Reims | Nill |
               EMEA| Henriot| Paul| Medium|2004|
51100| France|
                                                              High|
2|Reims Collectables| Canada|
     10395
                   33 69.12
                                           1|2280.96|2005-03-17
00:00:00| Shipped| 1| 3|
                               2005|Classic Cars| 207| S12 1108|
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... | Nill | Paris | Nill |
75508| Francel
              EMEA I
                         Da Cunha| Daniel| Small|2005|
3| Lyon Souveniers | Mexico |
| 10107| 21| 100.0| 1| 3036.6|2003-02-24
| 00:00:00| Shipped| 1| 2| 2003| Motorcycles| 150| S12_2823| Land
of Toys Inc.| 2125557818|897 Long Airport ...| Nill| NYC| NY|
                          Yu| Kwai| Medium|2003| Medium|
      USA|
10022|
                NA
1 | Land of Toys Inc. | USA |
                   50 | 100.0 |
                                           4| 8284.0|2003-05-07
     10121
00:00:00| Shipped| 2| 5| 2003| Motorcycles| 150| S12_2823|Reims
Collectables | 26.47.1555 | 59 rue de l'Abbaye | Nill | Reims | Nill |
51100 | France | EMEA | Henriot | Paul | Large | 2003 | High |
```

```
2|Reims Collectables| Canada|
                                   1 | 2711.2 | 2003-07-01
    10134 l
                 201
                      100.0
00:00:00| Shipped|
                3|
                      7 I
                         2003 | Motorcycles | 150 |
                                           S12 2823|
                                                   Lyon
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... |
                                     Nill|Paris| Nill|
75508| Francel
                    Da Cunhal
                                        Small|2003| Medium|
            EMEA
                                 Daniell
   Lyon Souveniers | Mexico |
    10329
                      100.0
                                   6|3542.64|2004-11-15
00:00:00| Shipped|
                41
                     111
                         2004 | Motorcycles | 150 |
                                           S12 2823 | Land
of Toys Inc.
            2125557818|897 Long Airport ...|
                                      Nill | NYC|
                                  Kwail Medium | 2004 | Medium |
100221
      USAl
              NA|
                        Yııl
1 | Land of Toys Inc. |
                USAI
__+____
_____
______
----+
only showing top 20 rows
```

w) Create a UDF that categorizes the sales values (SALES) into custom buckets like "Low", "Medium", "High". Apply this UDF to the DataFrame and calculate the count of orders in each category per COUNTRY.

```
+-----+
| COUNTRY|SALES_CATEGORY|count|
+-----+
|Philippines| Low| 7|
| Norway| Medium| 22|
| USA| Low| 319|
| Austria| Low| 17|
```

```
Denmark
                        Lowl
                                18|
      Canada
                       High|
                                18|
      Canada
                     Medium
                                24|
      Canada
                        Low
                                28|
                       High|
       Italy|
                                321
|Switzerland|
                     Medium|
                                12|
     Ireland|
                       High|
                                8|
     Ireland
                        Low
                                 7|
|Philippines|
                       High|
                                 9|
  Australia
                       High|
                                57|
   Singapore
                                28|
                        Low
     Finland|
                                34|
                     Medium|
      Norway
                       High|
                                33|
         USA
                     Medium
                               338|
       Italy|
                        Low
                                38|
         USAI
                       High|
                              347
only showing top 20 rows
```

x) Create a Python UDF to calculate discounts for specific product lines. For example, give a 10% discount for Classic Cars and 5% for Motorcycles. Apply this UDF to derive new discounted sales values.

```
[124]: from pyspark.sql.functions import udf, col
from pyspark.sql.types import FloatType

def calculate_discounted_sales(product_line, sales):
    if product_line == 'Classic Cars':
        discount = 0.10
    elif product_line == 'Motorcycles':
        discount = 0.05
    else:
        discount = 0.0
    return float(sales * (1 - discount))

calculate_discounted_sales_udf = udf(calculate_discounted_sales, FloatType())

sales_df2 = sales_df2.withColumn(
    'DISCOUNT_SALES',
    calculate_discounted_sales_udf(col('PRODUCTLINE'), col('SALES'))
)

sales_df2.select(col('PRODUCTLINE'), col('SALES'), col('DISCOUNT_SALES')).show()
```

```
|PRODUCTLINE| SALES|DISCOUNT_SALES|
+----+
|Motorcycles| 2871.0|
                           2727.45
|Motorcycles| 2765.9|
                          2627.605
|Motorcycles|3884.34|
                          3690.123|
|Motorcycles| 3746.7|
                           3559.365
|Motorcycles|5205.27|
                          4945.0063|
|Motorcycles|3479.76|
                           3305.772
|Motorcycles|2497.77|
                          2372.8816
|Motorcycles|5512.32|
                           5236.704|
|Motorcycles|2168.54|
                           2060.113|
|Motorcycles|4708.44|
                           4473.018
|Motorcycles|3965.66|
                           3767.377
|Motorcycles|2333.12|
                           2216.464|
|Motorcycles|3188.64|
                           3029.2081
|Motorcycles|3676.76|
                           3492.922|
|Motorcycles|4177.35|
                          3968.4824
|Motorcycles|4099.68|
                           3894.696|
|Motorcycles|2597.39|
                          2467.5205
|Motorcycles|4394.38|
                           4174.661
|Motorcycles|4358.04|
                          4140.138|
|Motorcycles|4396.14|
                          4176.333|
only showing top 20 rows
```

z) How do you implement a cumulative distribution function (CDF) over the SALES value for each CUSTOMERNAME? What insights can you gather from analyzing the CDF distribution for each customer?

```
|Suominen Souveniers|2632.89| 0.266666666666666666
|Suominen Souveniers| 2773.8|
                                                0.31
|Suominen Souveniers|2775.08|
                                0.333333333333333333
|Suominen Souveniers|2817.87| 0.366666666666664|
|Suominen Souveniers | 2851.84 |
                                                0.41
|Suominen Souveniers | 2931.98 |
                               0.4333333333333335|
|Suominen Souveniers|3128.65|
                                0.466666666666667
|Suominen Souveniers|3288.82|
                                                0.51
|Suominen Souveniers|3595.62|
                                0.533333333333333333
|Suominen Souveniers|3686.54|
                                0.5666666666666671
|Suominen Souveniers| 3784.8|
                                                0.6
|Suominen Souveniers| 4068.7|
                                0.63333333333333333
|Suominen Souveniers | 4142.64 |
                                0.6666666666666666
|Suominen Souveniers | 4157.73 |
                                                0.71
|Suominen Souveniers | 4381.25 |
                                0.733333333333333333
|Suominen Souveniers| 4836.5|
                                0.7666666666666671
|Suominen Souveniers|5154.41|
                                                0.81
|Suominen Souveniers|5500.44|
                                0.8333333333333341
|Suominen Souveniers|5938.53|
                                0.866666666666667|
|Suominen Souveniers | 6287.66|
                                                0.91
|Suominen Souveniers | 6576.5|
                                0.93333333333333333
|Suominen Souveniers | 6756.0|
                                0.966666666666667
|Suominen Souveniers | 10606.2| |
| Amica Models & Co.| 577.6|0.038461538461538464|
| Amica Models & Co. | 1381.05 | 0.07692307692307693 |
| Amica Models & Co. | 1557.36 | 0.11538461538461539 |
| Amica Models & Co. | 1574.0 | 0.15384615384615385 |
| Amica Models & Co. | 1656.69 | 0.19230769230769232 |
| Amica Models & Co. | 1921.92 | 0.23076923076923078 |
| Amica Models & Co. | 2084.81 |
                                0.26923076923076921
| Amica Models & Co. | 2137.05 |
                                0.3076923076923077
| Amica Models & Co. | 2418.24 | 0.34615384615384615 |
| Amica Models & Co. | 2800.08 | 0.38461538461538464 |
| Amica Models & Co. | 2819.28 |
                                0.4230769230769231
| Amica Models & Co. | 2941.89 | 0.46153846153846156 |
| Amica Models & Co. | 2954.53 |
                                                0.51
| Amica Models & Co. | 3006.43 |
                                0.5384615384615384
| Amica Models & Co. | 3474.46 |
                                0.57692307692307691
| Amica Models & Co. | 3668.6|
                                0.6153846153846154
| Amica Models & Co. | 3704.05 |
                                0.6538461538461539
| Amica Models & Co. | 4242.24 |
                                0.6923076923076923
| Amica Models & Co. | 4455.0|
                                0.7307692307692307
| Amica Models & Co. | 4750.8|
                                0.7692307692307693
+----+
```

only showing top 50 rows

Insights

- \* CDF values show how sales are distributed within each customer group.
- \* Sales Concentration for a customer like "Suominen Souveniers," we can see that a significant
- \* For example, 0.7 CDF corresponds to a sales value of 4381.25, meaning that 70% of the sales

## [95]: sales df1.printSchema()

## root

- |-- ORDERNUMBER: integer (nullable = true)
- |-- QUANTITYORDERED: integer (nullable = true)
- |-- PRICEEACH: double (nullable = true)
- |-- ORDERLINENUMBER: integer (nullable = true)
- |-- SALES: double (nullable = true)
- |-- ORDERDATE: timestamp (nullable = true)
- |-- STATUS: string (nullable = false)
- |-- QTR\_ID: integer (nullable = true)
- |-- MONTH\_ID: integer (nullable = true)
- |-- YEAR\_ID: integer (nullable = true)
- |-- PRODUCTLINE: string (nullable = false)
- |-- MSRP: integer (nullable = true)
- |-- PRODUCTCODE: string (nullable = false)
- |-- CUSTOMERNAME: string (nullable = false)
- |-- PHONE: string (nullable = false)
- |-- ADDRESSLINE1: string (nullable = false)
- |-- ADDRESSLINE2: string (nullable = false)
- |-- CITY: string (nullable = false)
- |-- STATE: string (nullable = false)
- |-- POSTALCODE: string (nullable = false)
- |-- COUNTRY: string (nullable = false)
- |-- TERRITORY: string (nullable = false)
- |-- CONTACTLASTNAME: string (nullable = false)
- |-- CONTACTFIRSTNAME: string (nullable = false)
- |-- DEALSIZE: string (nullable = false)
- |-- YEAR: integer (nullable = true)
- |-- CATEGORY: string (nullable = false)

#### [96]: sales\_df1.show(10)

|ORDERNUMBER|QUANTITYORDERED|PRICEEACH|ORDERLINENUMBER| SALES|
ORDERDATE| STATUS|QTR\_ID|MONTH\_ID|YEAR\_ID|PRODUCTLINE|MSRP|PRODUCTCODE|
CUSTOMERNAME| PHONE| ADDRESSLINE1|ADDRESSLINE2| CITY|ST
ATE|POSTALCODE|COUNTRY|TERRITORY|CONTACTLASTNAME|CONTACTFIRSTNAME|DEALSIZE|YEAR|
CATEGORY|

```
+-----
__+____
______
                     30 l
                           95.7
                                           2 | 2871.0 | 2003-02-24
     101071
00:00:00|Shipped| 1|
                         2|
                             2003|Motorcycles| 95|
                                                 S10_1678| Land of
             2125557818|897 Long Airport ...|
                                           Nill|
                                                        NYCl
      10022| USA|
                    NA |
                                    Yu|
                                                Kwai| Small|2003|
Mediuml
                     34|
                                           5| 2765.9|2003-05-07
     10121
                         81.35|
00:00:00|Shipped| 2|
                       5| 2003|Motorcycles| 95| S10_1678| Reims
Collectables | 26.47.1555 | 59 rue de l'Abbaye | Nill |
                                                          Reims
        51100 | France | EMEA |
                                  Henriot
                                                  Paull
Small|2003| Medium|
     101341
                     41 l
                          94.741
                                           2|3884.34|2003-07-01
00:00:00|Shipped| 3|
                        7| 2003|Motorcycles| 95| S10_1678|
                                                             Lyon
Souveniers | +33 1 46 62 7555 | 27 rue du Colonel... |
                                            Nill
                                                       Paris
        75508| France|
                                 Da Cunha|
                      EMEA I
                                                Daniel|
Medium | 2003 | Medium |
     10145 l
                     45 83.26
                                           6 | 3746.7 | 2003 - 08 - 25
00:00:00|Shipped| 3| 8| 2003|Motorcycles| 95|
Toys4GrownUps.com| 6265557265| 78934 Hillside Dr.|
                             2003|Motorcycles| 95| S10_1678|
               90003| USA|
Pasadena
        CAI
                                 NAI
                                            Young|
                                                         Juliel
Medium | 2003 | Medium |
                     49|
                          100.0|
                                          14|5205.27|2003-10-10
     10159|
                             2003|Motorcycles| 95| S10_1678|Corporate
00:00:00|Shipped| 4|
                       10|
            6505551386|
                       7734 Strong St.|
Gift Id...|
                                        Nill|San Francisco|
CAI
              USAl
                       NA I
                                  Brownl
                                                Julie | Medium | 2003 |
       Nill
High|
     10168|
                     36|
                         96.66
                                          1|3479.76|2003-10-28
00:00:00|Shipped|
                4 | 10 | 2003 | Motorcycles | 95 | S10_1678 | Technics
              6505556809| 9408 Furth Circle|
Stores Inc. |
                                              Nill | Burlingame |
CAI
      94217| USA|
                    NA |
                                Hirano|
                                                Juri | Medium | 2003 |
Medium |
     10180|
                     29 | 86.13 |
                                           9|2497.77|2003-11-11
00:00:00|Shipped| 4| 11| 2003|Motorcycles| 95| S10_1678|Daedalus
             20.16.1555|184, chausse de T...|
                                          Nill|
Nill
        59000| France|
                    EMEA |
                                  Rance
Small|2003|
             Low |
                                           1|5512.32|2003-11-18
10188 l
                     48|
                         100.0
00:00:00|Shipped| 4|
                        11| 2003|Motorcycles| 95| S10_1678|
Herkku Gifts
           +47 2267 3215|Drammen 121, PR 7...|
                                              Nill|
                                                        Bergen
                       EMEA |
                                  Oeztan
       N 5804| Norway|
                                                Vevsel
Medium | 2003 |
             High|
     10201
                     22|
                          98.57
                                           2|2168.54|2003-12-01
00:00:00|Shipped|
                             2003|Motorcycles| 95| S10_1678|
                 4|
                        12|
                                                             Mini
Wheels Co. | 6505555787 | 5557 North Pendal... | Nill | San Francisco |
```

CAI USA| Murphy| NillNA| Julie| Small|2003| Lowl 41| 10211| 100.0| 14|4708.44|2004-01-15 00:00:00|Shipped| 1| 1| 2004|Motorcycles| 95| S10\_1678| Auto Canal Petit (1) 47.55.6555 | 25, rue Lauriston Nill| Paris| Nill| 75016| France| EMEA | Perrier| Dominique Medium | 2004 | High| +-----\_\_+\_\_\_\_ \_\_\_\_\_ only showing top 10 rows

[]: