

Exercise 3

SparkSQL

Prior Knowledge

Unix Command Line Shell

Simple Python

Apache Spark in Jupyter (from previous exercise)

Learning Objectives

SparkSQL

Reading CSV files in Spark

Software Requirements

- Apache Spark
- Python 3.7
- Jupyter

1. Let's create a new directory for our work. Type the following at the terminal window on your virtual machine.

```
mkdir sql
cd sql
```

2. The data you need for this exercise is stored in the github repository for this course. You can link them into the current directory for easy access

```
ln -s ~/BigData/datafiles/wind/2015/*.csv .
```

3. You can check that the links to the files are there by listing the directory

```
ls -la
```

4. Ensure you have activated your bigdata conda environment and start Jupyter:

```
conda activate bigdata
jupyter notebook
```

5. Give the notebook a useful name e.g. wind.ipynb

6. Now create a cell with our line to configure tab completion:

```
%config IPCompleter.greedy=True
```

7. Run that cell.

- Now, create a new cell which will have our main code in it.

Type the following into the new cell (you don't need to type the comments):

```
# spark initialisation, needs to be done before later imports
import findspark
findspark.init()

#where the data is stored
datafiles='~/BigData/sql/*.csv'

#initialisation of spark sql session
from pyspark.sql import Row, SparkSession
spark = SparkSession.builder.getOrCreate()

# read the wind data from CSV files
df = spark.read.csv(datafiles, header=True)

# show the top 5 rows
df.show(5)
```

The df object we have is not an RDD, but instead a DataFrame. This is basically a SQL motivated construct that is similar to a Pandas or R dataframe (but not exactly the same!)

Run the cell. You should see:

```
In [1]: %config IPCompleter.greedy=True

In [2]:
import findspark
findspark.init()
from pyspark.sql import Row, SparkSession
spark = SparkSession.builder.getOrCreate()
df = spark.read.csv('/home/big/sql/*.csv', header=True)
df.show(5)

+-----+-----+-----+-----+-----+
|Station_ID|Station_Name|Location_Label|Interval_Minutes|Interval_End_Time|Wind_Velocity_Mtr_Sec|Wind_Direction_Variance_Deg|Wind_Direction_Deg|Ambient_Temperature_Deg_C|Global_Horizontal_Irradiance|
+-----+-----+-----+-----+-----+
|SF15|Warnerville Switc...|Warnerville|5|2015-01-05 00:05|0.92|0.81|148.5|0.92|1.628|
|SF15|Warnerville Switc...|Warnerville|5|2015-01-05 00:10|0.717|9.4|151.1|0.717|1.519|
|SF15|Warnerville Switc...|Warnerville|5|2015-01-05 00:15|0.627|8.7|142.7|0.627|1.482|
|SF15|Warnerville Switc...|Warnerville|5|2015-01-05 00:20|0.5|6.895|141.8|0.5|1.985|
|SF15|Warnerville Switc...|Warnerville|5|2015-01-05 00:25|0.727|8.64|144.7|0.727|1.903|
+-----+-----+-----+-----+-----+
only showing top 5 rows
```

This is data from weather stations in San Francisco showing the wind speed, direction and temperature throughout 2015.

- Before we do this as SQL, we are going to look at the data using the same map/reduce model we used previously. To do this, we will convert the DataFrame into an RDD, allowing us to do functional programming on it (map/reduce/etc). Create a new cell below the one where you have

initialised and loaded in the data, then type the following:

Note that this doesn't copy the data, but just exposes the rdd which is already hiding inside the dataframe.

```
winds = df.rdd
```

10. Let's do the normal step of mapping the data into a simple <K,V> pair.
Each column in the row can be accessed by the syntax e.g. row.Station_ID

We can therefore map our RDD with the following:

```
mapped = winds.map(lambda s: (s.Station_ID, s.Wind_Velocity_Mtr_Sec))
```

11. We need to clean up the data (remove any nulls)

```
filtered = mapped.filter(lambda s:s[1]!=None)
```

12. We can simply calculate the maximum values with this reducer:

```
maxes = filtered.reduceByKey(lambda a, b: a if (a>b) else b)
```

13. And once again collect / print:

```
for (k,v) in maxes.collect(): print (k,v)
```

14. Now run the cell. You should see something like:

```
: winds=df.rdd
mapped=winds.map(lambda s: (s.Station_ID,s.Wind_Velocity_Mtr_Sec))
maxes=mapped.reduceByKey(lambda a,b: a if a>b else b)
for (k,v) in maxes.collect():
    print(k,v)

(u'SF36', u'9.71')
(u'SF18', u'9.85')
(u'SF37', u'7.079')
(u'SF04', u'9.92')
(u'SF15', u'7.92')
(u'SF17', u'5.767')
```

15. You can also turn the response of a collect into a Python dictionary, which is handy. Try these (in new cells):

```
print(maxes.collectAsMap())
print(maxes.collectAsMap()['SF04'])
```

PART B – Using SQL

16. There is an easier way to do all this if you are willing to write some SQL.
First, in a new cell, we need to give our DataFrame a table name:
`df.registerTempTable('wind')`

17. Now we can use a simple SQL statement against our data.

```
spark.sql("SELECT Station_ID, avg(Wind_Velocity_Mtr_Sec) as
avg,max(Wind_Velocity_Mtr_Sec) as max from wind group by
Station_ID").show()
```

Now run the cell and you should see something like:

Station_ID	avg	max
SF37	2.260403505500663	7.079
SF15	1.8214145677504483	7.92
SF04	2.300981748124102	9.92
SF17	0.5183500253485376	5.767
SF18	2.2202234391695437	9.85
SF36	2.464172530911313	9.71

18. One thing you also might like is that you can convert from a Spark Dataframe to a Pandas dataframe just by calling **toPandas()**

Note that when you do this, you are collecting the results back from a cluster to a single server (the master).

19. Recap. We have:

- a. Used Spark to read in CSV files
- b. Explored Map/Reduce on those CSV files
- c. Used SQL to query the data.

20. We are going to make this into a standalone program now. Copy the python code (from the initialisation and sql cells) and paste into a file called `wind.py`

You can use Atom, PyCharm, nano or some other editor, or even just create a new text file in Jupyter .

21. It should look something like this (but you will need to make sure the location of the datafiles is correct for your machine):

```
1 import findspark
2 findspark.init()
3
4 datafiles='/home/big/sql/*.csv'
5
6 from pyspark.sql import Row, SparkSession
7 spark = SparkSession.builder.getOrCreate()
8
9 df=spark.read.csv(datafiles,header=True)
10
11 df.registerTempTable('wind')
12 spark.sql("SELECT Station_ID, avg(Wind_Velocity_Mtr_Sec) as avg,max(Wind_Velocity_Mtr_Sec) as max from wind group by
Station_ID").show()
13
14
```

22. Try running it as a standalone program:

```
spark-submit wind.py
```

23. You should see lots of log ending like this:

```
big@big:~/sql
17/12/08 13:15:32 INFO DAGScheduler: Job 5 finished: showString at NativeMethodAccessImpl.java:0, took 0.967812 s
17/12/08 13:15:32 INFO CodeGenerator: Code generated in 37.882479 ms
+-----+-----+
|Station_ID|      avg|  max|
+-----+-----+-----+
|  SF37| 2.260403505500663| 7.079|
|  SF15|1.8214145677504483| 7.92|
|  SF04| 2.300981748124102| 9.92|
|  SF17|0.5183500253485376| 5.767|
|  SF18|2.2202234391695437| 9.85|
|  SF36| 2.464172530911313| 9.71|
+-----+-----+-----+
17/12/08 13:15:32 INFO SparkContext: Invoking stop() from shutdown hook
17/12/08 13:15:32 INFO SparkUI: Stopped Spark web UI at http://10.0.2.15:4041
17/12/08 13:15:32 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
17/12/08 13:15:32 INFO MemoryStore: MemoryStore cleared
17/12/08 13:15:32 INFO BlockManager: BlockManager stopped
17/12/08 13:15:32 INFO BlockManagerMaster: BlockManagerMaster stopped
17/12/08 13:15:32 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
17/12/08 13:15:32 INFO SparkContext: Successfully stopped SparkContext
17/12/08 13:15:32 INFO ShutdownHookManager: Shutdown hook called
17/12/08 13:15:32 INFO ShutdownHookManager: Deleting directory /tmp/spark-3fee63e4-d99f-4c30-932a-9d0bd9bba9d4/pyspark-1823c4dd-ab75-498b-9d93-b37af95f3430
17/12/08 13:15:32 INFO ShutdownHookManager: Deleting directory /tmp/spark-3fee63e4-d99f-4c30-932a-9d0bd9bba9d4
big@big:~/sql$ $
```

In the next exercise we are going to run this code in the cloud using AWS.
That's all for now.

Congratulations, this lab is complete.