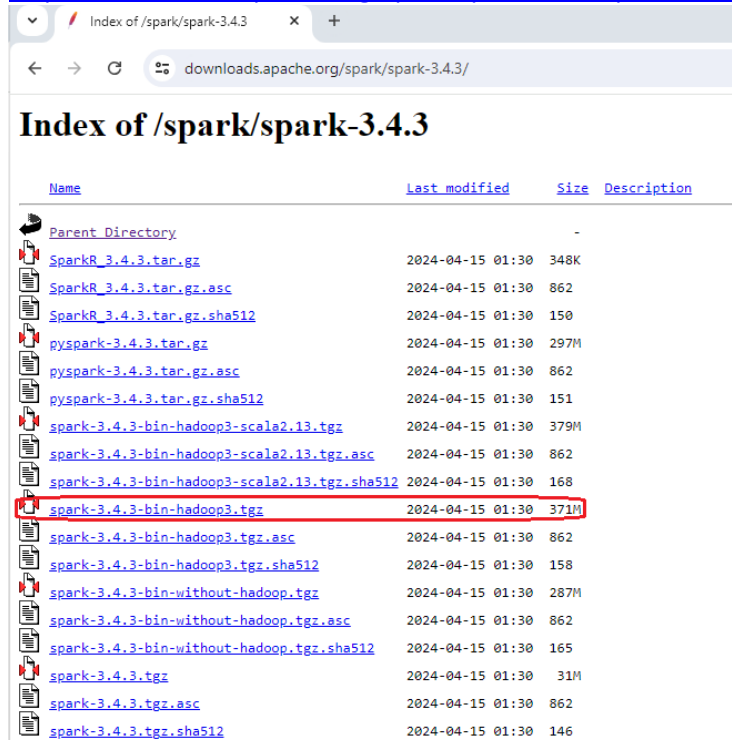# Tutorial 4 (Part II)

## Apache Spark in Jupyter

### Install and test Apache Spark in VirtualBox VM

1) Download Apache Spark 3.4.3 prebuilt for Hadoop3 from the link below in the **Downloads** folder by typing the web address in the Firefox browser

   https://downloads.apache.org/spark/spark-3.4.3/spark-3.4.3-bin-hadoop3.tgz



   After download completion, unzip it using the following command.

   Check the download folder is present or not using (**$ls**) command. Do not write $ sign as it showed prompt in Ubuntu shell

   **$cd Downloads**

   **$sudo tar -xvf spark-3.4.3-bin-hadoop3.tgz**

   <mark>OR use command line to download if you have any difficulty in downloading with the Mozilla Firefox browser.</mark> The below highlighted commands do not need to execute if you have already downloaded Apache Spark. If you have difficulty in downloading using browser, you can use terminal or shell to download Apache Spark.

   Use the following command

   **$cd /home/hduser/Downloads**

   **$wget** https://downloads.apache.org/spark/spark-3.4.3/spark-3.4.3-bin-hadoop3.tgz

   **$ls -l**

   You can check the zip folder after this download and now you can unzip the folder as mentioned below

   **$sudo tar -xvf spark-3.4.3-bin-hadoop3.tgz**

2) Install Apache Spark under /**usr/local** by running the commands
   **$sudo mv ./spark-3.4.3-bin-hadoop3 /usr/local**

   **$cd /usr/local**

3) Create a symbolic link called **spark** to the spark-3.4.3-bin-hadoop3.2:

   ```
   $sudo ln -sf ./spark-3.4.3-bin-hadoop3 spark
   ```

4) Change the ownership of the files in the **spark** directory so that the group is assigned to **hadoop** and the owner is **hduser**:

   ```
   $sudo chown -R hduser:hadoopgroup spark*
   ```
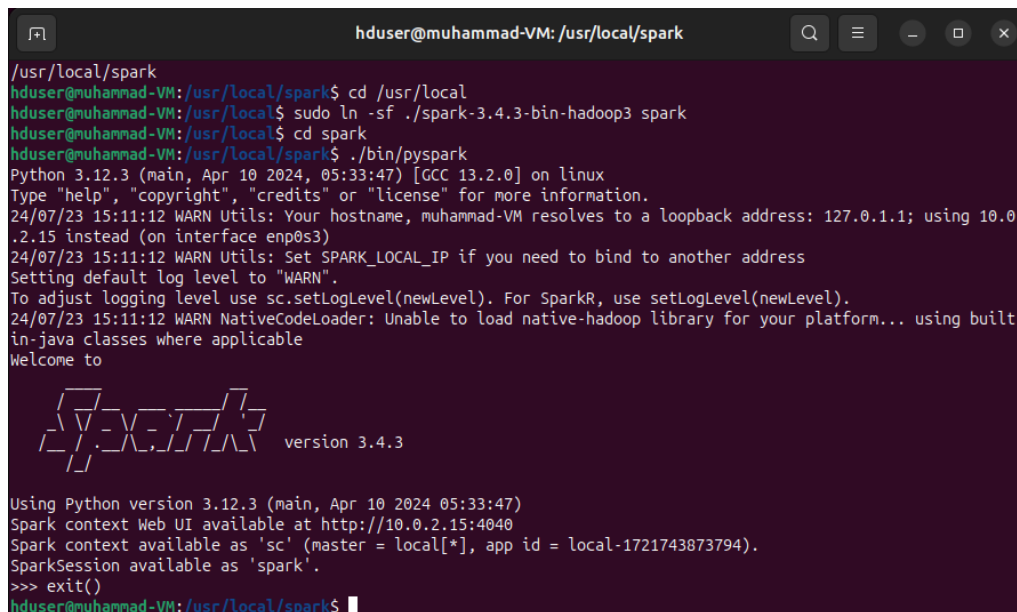
5) Test the installation at the path '**/usr/local/spark**'

   ```
   $cd /usr/local/spark
   ```
   and then
   ```
   $./bin/pyspark
   ```

   Did you get the Python shell as shown in the screenshot below? OK. If you did not get the python shell as mentioned below, then check the previous steps again.



6) Install Anaconda using the following command:

   ```
   $cd /usr/local
   $sudo wget https://repo.anaconda.com/archive/Anaconda3-2024.06-1-Linux-x86_64.sh
   $sudo bash Anaconda3-2024.06-1-Linux-x86_64.sh
   ```

Press the Enter Key and use arrow key to check the terms and conditions. Use arrow key to move to the end of this document. Then accept the license conditions as "yes" and choose option for **/usr/local/anaconda3** as the install directory and the following screenshot will appear as



When the installation of anaconda3 is completed, you will get the option for yes/ no for running anaconda.



7) Open a new terminal and add Spark environment variables to the **.bashrc** file from /home/hduser:

**$cd /home/hduser**

**$nano ./.bashrc**

Add the following lines in the ./.bashrc file using **nano** editor at the end of the file. Save the file using nano editor.

```
# Spark configuration
export SPARK_HOME=/usr/local/spark
export PATH=$PATH:$SPARK_HOME/bin
export PYSPARK_PYTHON=/usr/local/anaconda3/bin/python3
export PYSPARK_DRIVER_PYTHON=jupyter
export PYSPARK_PYTHON=python3
export PYSPARK_DRIVER_PYTHON_OPTS="notebook"
```

8) Save the above commands and source your **bashrc** with the command to update:

**$source ./.bashrc**

9) Test to see if Spark launches with the command:

**$pyspark**

**NOTE**: A web page is launched instead of the ripple launched earlier since we have opted (in the bashrc file) to launch a Juypter notebook and we will write our python code there.
If you see some errors like as mentioned on the below screenshot, then install jupyter installer



Use the command to launch **pyspark** and you should see a terminal in the browser will be opened. If you found an error like, env: 'jupyter': No such file or directory. Then install the jupyter notebook by using the command

**$sudo apt install jupyter**

Write 'y' when asked for the installation.

10) After installation of jupyter, write again **pyspark**

**$pyspark**

The jupyter notebook started automatically using google chrome browser (Installation of Google chrome is provided in Appendix A) or you can copy the link address from the terminal and paste in the Mozilla Firefox.

11) Open a new terminal and create a working directory (**Lab04**) for the Spark.

```
$cd /home/hduser/Desktop
$pwd
$mkdir Lab04
$cd Lab04
$pwd
```

12) Now we copy some text, e.g. song lyrics from: http://www.songlyrics.com/ into a file **sample_lyric.txt** or Download the file **pg30123.txt** from Moodle. Store this file in Lab04 folder.

13) Start hadoop using the commands (**$start-dfs.sh** and **$start-yarn.sh**)

14) Download the zip file from Moodle, Tutorial_4_Spark_WordCount_Example.zip, unzip it and copy both files into Lab04 folder on Desktop using the commands.

```
$cd Downloads
$ls
$unzip ./Tutorial_4_Spark_WordCount_Example.zip -d
/home/hduser/Desktop/Lab04
```



Move "**pg30123.txt**" file from Lab04 folder on your local VM into the folder on hadoop distributed file system (hdfs) before launching the wordcount code in pyspark as shown in the above screenshot.

15) Now again launch Pyspark or you can move to the folder in Lab04 in the already opened jupyter notebook.

16) In the browser Jupyter interface navigate and open the Tutorial_4_Spark_WordCount.ipynb



17) Test the code and analyse what it does. Make sure the input file should be present on the hadoop distributed file system.



18) Practice Tutorial exercises provided for further exploration Apache spark on Moodle.

## Appendix A

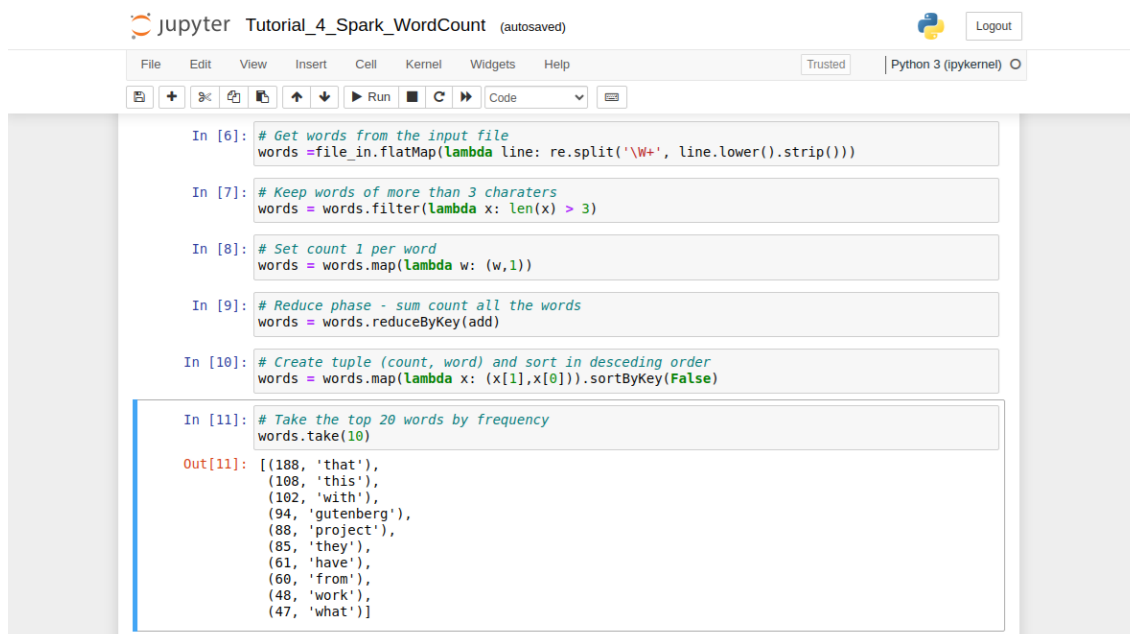If you would like to install google chrome, Follow the following steps as mentioned below

```
$sudo apt update
$ wget https://dl.google.com/linux/direct/google-chrome-stable_current_amd64.deb
$sudo dpkg -i google-chrome-stable_current_amd64.deb
$google-chrome
```

Add to your favorites when Google Chrome launches so that you won't need to use the terminal to launch it again.

**References:**

- https://spark.apache.org/examples.html
- https://www.cloudduggu.com/spark/