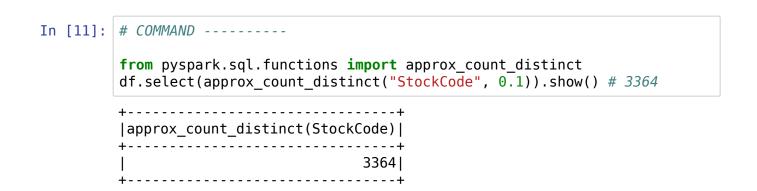
```
In [2]:
       from pyspark.sql import SparkSession
       import pyspark.sql.functions as F
       from pyspark.sql.types import *
       spark = SparkSession\
          .builder\
          .appName("chapter-07-aggregation")\
          .get0rCreate()
       import os
       SPARK BOOK DATA PATH = os.environ['SPARK BOOK DATA PATH']
In [3]:
      file path = SPARK BOOK DATA PATH + "/data/retail-data/all/*.csv"
       df = spark.read.format("csv")\
         .option("header", "true")\
        .option("inferSchema", "true")\
        .load(file path)\
        .coalesce(3)
       df.show(5,False)
       +-----
       -----
       |InvoiceNo|StockCode|Description
                                                     |Quantity|Invo
             |UnitPrice|CustomerID|Country
       iceDate
       +-----
       1536365
             |85123A
                       |WHITE HANGING HEART T-LIGHT HOLDER | 6
                                                            112/
       1/2010 8:26|2.55
                       |17850
                                |United Kingdom|
                       |WHITE METAL LANTERN
              |71053
                                                            |12/
       1536365
                                                     |6
       1/2010 8:26|3.39
                       |17850
                                |United Kingdom|
       |536365
                       |CREAM CUPID HEARTS COAT HANGER
              184406B
                                                     18
                                                            |12/
       1/2010 8:26|2.75
                       |17850
                                |United Kingdom|
       |536365
                       |KNITTED UNION FLAG HOT WATER BOTTLE|6
             184029G
                                                            112/
                              |United Kingdom|
       1/2010 8:26|3.39
                       |17850
       |536365
              184029E
                       |RED WOOLLY HOTTIE WHITE HEART.
                                                     16
                                                            |12/
       1/2010 8:26|3.39 |17850
                              |United Kingdom|
       -----+
       only showing top 5 rows
In [4]: df.cache()
Out[4]: DataFrame[InvoiceNo: string, StockCode: string, Description: string, Q
      uantity: int, InvoiceDate: string, UnitPrice: double, CustomerID: int,
      Country: string]
In [5]:
      df.createOrReplaceTempView("dfTable")
```

```
In [6]: # COMMAND -----
        from pyspark.sql.functions import count
        df.select(count("StockCode")).show() # 541909
        +----+
        |count(StockCode)|
                  541909|
In [7]: # COMMAND -----
        from pyspark.sql.functions import countDistinct,expr,desc
        df.select(countDistinct("StockCode")).show() # 4070
In [8]:
        |count(DISTINCT StockCode)|
                             4070|
       df.groupBy("StockCode").agg(expr("count(StockCode)")).show(10, False)
In [9]:
        +----+
        |StockCode|count(StockCode)|
        122728
                 1810
        21889
                  1607
        190210B
                  17
                  1296
        21259
        121894
                 1135
        21452
                  |200
        22121
                  |141
        190022
                  |21
        21249
                  1119
        |90143
                 |22
        only showing top 10 rows
```

```
df.where("StockCode in ('90026D', '90210B') ").orderBy("StockCode", des
+-----
-----+
|InvoiceNo|StockCode|Description
                                            |Quantity|Invoic
        |UnitPrice|CustomerID|Country
-----+
                |GLASS BEAD HOOP NECKLACE AMETHYST|1
1548545
        190026D
                                                   13/31/2
011 19:12 |8.5
                         |United Kingdom|
                |13118
|544463
        |90026D
                |GLASS BEAD HOOP NECKLACE AMETHYST|1
                                                   |2/20/2
011 14:31 |8.5
                12988
                         |United Kingdom|
|577315
        190026D
                |GLASS BEAD HOOP NECKLACE AMETHYST|1
                                                   111/18/
                         |United Kingdom|
2011 13:25 | 8.5
                17811
                |CLEAR ACRYLIC FACETED BANGLE
|568787
        190210B
                                           |6
                                                   19/29/2
011 9:19
        12.95
                         |United Kingdom|
                13741
                |CLEAR ACRYLIC FACETED BANGLE
1538071
        |90210B
                                           |1
                                                   |12/9/2
        12.96
010 14:09
                Inull
                         |United Kingdom|
                |CLEAR ACRYLIC FACETED BANGLE
1581434
        190210B
                                           |10
                                                   12/8/2
011 16:10
                         |United Kingdom|
        11.0
                13599
|537045
        190210B
                |CLEAR ACRYLIC FACETED BANGLE
                                           |1
                                                   |12/5/2
010 10:54 | 2.95
                115038
                         |United Kingdom|
|538661
        190210B
                |CLEAR ACRYLIC FACETED BANGLE
                                           |12
                                                   |12/13/
                         |United Kingdom|
2010 15:42 | 1.25
                15194
                |CLEAR ACRYLIC FACETED BANGLE
1536409
        190210B
                                           |1
                                                   12/1/2
                         |United Kingdom|
010 11:45 | 2.95
                17908
|573102
       190210B
                |CLEAR ACRYLIC FACETED BANGLE
                                           |10
                                                   |10/27/
2011 14:52|1.0
                13266
                       |United Kingdom|
-----
```



```
In [12]: | # COMMAND -----
       from pyspark.sql.functions import first, last
       df.select(first("StockCode"), last("StockCode")).show()
       +----+
       |first(StockCode, false)|last(StockCode, false)|
                     85123A|
        ----+
In [13]:
       # COMMAND -----
       from pyspark.sql.functions import min, max
       df.select(min("Quantity"), max("Quantity")).show()
       +----+
       |min(Quantity)|max(Quantity)|
       +----+
            -80995| 80995|
       +----+
In [14]: # COMMAND -----
       from pyspark.sql.functions import sum
       df.select(sum("Quantity")).show() # 5176450
       |sum(Quantity)|
       +----+
            5176450
In [15]:
       # COMMAND -----
       from pyspark.sql.functions import sumDistinct
       df.select(sumDistinct("Quantity")).show() # 29310
       +-----+
       |sum(DISTINCT Quantity)|
                     293101
       +----+
```

```
In [16]: | # COMMAND -----
      from pyspark.sql.functions import sum, count, avg, expr
      df.select(
         count("Quantity").alias("total transactions"),
         sum("Quantity").alias("total_purchases"),
         avg("Quantity").alias("avg purchases"),
         expr("mean(Quantity)").alias("mean purchases"))\
        .selectExpr(
         "total purchases/total transactions",
         "avg_purchases",
         "mean purchases").show()
      +-----
      |(total purchases / total transactions)| avg purchases| mean purcha
      9.55224954743324|9.55224954743324|9.55224954743
      3241
      In [17]: # COMMAND -----
      from pyspark.sql.functions import var pop, stddev pop
      from pyspark.sql.functions import var_samp, stddev_samp
      df.select(var pop("Quantity"),
             var samp("Quantity"),
             stddev pop("Quantity"),
             stddev_samp("Quantity")).show()
      +-----
      |var_pop(Quantity)|var_samp(Quantity)|stddev_pop(Quantity)|stddev_samp
      (Quantity)|
      |47559.30364660928| 47559.39140929898| 218.08095663447847|
                                                  218.081
      1578502347|
      +-----
      ----+
```

```
In [18]: | # COMMAND -----
       from pyspark.sql.functions import skewness, kurtosis
       df.select(skewness("Quantity"), kurtosis("Quantity")).show()
       +----+
         skewness(Quantity)|kurtosis(Quantity)|
       +----+
       |-0.26407557610528154| 119768.054955306|
       +----+
In [19]:
       # COMMAND -----
       from pyspark.sql.functions import corr, covar pop, covar samp
       df.select(
             corr("InvoiceNo", "Quantity"),
             covar_samp("InvoiceNo", "Quantity"),
             covar_pop("InvoiceNo", "Quantity")).show()
       |corr(InvoiceNo, Quantity)|covar samp(InvoiceNo, Quantity)|covar pop(I
       nvoiceNo, Quantity)|
       +-----
       ----+
           4.912186085639875E-4|
                                      1052.7280543916152
       1052.726077875511
       +-----
       -----+
In [22]: # COMMAND -----
       from pyspark.sql.functions import collect set, collect list
       df.agg(collect set("Country"), collect list("Country")).show(10, False)
       IOPub data rate exceeded.
       The notebook server will temporarily stop sending output
       to the client in order to avoid crashing it.
       To change this limit, set the config variable
       `--NotebookApp.iopub data rate limit`.
       Current values:
       NotebookApp.iopub data rate limit=1000000.0 (bytes/sec)
       NotebookApp.rate limit window=3.0 (secs)
```

```
|InvoiceNo|quan|count(Quantity)|
    536596|
                                  14|
    5369381
               14|
    537252|
                                   1|
               1|
                                  20|
    537691|
               20|
    538041|
                                   1|
                1|
                                  26|
    538184
               26|
    538517|
               53|
                                  53|
                                  191
    5388791
               19|
    539275|
                6|
                                   6|
    539630|
               12|
                                  12|
    540499|
               241
                                  24|
    540540|
               22|
                                  22|
   C540850|
                                   1|
                1|
    5409761
               481
                                  481
    541432|
                4|
                                   4|
    541518 | 101 |
                                 101|
    541783|
               35|
                                  35|
                                   9|
    542026|
                91
                                   61
    542375|
                61
   C542604|
                8|
                                   8|
```

only showing top 20 rows

```
+----+
            avg(Quantity)|stddev pop(Quantity)|
 -----+
                            1.1180339887498947
   536596 L
                       1.5
   536938|33.142857142857146|
                            20.698023172885524
                      31.0|
                             5.597097462078001
   5376911
                      8.15|
                      30.01
   538041
                                          0.01
   538184 | 12.076923076923077 |
                             8.142590198943392|
   538517 | 3.0377358490566038 |
                            2.3946659604837897
   538879 | 21.157894736842106 |
                            11.811070444356483
   5392751
                      26.01
                             12.806248474865697
   539630 | 20.3333333333333332 |
                            10.225241100118645
   540499|
                            2.6653642652865788
                      3.75
   540540 | 2.1363636363636362 |
                            1.0572457590557278
  C5408501
                      -1.0|
                                          0.01
   540976 | 10.520833333333334 |
                             6.496760677872902
   5414321
                     12.25|
                            10.825317547305483
   541518 | 23.10891089108911 |
                            20.550782784878713|
   541783 | 11.314285714285715 |
                             8.467657556242811
   542026 | 7.666666666666667 |
                             4.853406592853679
   5423751
                      8.01
                            3.4641016151377544
  C5426041
                     -8.01
                            15.173990905493518
 -----+
```

only showing top 20 rows

```
In [25]: # COMMAND -----
```

```
from pyspark.sql.functions import col, to_date
dfWithDate = df.withColumn("date", to_date(col("InvoiceDate"), "MM/dd/y
dfWithDate.createOrReplaceTempView("dfWithDate")
```

```
dfWithDate.show(10)
In [26]:
       ----+
        |InvoiceNo|StockCode|
                               Description|Quantity| InvoiceDate|Unit
       Price|CustomerID|
                                       datel
                           Country
       +----
        ----+------
                   85123A|WHITE HANGING HEA...|
           536365|
                                                 6|12/1/2010 8:26|
                17850|United Kingdom|2010-12-01|
       2.55
           536365|
                    71053| WHITE METAL LANTERN|
                                                 6|12/1/2010 8:26|
                17850|United Kingdom|2010-12-01|
       3.391
                   84406B|CREAM CUPID HEART...|
           5363651
                                                 8|12/1/2010 8:26|
                17850|United Kingdom|2010-12-01|
       2.75
                   84029G|KNITTED UNION FLA...|
                                                 6|12/1/2010 8:26|
           536365|
       3.39|
                17850|United Kingdom|2010-12-01|
                   84029E|RED WOOLLY HOTTIE...|
           5363651
                                                 6|12/1/2010 8:26|
                17850|United Kingdom|2010-12-01|
       3.39
                    22752|SET 7 BABUSHKA NE...|
                                                 2|12/1/2010 8:26|
           5363651
                17850|United Kingdom|2010-12-01|
       7.65
           536365|
                    21730 | GLASS STAR FROSTE...|
                                                 6|12/1/2010 8:26|
                17850|United Kingdom|2010-12-01|
       4.251
           5363661
                    22633 | HAND WARMER UNION...|
                                                 6|12/1/2010 8:28|
       1.85|
                17850|United Kingdom|2010-12-01|
                    22632 | HAND WARMER RED P...|
                                                 6|12/1/2010 8:28|
           5363661
                17850|United Kingdom|2010-12-01|
        1.85|
           536367|
                    84879|ASSORTED COLOUR B...|
                                                32|12/1/2010 8:34|
                13047|United Kingdom|2010-12-01|
        1.69|
       ----+
       only showing top 10 rows
In [27]: # COMMAND -----
        from pyspark.sql.window import Window
        from pyspark.sql.functions import desc
       windowSpec = Window\
         .partitionBy("CustomerId", "date")\
         .orderBy(desc("Quantity"))\
         .rowsBetween(Window.unboundedPreceding, Window.currentRow)
In [28]:
       # COMMAND -----
        from pyspark.sql.functions import max
        maxPurchaseQuantity = max(col("Quantity")).over(windowSpec)
In [29]: # COMMAND -----
        from pyspark.sql.functions import dense rank, rank
        purchaseDenseRank = dense rank().over(windowSpec)
        purchaseRank = rank().over(windowSpec)
```

```
In [30]: # COMMAND -----
from pyspark.sql.functions import col

dfWithDate.where("CustomerId IS NOT NULL").orderBy("CustomerId")\
    .select(
    col("CustomerId"),
    col("date"),
    col("Quantity"),
    purchaseRank.alias("quantityRank"),
    purchaseDenseRank.alias("quantityDenseRank"),
    maxPurchaseQuantity.alias("maxPurchaseQuantity")).show()
```

	+ +	+		
	omerId  date uantity	Quantity quar	ntityRank quantit	yDenseRank maxPurc
+		+		
	12346   2011 - 01 - 18	74215	1	1
1	12346 2011-01-18	-74215	2	2
74215	12347 2010-12-07	36	1	1
36  	12347   2010 - 12 - 07	30	2	2
36  	12347   2010 - 12 - 07	24	3	3
36  	12347 2010-12-07	12	4	4
36  	12347 2010-12-07	12	4	4
36  	12347 2010-12-07	12	4	4
36  	12347 2010-12-07	12	4	4
36  	12347 2010-12-07	12	4	4
36  	12347 2010-12-07	-	4	4
36  	12347 2010-12-07	-	4	4
36	·		·	•
36	12347 2010-12-07		4	4
 36	12347   2010 - 12 - 07		4	4
 36	12347   2010 - 12 - 07		4	4
 36	12347   2010 - 12 - 07	12	4	4
 36	12347   2010 - 12 - 07	12	4	4
 36	12347   2010 - 12 - 07	12	4	4
	12347   2010 - 12 - 07	[ 6	17	5

```
In [31]: # COMMAND -----

dfNoNull = dfWithDate.drop()
dfNoNull.createOrReplaceTempView("dfNoNull")
```

```
In [32]: # COMMAND ------

rolledUpDF = dfNoNull.rollup("Date", "Country").agg(sum("Quantity"))\
    .selectExpr("Date", "Country", "`sum(Quantity)` as total_quantity")\
    .orderBy("Date")
rolledUpDF.show()
```

+	<b></b>	++					
Date	Country	total_quantity					
+null		++   5176450					
2010-12-01							
2010-12-01		•					
2010-12-01							
2010-12-01							
2010-12-01		117					
2010-12-01	null	26814					
2010-12-01	United Kingdom	23949					
2010-12-01	Norway	1852					
2010-12-02	null	21023					
2010-12-02	Germany	146					
2010-12-02	EIRE	4					
2010-12-02	United Kingdom	20873					
2010-12-03	France	239					
2010-12-03	Italy	164					
2010-12-03							
2010-12-03		110					
2010-12-03							
2010-12-03							
2010-12-03	null	14830					
++							
only showing top 20 rows							

++-			+				
Date	Country	<pre>Country sum(Quantity) count(Quantity) </pre>					
++							
null	Japan	25218	358				
null	Australia	83653	1259				
null	Portugal	16180	1519				
null	null	5176450	541909				
null	RSA	352	58				
null	Finland	10666	695				
nullil	Inited Arab Emirates	982	68				
null	Singapore	5234	229				
null	Unspecified	3300	446				
null	Germany	117448	9495				
null	Channel Islands	9479	758				
null	USA	1034	291				
null	Hong Kong	4769	288				
null	Denmark	8188	389				
null	Czech Republic	592	30				
null	European Community	497	61				
null	Lebanon	386	45				
null	Spain	26824	2533				
null	Cyprus	6317	622				
null	Norway	19247	1086				
++							
only showing top 20 rows							

```
In [34]: # COMMAND ------
pivoted = dfWithDate.groupBy("date").pivot("Country").sum()
# COMMAND ------
In [41]: # pivoted.show()
```

```
In [ ]:
```