

# Exercise 7

*Get started with Apache Spark and Python*

## Prior Knowledge

Unix Command Line Shell

Simple Python

## Learning Objectives

Understand the Spark system

Understand the Jupyter Notebook model

WordCount!

## Software Requirements

(see separate document for installation of these)

- Apache Spark 3.0.0
- Python 3.8.x
- Jupyter notebooks

## Part A. Spark Python Shell (pySpark)

1. We are going to do a wordcount against a set of books downloaded from Project Gutenberg. Wordcount is the definitive Big Data program (sort of Hello World for Big Data) and it is frankly embarrassing that we haven't done one yet.
2. Apache Spark has a useful Python shell, which we can use to interactively test and run code.
3. We are going to start by using data in HDFS, so *we need to ensure HDFS is running*. (Follow the instructions from the Hadoop lab).
4. Let's load some books into HDFS. In a terminal window (Ctrl-Alt-T)

```
hadoop fs -mkdir -p /user/oxclo/books  
hadoop fs -put ~/datafiles/books/* /user/oxclo/books/
```

5. Let's make a directory for our code:

```
mkdir ~/pse  
cd ~/pse
```

6. Now start the Spark Python command line tool –

```
pyspark
```

7. In the command-line you will see something like

```
[I 13:53:23.865 NotebookApp] Serving notebooks from local directory: /home/oxclo/pse
[I 13:53:23.866 NotebookApp] 0 active kernels
[I 13:53:23.866 NotebookApp] The Jupyter Notebook is running at:
http://localhost:8888/?token=fd655aab32ed4840ceb47b8b7392b1243a27f56350888a91
[I 13:53:23.866 NotebookApp] Use Control-C to stop this server and shut down all
kernels (twice to skip confirmation).
[C 13:53:23.868 NotebookApp]

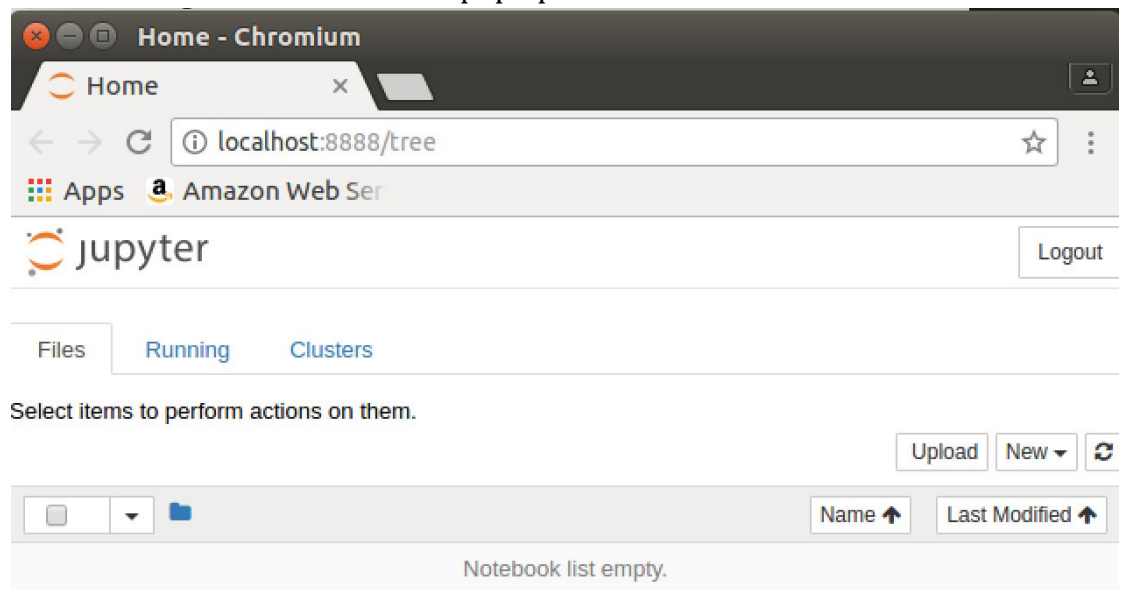
Copy/paste this URL into your browser when you connect for the first time,
to login with a token:
http://localhost:8888/?token=fd655aab32ed4840ceb47b8b7392b1243a27f56350888a91
```

8. Jupyter is starting. This is because I have preconfigured the following environment variables in `~/.bashrc`

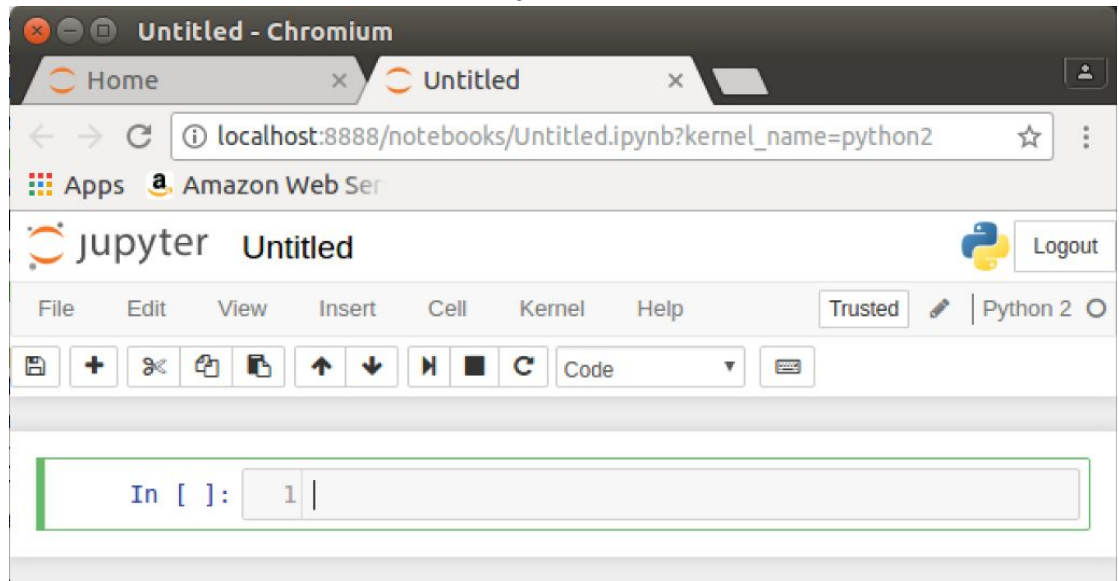
```
export PYSPARK_DRIVER_PYTHON=jupyter
export PYSPARK_DRIVER_PYTHON_OPTS='notebook'
```

If you want to use the more basic version of pyspark, you'll have to unset these.

9. And then a browser window will pop up.



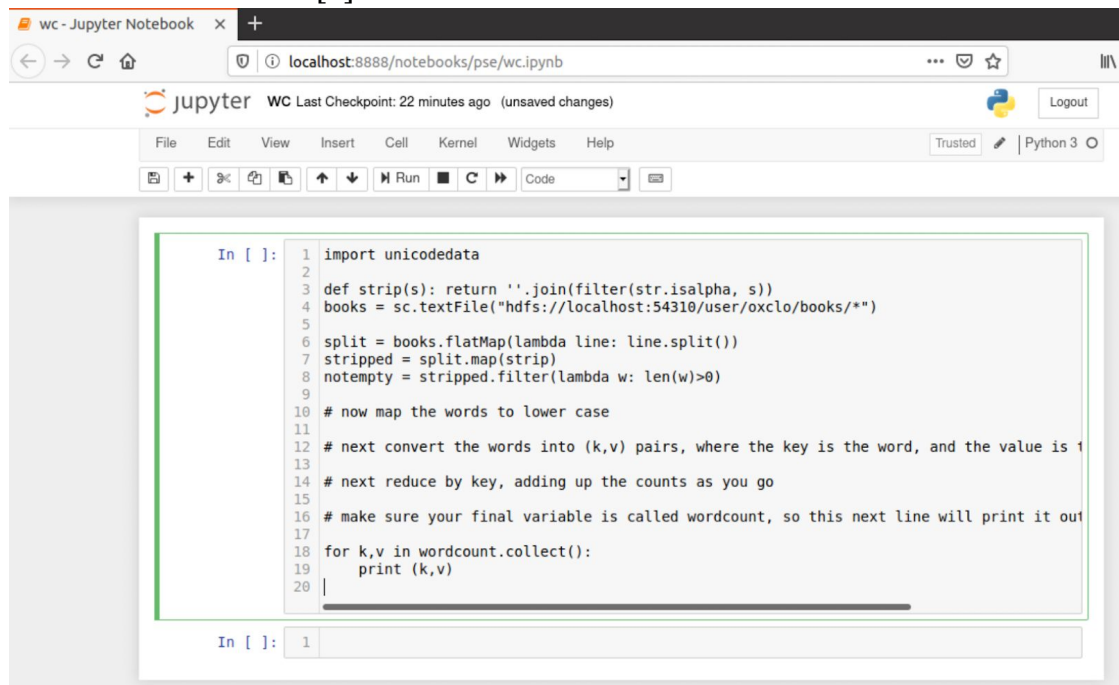
10. Use the **New** button to create a new Python3 notebook:



11. There is a starter of the code you need in the following URL:

<https://freo.me/first-notebook>

Paste that into the cell [1] so it looks like this:



12. There are some aspects that are not filled in that you need to write.  
Basically this is a data-processing pipeline (also a directed acyclic graph)

13. *Let's look at the parts that are there already.*

14. We already have a `SparkContext` object defined in the notebook (in a program you need to define one, which we will see later)
15. We want to remove any non-alphanumeric characters. This is a function that will do that:

```
def strip(s): return ''.join(filter(str.isalpha, s))
```

16. With the preliminaries over, the next line loads the data in:

```
books =  
sc.textFile("hdfs://localhost:54310/user/oxclo/books/*")
```

17. Then splits the lines into separate words

```
split = books.flatMap(lambda line: line.split())
```

18. Removes non-alpha characters:

```
stripped = asc.map(strip)
```

and removes empty items:

```
notempty = stripped.filter(lambda w: len(w)>0)
```

19. Now it is time for you to do something!

Convert all the words to lower case, using a map operation. In python, if `str` is a string, then `str.lower()` is the same string in lower case.


20. Now you need to get ready for a reduce. In order to do a reduce, we need some form of *key, value* pairs. I recommend using *tuples* which are simply `(k,v)` in Python (the brackets group the items into a tuple).

21. Remembering how reduce works, we need each word to have a count. Before reducing, that count is 1. So we need a lambda that takes a word `w` and returns `(w,1)`

22. Now we can do a reduce that adds all those counts together.

23. Finally, we need to collect the results and print them. In Spark, they may be distributed across different RDD partitions on different machines, so the `collect()` method brings them together.

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

24. Try running the cell, by clicking 

25. Be patient. I suggest you look at the command window and wait until you see spark start working.

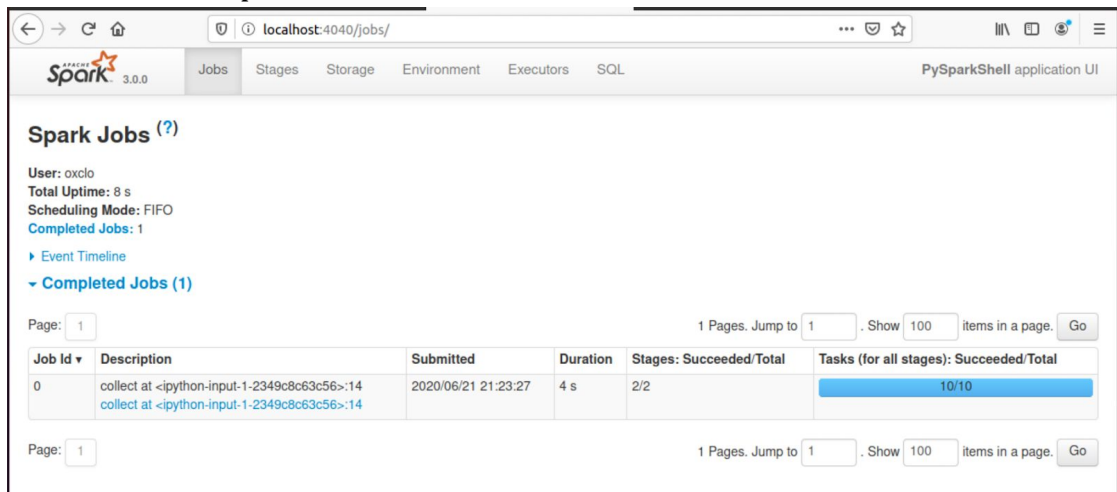
26. You should see a word count appear below cell 1:

```
systematic 7  
parallelogram 1  
sowell 1  
presnya 1  
four 265  
conjuring 1  
chamberagain 1  
marching 32  
sevens 4  
awistocwacy 1  
trotat 1  
canes 1  
shipnets 1  
understandthat 2  
lorn 16  
lore 1  
inwards 2  
wickam 62  
utterand 1  
alright 1
```

27. Congratulations!

28. While the pyspark is still running browse to <http://localhost:4040>

29. You will see the Spark web console:



The screenshot shows the Spark web console interface. At the top, there's a navigation bar with tabs for Jobs, Stages, Storage, Environment, Executors, and SQL. The 'Jobs' tab is selected. Below the navigation bar, the 'Spark Jobs (?)' section is visible. It shows the user 'oxclo', total uptime of 8 s, scheduling mode of FIFO, and 1 completed job. A table titled 'Completed Jobs (1)' displays the details of the job. The table has columns for Job Id, Description, Submitted, Duration, Stages: Succeeded/Total, and Tasks (for all stages): Succeeded/Total. The job details are as follows:

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
0	collect at <ipython-input-1-2349c8c63c56>:14 collect at <ipython-input-1-2349c8c63c56>:14	2020/06/21 21:23:27	4 s	2/2	10/10

30. Click on the blue link “collect at ipython-input”

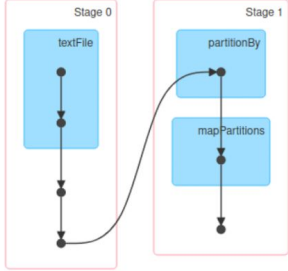
This shows you how Spark converted your code into stages:

Spark 3.0.0 Jobs Stages Storage Environment Executors SQL PySparkShell application UI

### Details for Job 0

Status: SUCCEEDED  
Completed Stages: 2

▶ Event Timeline  
▼ DAG Visualization



▼ Completed Stages (2)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
1	<a href="#">collect at &lt;ipython-input-1-2349c8c63c56&gt;:14</a> +details	2020/06/21 21:23:31	0.2 s	5/5			860.6 KiB	
0	<a href="#">reduceByKey at &lt;ipython-input-1-2349c8c63c56&gt;:12</a> +details	2020/06/21 21:23:27	4 s	5/5	10.6 MiB			860.6 KiB

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

### 31. Click on Stage 0

## Details for Stage 0 (Attempt 0)

**Total Time Across All Tasks:** 11 s

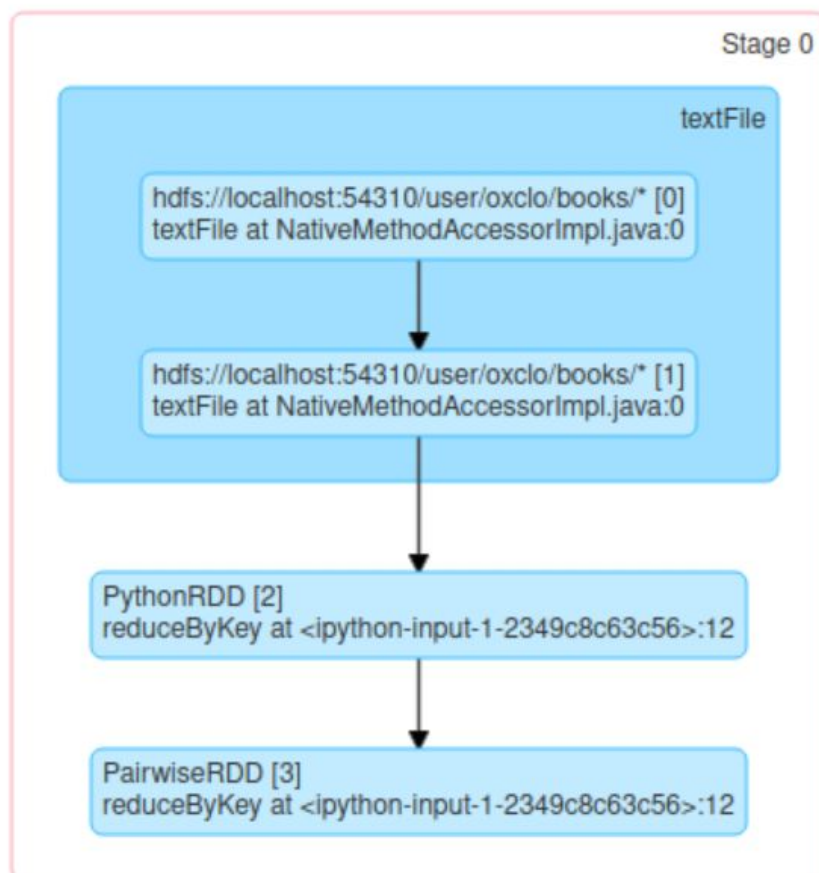
**Locality Level Summary:** Process local: 5

**Input Size / Records:** 10.6 MiB / 224113

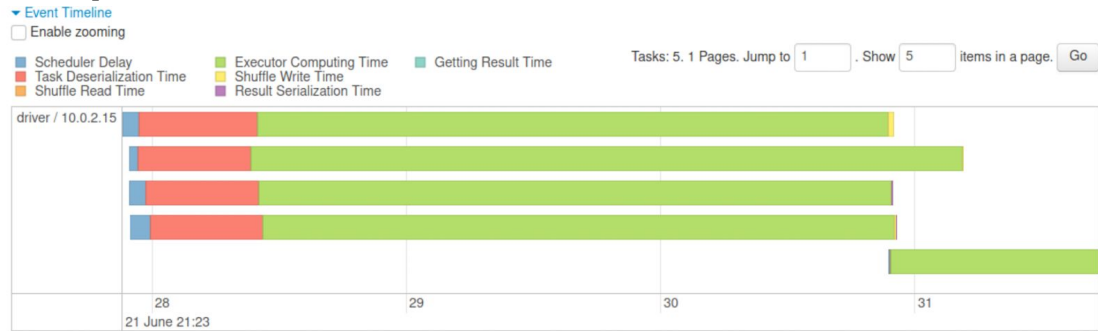
**Shuffle Write Size / Records:** 860.6 KiB / 330

**Associated Job Ids:** 0

#### ▼ DAG Visualization



### 32. And expand the Event Timeline:



#### Summary Metrics for 5 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	0.8 s	2 s	2 s	3 s	3 s
GC Time	0.0 ms	10.0 ms	10.0 ms	10.0 ms	10.0 ms
Input Size / Records	1.8 MiB / 37272	1.8 MiB / 38588	1.9 MiB / 39566	2 MiB / 43680	3.1 MiB / 65007
Shuffle Write Size / Records	140.4 KiB / 65	168.1 KiB / 65	174.4 KiB / 65	176.4 KiB / 65	201.3 KiB / 70

### 33. Make sure your code is saved from the notebook.

### 34. Quit the notebook shell by typing Ctrl-C on the command line, and then Y Also close the notebook windows in the browser.

### 35. Now let's run the same code as a "job" instead of interactively.

### 36. From <http://freo.me/oxclo-wc-py> copy the code into a file `wc-job.py`

### 37. You will notice that there is a bunch of "setup" code that we didn't need in the pyspark command line tool. That is because pyspark assumes you want all this and does it for you.

### 38. The default environment in your terminal session is telling Pyspark to use Jupyter. We need to stop that. Type:

```
unset PYSPARK_DRIVER_PYTHON
unset PYSPARK_DRIVER_PYTHON_OPTS
```

Note that this will reset when you start a new terminal window, so when you want to go back to using Jupyter, just kill this window and start another.

### 39. We run jobs locally on a single node directly on Spark:

The `local[*]` indicates to use as many threads as you have cores on your system:

```
~/spark/bin/spark-submit --master local[*] wc-job.py
"hdfs://localhost:54310/user/oxclo/books/*"
```

### 40. Congratulations, the lab is complete!

## Extension



41. Re load the code into the Jupyter notebook and now improve it to show the wordcount in descending order, starting with the most common words.