```
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-09-data-src")\
           8      .getOrCreate()
           9
          10  import os
          11  SPARK_BOOK_DATA_PATH = os.environ['SPARK_BOOK_DATA_PATH']
```

```
In [3]:    1  spark
```

Out[3]: **SparkSession - hive**

**SparkContext**

Spark UI (http://172.17.0.1:4053)

**Version**

 v3.0.1

**Master**

 local[*]

**AppName**

 PySparkShell

## CSV

```
In [26]:   1  file_path = SPARK_BOOK_DATA_PATH + "/data/flight-data/csv/2010-summa
           2
           3  csvFile = spark.read.format("csv")\
           4    .option("header", "true")\
           5    .option("mode", "FAILFAST")\
           6    .option("inferSchema", "true")\
           7    .load(file_path)
```

```
In [27]:   1  csvFile.show(5)
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|    United States|            Romania|    1|
|    United States|            Ireland|  264|
|    United States|              India|   69|
|            Egypt|      United States|   24|
|Equatorial Guinea|      United States|    1|
+-----------------+-------------------+-----+
only showing top 5 rows
```

```
In [6]:   1  # COMMAND ----------
          2
          3  csvFile.write.format("csv").mode("overwrite").option("sep", "\t")\
          4    .save("/tmp/my-tsv-file.tsv")
```

## Json

```
In [8]:   1  # COMMAND ----------
          2
          3  file_path = SPARK_BOOK_DATA_PATH + "/data/flight-data/json/2010-summ
          4  jsonFile = spark.read.option("mode","FAILFAST").option("inferSchema'
          5  jsonFile.show(5)
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|    United States|            Romania|    1|
|    United States|            Ireland|  264|
|    United States|              India|   69|
|            Egypt|      United States|   24|
|Equatorial Guinea|      United States|    1|
+-----------------+-------------------+-----+
only showing top 5 rows
```

```
In [9]:   1  jsonFile.printSchema()
```

```
root
 |-- DEST_COUNTRY_NAME: string (nullable = true)
 |-- ORIGIN_COUNTRY_NAME: string (nullable = true)
 |-- count: long (nullable = true)
```

```
In [10]:   1  # COMMAND ----------
           2
           3  jsonFile.write.format("json").mode("overwrite").save("/tmp/my-json-1
```

## Parquet

```
In [11]:   1  # COMMAND ----------
           2  file_path = SPARK_BOOK_DATA_PATH + "/data/flight-data/parquet/2010-
           3  df = spark.read.format("parquet").load(file_path)
           4
           5  df.show(5)
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|    United States|            Romania|    1|
|    United States|            Ireland|  264|
|    United States|              India|   69|
|            Egypt|      United States|   24|
|Equatorial Guinea|      United States|    1|
+-----------------+-------------------+-----+
only showing top 5 rows
```

```
In [12]:   1  # COMMAND ----------
           2
           3  df.write.format("parquet").mode("overwrite")\
           4    .save("/tmp/my-parquet-file.parquet")
```

## Orc

```
In [13]:   1  # COMMAND ----------
           2  file_path = SPARK_BOOK_DATA_PATH + "/data/flight-data/orc/2010-summa
           3  df = spark.read.format("orc").load(file_path)
           4
           5  df.show(5)
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|    United States|            Romania|    1|
|    United States|            Ireland|  264|
|    United States|              India|   69|
|            Egypt|      United States|   24|
|Equatorial Guinea|      United States|    1|
+-----------------+-------------------+-----+
only showing top 5 rows
```

```
In [14]:   1  # COMMAND ----------
           2
           3  df.write.format("orc").mode("overwrite").save("/tmp/my-json-file.orc
```

## Database - Sqlite

https://intellipaat.com/community/9608/how-to-load-table-from-sqllite-db-file-from-pyspark
(https://intellipaat.com/community/9608/how-to-load-table-from-sqllite-db-file-from-pyspark)

https://repo1.maven.org/maven2/org/xerial/sqlite-jdbc/3.27.2.1/
(https://repo1.maven.org/maven2/org/xerial/sqlite-jdbc/3.27.2.1/)

/home/wengong/spark/spark-3.0.1-bin-hadoop2.7/jars/sqlite-jdbc-3.27.2.1.jar

In [7]:
```
1  # COMMAND ----------
2  file_path = SPARK_BOOK_DATA_PATH + "/data/flight-data/jdbc/my-sqlite
3  driver = "org.sqlite.JDBC"
4  path = file_path
5  url = "jdbc:sqlite:" + path
6  tablename = "flight_info"
```

In [8]:
```
1  file_path
```

Out[8]: '/home/wengong/spark_data//data/flight-data/jdbc/my-sqlite.db'

In [9]:
```
1  # COMMAND ----------
2
3  dbDataFrame = spark.read.format("jdbc")\
4      .option("url", url)\
5      .option("dbtable", tablename)\
6      .option("driver",  driver)\
7      .load()
```

In [10]:
```
1  dbDataFrame.show(5)
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|    United States|            Romania|    1|
|    United States|            Ireland|  264|
|    United States|              India|   69|
|            Egypt|      United States|   24|
|Equatorial Guinea|      United States|    1|
+-----------------+-------------------+-----+
only showing top 5 rows
```

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

```
root
 |-- DEST_COUNTRY_NAME: string (nullable = true)
 |-- ORIGIN_COUNTRY_NAME: string (nullable = true)
 |-- count: decimal(20,0) (nullable = true)
```

```
1  # COMMAND ----------
2
3  pgDF = spark.read.format("jdbc")\
4    .option("driver", "org.postgresql.Driver")\
5    .option("url", "jdbc:postgresql://database_server")\
6    .option("dbtable", "schema.tablename")\
7    .option("user", "username").option("password", "my-secret-
   password").load()
```

```
In [12]:    1  # COMMAND ----------
            2
            3  dbDataFrame.filter("DEST_COUNTRY_NAME in ('Anguilla', 'Sweden')").e>
```

```
== Physical Plan ==
*(1) Scan JDBCRelation(flight_info) [numPartitions=1] [DEST_COUNTRY_NA
ME#44,ORIGIN_COUNTRY_NAME#45,count#46] PushedFilters: [*In(DEST_COUNTR
Y_NAME, [Anguilla,Sweden])], ReadSchema: struct<DEST_COUNTRY_NAME:stri
ng,ORIGIN_COUNTRY_NAME:string,count:decimal(20,0)>
```

```
In [13]:    1  # COMMAND ----------
            2
            3  pushdownQuery = """(SELECT DISTINCT(DEST_COUNTRY_NAME) FROM flight_i
            4    AS flight_info"""
            5  dbDataFrame = spark.read.format("jdbc")\
            6    .option("url", url).option("dbtable", pushdownQuery).option("drive
            7    .load()
```

```
In [16]:    1  dbDataFrame.show(3)
```

```
+-----------------+
|DEST_COUNTRY_NAME|
+-----------------+
|    United States|
|            Egypt|
|Equatorial Guinea|
+-----------------+
only showing top 3 rows
```

```
In [17]:    1  dbDataFrame.explain()
```

```
== Physical Plan ==
*(1) Scan JDBCRelation((SELECT DISTINCT(DEST_COUNTRY_NAME) FROM flight
_info)
  AS flight_info) [numPartitions=1] [DEST_COUNTRY_NAME#63] PushedFilte
rs: [], ReadSchema: struct<DEST_COUNTRY_NAME:string>
```

```
In [18]:    1  # COMMAND ----------
            2
            3  dbDataFrame = spark.read.format("jdbc")\
            4    .option("url", url).option("dbtable", tablename).option("driver",
            5    .option("numPartitions", 10).load()
```

```
In [19]:  1  dbDataFrame.show(3)
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|    United States|            Romania|    1|
|    United States|            Ireland|  264|
|    United States|              India|   69|
+-----------------+-------------------+-----+
only showing top 3 rows
```

```
In [20]:  1  dbDataFrame.explain()
```

```
== Physical Plan ==
*(1) Scan JDBCRelation(flight_info) [numPartitions=1] [DEST_COUNTRY_NA
ME#75,ORIGIN_COUNTRY_NAME#76,count#77] PushedFilters: [], ReadSchema:
struct<DEST_COUNTRY_NAME:string,ORIGIN_COUNTRY_NAME:string,count:decim
al(20,0)>
```

```
In [21]:  1  # COMMAND ----------
          2
          3  props = {"driver":"org.sqlite.JDBC"}
          4  predicates = [
          5    "DEST_COUNTRY_NAME = 'Sweden' OR ORIGIN_COUNTRY_NAME = 'Sweden'",
          6    "DEST_COUNTRY_NAME = 'Anguilla' OR ORIGIN_COUNTRY_NAME = 'Anguilla
          7  spark.read.jdbc(url, tablename, predicates=predicates, properties=pi
```

```
+-----------------+-------------------+-----+
|DEST_COUNTRY_NAME|ORIGIN_COUNTRY_NAME|count|
+-----------------+-------------------+-----+
|           Sweden|      United States|   65|
|    United States|             Sweden|   73|
|         Anguilla|      United States|   21|
|    United States|           Anguilla|   20|
+-----------------+-------------------+-----+
```

```
In [22]:  1  spark.read.jdbc(url,tablename,predicates=predicates,properties=props
          2    .rdd.getNumPartitions() # 2
```

Out[22]: 2

```
In [23]:  1  # COMMAND ----------
          2
          3  props = {"driver":"org.sqlite.JDBC"}
          4  predicates = [
          5    "DEST_COUNTRY_NAME != 'Sweden' OR ORIGIN_COUNTRY_NAME != 'Sweden''
          6    "DEST_COUNTRY_NAME != 'Anguilla' OR ORIGIN_COUNTRY_NAME != 'Anguil
          7  spark.read.jdbc(url, tablename, predicates=predicates, properties=pi
```

Out[23]: 510

```
In [25]:   1  # COMMAND ----------
           2
           3  colName = "count"
           4  lowerBound = 0
           5  upperBound = 348113 # this is the max count in our database
           6  numPartitions = 10
           7
           8
           9  # COMMAND ----------
          10
          11  spark.read.jdbc(url, tablename, column=colName, properties=props,
          12                  lowerBound=lowerBound, upperBound=upperBound,
          13                  numPartitions=numPartitions).count() # 255
```

Out[25]: 255

```
In [28]:   1  # COMMAND ----------
           2
           3  newPath = "jdbc:sqlite://tmp/my-sqlite.db"
           4  csvFile.write.jdbc(newPath, tablename, mode="overwrite", properties=
```

```
In [29]:   1  # COMMAND ----------
           2
           3  spark.read.jdbc(newPath, tablename, properties=props).count() # 255
```

Out[29]: 255

```
In [30]:   1  # COMMAND ----------
           2
           3  csvFile.write.jdbc(newPath, tablename, mode="append", properties=pro
```

```
In [31]:   1  # COMMAND ----------
           2
           3  spark.read.jdbc(newPath, tablename, properties=props).count() # 510
```

Out[31]: 510

```
In [32]:   1  csvFile.limit(10).select("DEST_COUNTRY_NAME", "count").show()

         +-----------------+-----+
         |DEST_COUNTRY_NAME|count|
         +-----------------+-----+
         |    United States|    1|
         |    United States|  264|
         |    United States|   69|
         |            Egypt|   24|
         |Equatorial Guinea|    1|
         |    United States|   25|
         |    United States|   54|
         |       Costa Rica|  477|
         |          Senegal|   29|
         |    United States|   44|
         +-----------------+-----+
```

**write out data by partition**

In [33]:
```
1  # COMMAND ----------
2
3  csvFile.limit(10).select("DEST_COUNTRY_NAME", "count")\
4    .write.partitionBy("count").text("/tmp/five-csv-files2py.csv")
```

In [34]:
```
1  # COMMAND ----------
2
3  csvFile.limit(10).write.mode("overwrite").partitionBy("DEST_COUNTRY_
4    .save("/tmp/partitioned-files.parquet")
5
6
7  # COMMAND ----------
```

In [ ]:
```
1
```