

```
In [1]: 1 from IPython.display import display, clear_output
```

```
In [2]: 1 from pyspark.sql import SparkSession
2 import pyspark.sql.functions as F
3 from pyspark.sql.types import *
4
5 spark = SparkSession\
6     .builder\
7     .appName("chapter-03-tour")\
8     .getOrCreate()
9
10 import os
11 SPARK_BOOK_DATA_PATH = os.environ['SPARK_BOOK_DATA_PATH']
```

Spark SQL

```
In [4]: 1 file_path = SPARK_BOOK_DATA_PATH + "/data/retail-data/by-day/*.csv"
2
3 retail_df = spark.read.csv(file_path, header=True, inferSchema=True)
```

```
1 retail_df = spark.read\
2     .format("csv")\
3     .option("header", "true")\
4     .option("inferSchema", "true")\
5     .load(file_path)
```

```
In [5]: 1 retail_df.count()
```

```
Out[5]: 541909
```


In [7]: 1 display(retail_df.toPandas())

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	580538	23084	RABBIT NIGHT LIGHT	48	2011-12-05 08:38:00	1.79	14075.0	United Kingdom
1	580538	23077	DOUGHNUT LIP GLOSS	20	2011-12-05 08:38:00	1.25	14075.0	United Kingdom
2	580538	22906	12 MESSAGE CARDS WITH ENVELOPES	24	2011-12-05 08:38:00	1.65	14075.0	United Kingdom
3	580538	21914	BLUE HARMONICA IN BOX	24	2011-12-05 08:38:00	1.25	14075.0	United Kingdom
4	580538	22467	GUMBALL COAT RACK	6	2011-12-05 08:38:00	2.55	14075.0	United Kingdom
...
541904	543282	22849	BREAD BIN DINER STYLE MINT	1	2011-02-06 16:08:00	16.95	12956.0	United Kingdom
541905	543282	84879	ASSORTED COLOUR BIRD ORNAMENT	8	2011-02-06 16:08:00	1.69	12956.0	United Kingdom
541906	543282	84659A	WHITE TRAVEL ALARM CLOCK	1	2011-02-06 16:08:00	2.55	12956.0	United Kingdom
541907	543282	82484	WOOD BLACK BOARD ANT WHITE FINISH	1	2011-02-06 16:08:00	7.95	12956.0	United Kingdom
541908	543282	22168	ORGANISER WOOD ANTIQUE WHITE	1	2011-02-06 16:08:00	8.50	12956.0	United Kingdom

541909 rows × 8 columns

In [8]: 1 retail_df.createOrReplaceTempView("retail_table")
2 # retail_table is a SQL table for query

In [9]: 1 staticSchema = retail_df.schema

```
In [10]: 1 print(staticSchema)
```


```
StructType(List(StructField(InvoiceNo,StringType,true),StructField(StockCode,StringType,true),StructField(Description,StringType,true),StructField(Quantity,IntegerType,true),StructField(InvoiceDate,StringType,true),StructField(UnitPrice,DoubleType,true),StructField(CustomerID,DoubleType,true),StructField(Country,StringType,true)))
```

```
In [11]: 1 retail_df.printSchema()
```

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

```
In [12]: 1 display(retail_df.describe().toPandas())
```

	summary	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID
0	count	541909	541909	540455	541909	541909	541909	40
1	mean	559965.752026781	27623.240210938104	20713.0	9.55224954743324	None	4.611113626089641	15287.69057023
2	stddev	13428.417280796697	16799.737628427683	NaN	218.0811578502335	None	96.75985306117963	1713.60030332
3	min	536365	10002	4 PURPLE FLOCK DINNER CANDLES	-80995	2010-12-01 08:26:00	-11062.06	123
4	max	C581569	m	wrongly sold sets	80995	2011-12-09 12:50:00	38970.0	182



```
In [13]: 1 df = spark.sql("select * from retail_table limit 5")
```

```
In [14]: 1 df.show(truncate=False) # disable truncate to show description in full
```

```
+-----+-----+-----+-----+-----+-----+-----+
---+-----+
|InvoiceNo|StockCode|Description|Quantity|InvoiceDate|UnitPrice|Custome
rID|Country|
+-----+-----+-----+-----+-----+-----+-----+
---+-----+
|580538|23084|RABBIT NIGHT LIGHT|48|2011-12-05 08:38:00|1.79|14075.0
|United Kingdom|
|580538|23077|DOUGHNUT LIP GLOSS|20|2011-12-05 08:38:00|1.25|14075.0
|United Kingdom|
|580538|22906|12 MESSAGE CARDS WITH ENVELOPES|24|2011-12-05 08:38:00|1.65|14075.0
|United Kingdom|
|580538|21914|BLUE HARMONICA IN BOX|24|2011-12-05 08:38:00|1.25|14075.0
|United Kingdom|
|580538|22467|GUMBALL COAT RACK|6|2011-12-05 08:38:00|2.55|14075.0
|United Kingdom|
+-----+-----+-----+-----+-----+-----+-----+
---+-----+
```

In [15]:

```
1 # COMMAND -----
2
3 from pyspark.sql.functions import window, column, desc, col
4
5 (retail_df.selectExpr(
6     "CustomerId",
7     "(UnitPrice * Quantity) as total_cost",
8     "InvoiceDate")
9     .groupBy(col("CustomerId"), window(col("InvoiceDate"), "1 day"))
10    .sum("total_cost")
11    .sort(desc("sum(total_cost)"))
12    .show(5, False)
13 )
```

```
+-----+-----+-----+-----+
|CustomerId|window                                |sum(total_cost) |
+-----+-----+-----+-----+
|17450.0   |[2011-09-19 20:00:00, 2011-09-20 20:00:00]|71601.44         |
|null      |[2011-11-13 19:00:00, 2011-11-14 19:00:00]|55316.08         |
|null      |[2011-11-06 19:00:00, 2011-11-07 19:00:00]|42939.17         |
|null      |[2011-03-28 20:00:00, 2011-03-29 20:00:00]|33521.39999999998|
|null      |[2011-12-07 19:00:00, 2011-12-08 19:00:00]|31975.590000000007|
+-----+-----+-----+-----+
```

only showing top 5 rows

In [16]:

```
1 # col() can be omitted
2 (
3 retail_df.selectExpr(
4     "CustomerId",
5     "(UnitPrice * Quantity) as total_cost",
6     "InvoiceDate")
7 .groupBy("CustomerId", window("InvoiceDate", "1 day"))
8 .sum("total_cost")
9 .withColumnRenamed("sum(total_cost)", "sum_total_cost")
10 .withColumnRenamed("window", "InvoiceDateWindow")
11 .sort(desc("sum_total_cost"))
12 .withColumn("sum_total_cost", F.round("sum_total_cost",2))
13 .show(5, truncate=False)
14 )
```

```
+-----+-----+-----+-----+
|CustomerId|InvoiceDateWindow|sum_total_cost|
+-----+-----+-----+-----+
|17450.0   |[2011-09-19 20:00:00, 2011-09-20 20:00:00]|71601.44|
|null     |[2011-11-13 19:00:00, 2011-11-14 19:00:00]|55316.08|
|null     |[2011-11-06 19:00:00, 2011-11-07 19:00:00]|42939.17|
|null     |[2011-03-28 20:00:00, 2011-03-29 20:00:00]|33521.4 |
|null     |[2011-12-07 19:00:00, 2011-12-08 19:00:00]|31975.59|
+-----+-----+-----+-----+
```

only showing top 5 rows

Spark Streaming

In [17]:

```
1 ## Extract
2 streamingDataFrame = (spark
3     .readStream
4     .format("csv")
5     .schema(staticSchema)
6     .option("maxFilesPerTrigger", 1)
7     .option("header", "true")
8     .load(SPARK_BOOK_DATA_PATH + "/data/retail-data/by-day/*.csv")
9 )
```

```
In [18]: 1  ## Transform
2  # cost per day
3  purchaseByCustomerPerDay = (streamingDataFrame
4    .selectExpr(
5      "CustomerId",
6      "(UnitPrice * Quantity) as total_cost",
7      "InvoiceDate")
8    .groupBy(col("CustomerId"), window(col("InvoiceDate"), "1 day"))
9    .sum("total_cost")
10 )
```

```
In [19]: 1  ## Load
2  # store result into a SQL table `customer_purchases` is specified by `queryName`
3  (
4    purchaseByCustomerPerDay
5      .writeStream
6      .queryName("customer_purchases")
7      .format("memory")
8      .outputMode("complete")
9      .start()
10 )
```

Out[19]: <pyspark.sql.streaming.StreamingQuery at 0x7f452f78dd90>


```
In [26]: 1  ## Monitor stream
2  # use `Ctrl-Enter` to execute below cell repeatedly to see streaming result as more data are read
3  spark.sql("""
4  SELECT
5  *
6  FROM customer_purchases
7  ORDER BY `sum(total_cost)` DESC
8  """).show(5, False)
```

```
+-----+-----+-----+
|CustomerId|window|sum(total_cost)|
+-----+-----+-----+
|null|[2011-11-06 19:00:00, 2011-11-07 19:00:00]|42939.17|
|null|[2011-07-03 20:00:00, 2011-07-04 20:00:00]|13667.659999999993|
|18102.0|[2011-07-03 20:00:00, 2011-07-04 20:00:00]|13282.0|
|null|[2011-11-21 19:00:00, 2011-11-22 19:00:00]|13216.8899999999894|
|null|[2010-11-30 19:00:00, 2010-12-01 19:00:00]|12584.2999999999988|
+-----+-----+-----+
only showing top 5 rows
```

```
In [21]: 1  spark.sql("""
2  SELECT *
3  FROM customer_purchases
4  ORDER BY `sum(total_cost)` DESC
5  """).show(5,truncate=False)
```

```
+-----+-----+-----+
|CustomerId|window|sum(total_cost)|
+-----+-----+-----+
|12678.0|[2011-10-27 20:00:00, 2011-10-28 20:00:00]|8947.9600000000005|
|null|[2011-05-18 20:00:00, 2011-05-19 20:00:00]|4012.66000000000067|
|13694.0|[2011-10-27 20:00:00, 2011-10-28 20:00:00]|3304.0300000000001|
|null|[2011-10-27 20:00:00, 2011-10-28 20:00:00]|3270.9800000000003|
|13199.0|[2011-10-27 20:00:00, 2011-10-28 20:00:00]|1912.7999999999997|
+-----+-----+-----+
only showing top 5 rows
```

Spark ML Pipeline

```
In [30]: 1 from pyspark.sql.functions import date_format, col
2
3 preppedDataFrame = (retail_df
4   .na.fill(0)
5   .withColumn("day_of_week", date_format(col("InvoiceDate"), "EEEE"))
6   .coalesce(5)
7 )
8
9 preppedDataFrame.show(3, truncate=False)
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|InvoiceNo|StockCode|Description|Quantity|InvoiceDate|UnitPrice|CustomerID|Country|day_of_week|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|580538|23084|RABBIT NIGHT LIGHT|48|2011-12-05 08:38:00|1.79|14075.0|United Kingdom|Monday|
|580538|23077|DOUGHNUT LIP GLOSS|20|2011-12-05 08:38:00|1.25|14075.0|United Kingdom|Monday|
|580538|22906|12 MESSAGE CARDS WITH ENVELOPES|24|2011-12-05 08:38:00|1.65|14075.0|United Kingdom|Monday|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 3 rows
```

```
In [32]: 1 preppedDataFrame.where(F.isnull(F.col("InvoiceDate"))).show(4)
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|InvoiceNo|StockCode|Description|Quantity|InvoiceDate|UnitPrice|CustomerID|Country|day_of_week|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
```

```
In [33]: 1 # split data into (train,test)
2 trainDataFrame = preppedDataFrame.where("InvoiceDate < '2011-07-01'")
3
4 testDataFrame = preppedDataFrame.where("InvoiceDate >= '2011-07-01'")
```

```
In [34]: 1 trainDataFrame.show(3)
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	day_of_week
537226	22811	SET OF 6 T-LIGHTS...	6	2010-12-06 08:34:00	2.95	15987.0	United Kingdom	Monday
537226	21713	CITRONELLA CANDLE...	8	2010-12-06 08:34:00	2.1	15987.0	United Kingdom	Monday
537226	22927	GREEN GIANT GARDE...	2	2010-12-06 08:34:00	5.95	15987.0	United Kingdom	Monday

only showing top 3 rows

In [36]:

```
1 display(testDataFrame.toPandas())
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	day_of_week
0	580538	23084	RABBIT NIGHT LIGHT	48	2011-12-05 08:38:00	1.79	14075.0	United Kingdom	Monday
1	580538	23077	DOUGHNUT LIP GLOSS	20	2011-12-05 08:38:00	1.25	14075.0	United Kingdom	Monday
2	580538	22906	12 MESSAGE CARDS WITH ENVELOPES	24	2011-12-05 08:38:00	1.65	14075.0	United Kingdom	Monday
3	580538	21914	BLUE HARMONICA IN BOX	24	2011-12-05 08:38:00	1.25	14075.0	United Kingdom	Monday
4	580538	22467	GUMBALL COAT RACK	6	2011-12-05 08:38:00	2.55	14075.0	United Kingdom	Monday
...
296001	562595	84818	DANISH ROSE PHOTO FRAME	24	2011-08-07 15:52:00	0.79	17602.0	United Kingdom	Sunday
296002	562595	47343A	FUSCHIA FLOWER PURSE WITH BEADS	12	2011-08-07 15:52:00	0.83	17602.0	United Kingdom	Sunday
296003	562595	15044C	PURPLE PAPER PARASOL	6	2011-08-07 15:52:00	2.95	17602.0	United Kingdom	Sunday
296004	562595	15044D	RED PAPER PARASOL	12	2011-08-07 15:52:00	2.95	17602.0	United Kingdom	Sunday
296005	562595	15044A	PINK PAPER PARASOL	6	2011-08-07 15:52:00	2.95	17602.0	United Kingdom	Sunday

296006 rows × 9 columns

```
In [37]: 1 # pre-processing features
2
3 from pyspark.ml.feature import StringIndexer
4 from pyspark.ml.feature import OneHotEncoder
5 from pyspark.ml.feature import VectorAssembler
6
7 indexer = StringIndexer()\
8     .setInputCol("day_of_week")\
9     .setOutputCol("day_of_week_index")
10
11 encoder = OneHotEncoder()\
12     .setInputCol("day_of_week_index")\
13     .setOutputCol("day_of_week_encoded")
14
15 vectorAssembler = VectorAssembler()\
16     .setInputCols(["UnitPrice", "Quantity", "day_of_week_encoded"])\
17     .setOutputCol("features")
```

```
In [38]: 1 # setup pipeline
2
3 from pyspark.ml import Pipeline
4
5 transformationPipeline = Pipeline()\
6     .setStages([indexer, encoder, vectorAssembler])
```

```
In [39]: 1 # run pipeline on train data
2
3 fittedPipeline = transformationPipeline.fit(trainDataFrame)
```

```
In [40]: 1 # verify on train data
2
3 transformedTraining = fittedPipeline.transform(trainDataFrame)
```

```
In [41]: 1 transformedTraining.show(3, truncate=False, vertical=True)
```

```
-RECORD 0-----
InvoiceNo      | 537226
StockCode      | 22811
Description    | SET OF 6 T-LIGHTS CACTI
Quantity       | 6
InvoiceDate    | 2010-12-06 08:34:00
UnitPrice      | 2.95
CustomerID     | 15987.0
Country        | United Kingdom
day_of_week    | Monday
day_of_week_index | 2.0
day_of_week_encoded | (5,[2],[1.0])
features       | (7,[0,1,4],[2.95,6.0,1.0])
-RECORD 1-----
InvoiceNo      | 537226
StockCode      | 21713
Description    | CITRONELLA CANDLE FLOWERPOT
Quantity       | 8
InvoiceDate    | 2010-12-06 08:34:00
UnitPrice      | 2.1
```

Spark ML Clustering

```
In [42]: 1 # COMMAND -----
2
3 from pyspark.ml.clustering import KMeans
4
5 kmeans = KMeans()\
6     .setK(20)\
7     .setSeed(10)
```

```
In [43]: 1 # COMMAND -----
2
3 kmModel = kmeans.fit(transformedTraining)
```

```
In [44]: 1 type(kmModel)
```

```
Out[44]: pyspark.ml.clustering.KMeansModel
```

In [45]: 1 kmModel.summary

Out[45]: <pyspark.ml.clustering.KMeansSummary at 0x7f452f764a60>

In [46]: 1 spark.stop()

In []: 1