# BU.330.740 Large Scale Computing with Hadoop
## Lab 4. Simple text mining using Python Spark

Learning Goal: try Bag of Words model (aka. wordcount) in Python Spark on AWS Hadoop cluster and Jupyter notebook. Refer to lecture 4 notes on page 36-37.
Required Skills: understand basic Spark syntax

1. **Prepare:** download an e-book *Pride and Prejudice* by Jane Austen in txt format from http://www.gutenberg.org/ebooks/1342 and save in your local folder.
2. **S3:** login into AWS Educate account, go to AWS Management Console, create a S3 bucket for lab4. Create 1 folder for input text file. Upload the e-book into input text file folder, using all the default settings.
3. **EC2:** if you have an EC2 key pair, you can reuse it. Otherwise, create a new key pair.
4. **EMR:** Be careful this step is slightly different from Lab 2 and Lab 3. DO NOT CLONE the previous cluster.
   Go to EMR, click **Create Cluster**, then **Go to advanced options**.



   At Step 1, in **Release emr-5.31.0** section, make the following selections (Hadoop, Spark and Livy):
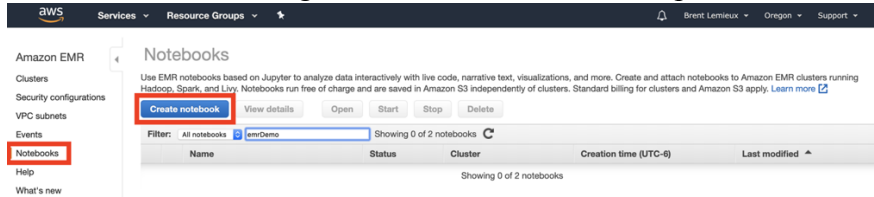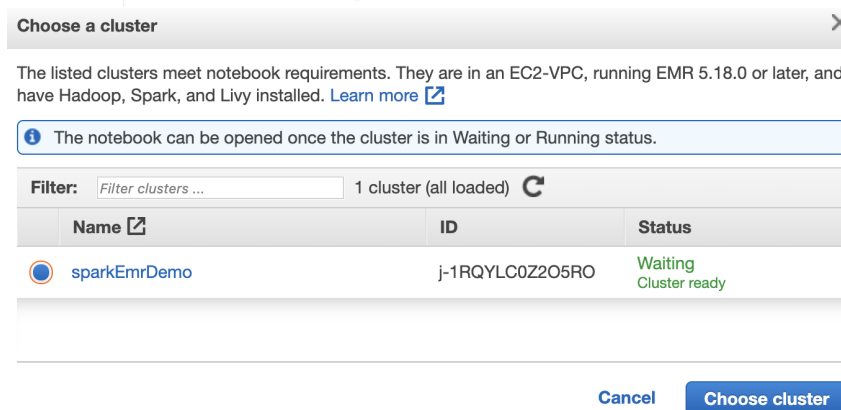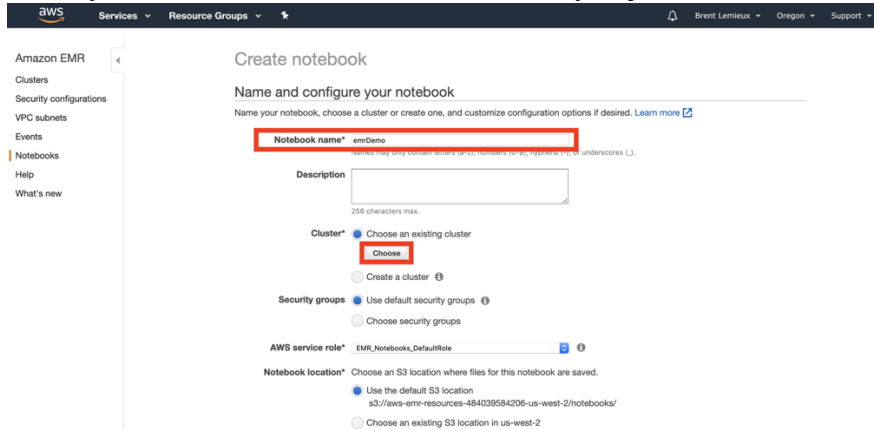


   Click **Next**. Do not change anything at Step 2, and click **Next**. Name your cluster at Step 3, then click **Next**. Select the key pair you have in Step 4 and click **Create cluster**. Your cluster will take a few minutes to start, until you see the status turns to **Waiting**.

5. **Add step:** <u>We will use Jupyter notebook to execute Python Spark scripts. The Jupyter notebook will run on the Hadoop cluster we just created.</u>
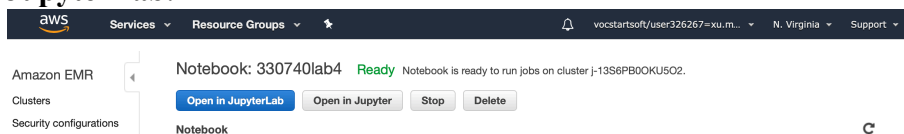On EMR console, navigate to **Notebooks** in the left panel. Click **Create notebook**.
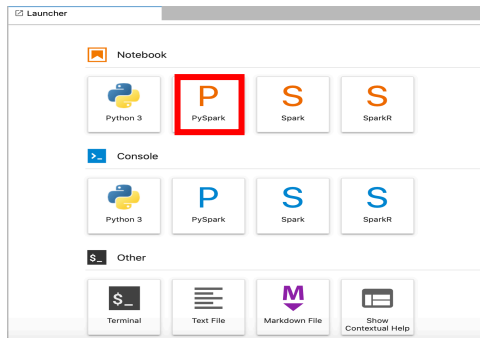


Name your notebook and choose the cluster you just created.





Then click **Create notebook**. Wait till the status shows **Ready**. Click Open in **JupyterLab**.
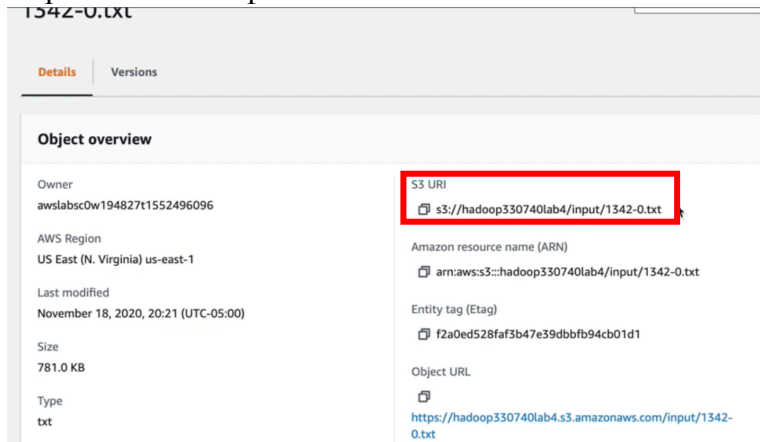


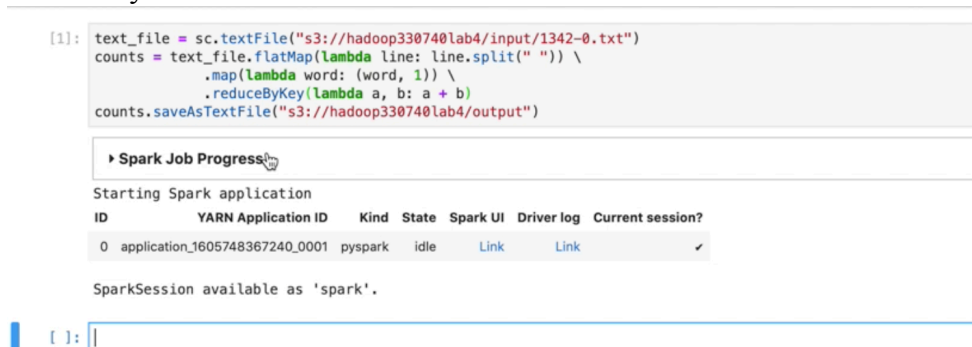Choose **Launcher->Notebook->PySpark**.

In the notebook, type in the following scripts:

```
text_file = sc.textFile("s3://bucket/folder/ebook.txt")
counts = text_file.flatMap(lambda line: line.split(" ")) \
            .map(lambda word: (word, 1)) \
            .reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("s3://bucket/output")
```

Replace the red part with your actual S3 path. Go to S3, select your S3 bucket, then your input folder, and your e-book file. Copy the **S3 URI,** and replace the text file path in red. The output path should be a folder in your S3 bucket. Note that I did not create this output folder in step 2.
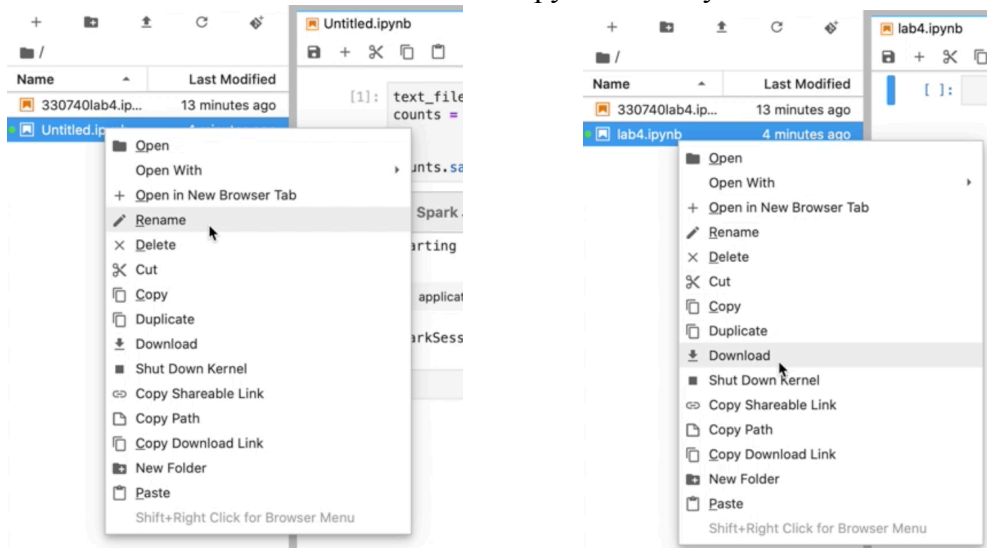


Here is my notebook screen:



Note that the input is an e-book, so it takes time to complete Spark application.

6. After the application is done, you can check and download the results in your output folder in S3.

You can also rename and download the .ipynb file for your record.



7. **DO NOT FORGET TO CLEAN UP RESOURCSES!!** Stop the notebook. Terminate the cluster (If terminate protection is on, turn it off and then terminate the cluster). Delete all S3 buckets under your account. And always double check.

Now that you are done with lab 4. You can go ahead to Assignment 2.

Reference:
https://towardsdatascience.com/getting-started-with-pyspark-on-amazon-emr-c85154b6b921