**Lab 0: Creating Single Node Spark/ PySpark Cluster**

Sanjana Mohile [002123793]

Northeastern University, Toronto

ALY 6110: Data Management and Big Data

Dr. Mohammad Shafiqul Islam

November 11th, 2022

## Introduction

The main aim of Lab 0 is to help us in creating an environment, where we can gain a hands-on experience with our big data knowledge. We would first install Windows Subsystem for Linux also commonly known as WSL2 for creating a single node Spark Cluster. To create this cluster we would install Java, Python3, create Jupyter Notebook, Scala, Spark with Hadoop, and finally activate spark. We would also test Spark and close the cluster. This lab gives an overall understanding of how to create environments and would help us further in learning this course.

## Installing WSL2

We would install the WSL2 using Windows Powershell. As we can see from figure 1 and figure 2, we were able to install the Ubuntu – 20.04 version for creating our environment.

```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Install the latest PowerShell for new features and improvements! https://aka.ms/PSWindows

PS C:\WINDOWS\system32> wsl --set-default-version 2
For information on key differences with WSL 2 please visit https://aka.ms/wsl2
The operation completed successfully.
PS C:\WINDOWS\system32> wsl --list --verbose
  NAME             STATE           VERSION
* Ubuntu-20.04     Running         2
PS C:\WINDOWS\system32> wsl --install
Copyright (c) Microsoft Corporation. All rights reserved.

Usage: wsl.exe [Argument] [Options...] [CommandLine]
```

*Figure 1: Installing WSL2*

```
PS C:\WINDOWS\system32> wsl --list --online
The following is a list of valid distributions that can be installed.
Install using 'wsl --install -d <Distro>'.

NAME             FRIENDLY NAME
Ubuntu           Ubuntu
Debian           Debian GNU/Linux
kali-linux       Kali Linux Rolling
openSUSE-42      openSUSE Leap 42
SLES-12          SUSE Linux Enterprise Server v12
Ubuntu-16.04     Ubuntu 16.04 LTS
Ubuntu-18.04     Ubuntu 18.04 LTS
Ubuntu-20.04     Ubuntu 20.04 LTS
```

*Figure 2: Successful installation of WSL2*

**Creating Cluster Spark**

1. Install Java

   We would be installing the Java Development Kit(JDK) for creating our environment. JDK is a cross-platform software development environment that has a collection of libraries needed to develop Java-based software applications. In our case, using only JRE (Java Runtime Environment) is not a wise choice as we are not interested in running only Java programs. The codes in figure 3 would be used to install Java –

```
sanjana_mohile@Sanjanaaaa:~$ sudo apt update && sudo apt upgrade
[sudo] password for sanjana_mohile:
Hit:1 http://archive.ubuntu.com/ubuntu focal InRelease
Get:2 http://archive.ubuntu.com/ubuntu focal-updates InRelease [114 kB]
Get:3 http://security.ubuntu.com/ubuntu focal-security InRelease [114 kB]
Get:4 http://archive.ubuntu.com/ubuntu focal-backports InRelease [108 kB]
Get:5 http://archive.ubuntu.com/ubuntu focal-updates/main amd64 Packages [2197 kB]
Fetched 2533 kB in 1s (3095 kB/s)
Reading package lists... Done
Building dependency tree
Reading state information... Done
All packages are up to date.
Reading package lists... Done
Building dependency tree
Reading state information... Done
Calculating upgrade... Done
0 upgraded, 0 newly installed, 0 to remove and 0 not upgraded.
sanjana_mohile@Sanjanaaaa:~$ sudo apt-get install openjdk-11-jre
Reading package lists... Done
Building dependency tree
Reading state information... Done
openjdk-11-jre is already the newest version (11.0.17+8-1ubuntu2~20.04).
0 upgraded, 0 newly installed, 0 to remove and 0 not upgraded.
sanjana_mohile@Sanjanaaaa:~$ sudo apt-get install openjdk-11-jdk
Reading package lists... Done
Building dependency tree
Reading state information... Done
openjdk-11-jdk is already the newest version (11.0.17+8-1ubuntu2~20.04).
0 upgraded, 0 newly installed, 0 to remove and 0 not upgraded.
```

*Figure 3: Installation of Java*

Now, we check the version of Java. This can be seen in figure 4.

```
sanjana_mohile@Sanjanaaaa:~$ java -version
openjdk version "11.0.17" 2022-10-18
OpenJDK Runtime Environment (build 11.0.17+8-post-Ubuntu-1ubuntu220.04)
OpenJDK 64-Bit Server VM (build 11.0.17+8-post-Ubuntu-1ubuntu220.04, mixed mode, sharing)
```

*Figure 4: Version of Java*

2. Installing Python

   PySpark is a Spark library that is written in Python, to run Python Applications along with Apache Spark capabilities. We can run multiple nodes parallelly using PySpark.
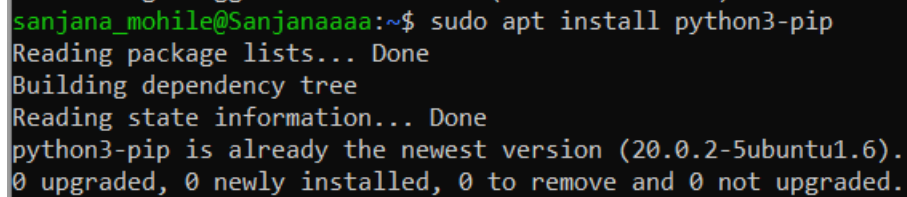
   We use the following code to install python using the Ubuntu terminal –

   ```
   $sudo apt update && upgrade

   $sudo apt install python3 python3-pip ipython3

   $sudo apt install python3-pip

   $pip3 install jupyter py4j pyspark
   ```
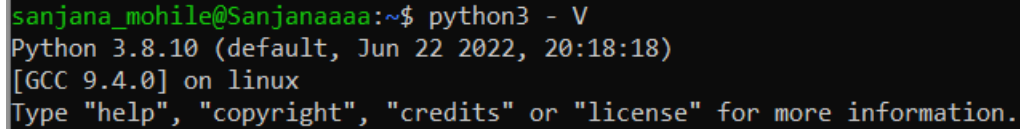


*Figure 5: Python Installation*

Figure 5 shows the installation of Python3. This is done in the Ubuntu terminal.



*Figure 6: Version of Python*

Figure 6 shows the version of Python installed on our device. As we can see, we installed the latest version of Python which was 3.8.10

*Figure 7: Installation of Jupyter notebook*

After the successful installation of python, we would now install the jupyter notebook. We use the command shown in figure 7 for its installation.

3. Installation of Scala



*Figure 8: Installing and unzipping Scala*

Figure 8 shows the command used to install Scala. Scala is an object-oriented as well as a functional programming language. Apache Spark is the most well-known open-source cluster-computing solution that is written in Scala. After the successful installation of Scala, we needed to extract the files from Scala as it was in TAR archive file (.tgz) format. This type of file is most used in Unix and Linux systems.

*Figure 9: Glimpse of bashrc file*

The bashrc file is a popular Linux distribution. We use "nano" to edit our bashrc file and add Scala to it using the following command –

*$nano ~/.bashrc*

Any changes we make in this file would automatically show when we open the file next time. Because we needed the file immediately, we use the following command -

*$source ~/.bashrc*

We can see figure 9, which shows how to add Scala in the bashrc file.



*Figure 10: Scala Version*

Now, our Scala has been successfully installed and is ready to use. As we can see, figure 10 shows the version of Scala installed in our bashrc file.

4. Install Spark with Hadoop

   We repeat the same steps we performed while installing Scala to install Spark with Hadoop. We download the TAR file of Spark with Hadoop and then unzip it. Then, we open our bashrc file again and add the same lines again to install Spark with Hadoop in the bashrc files.



*Figure 11: Downloading the TAR file and unzipping it*

*Figure 12: Successful installation of Spark with Hadoop.*

5. Installation of Jupyter Notebook

While installing Jupyter Notebook, I encountered an error. I was able to solve this problem and move ahead for creating clusters.



*Figure 13: Error in installation of Jupyter Notebook*

After qualitative research, I understood the error and was able to install new libraries needed to support the command and get rid of the error.



*Figure 14: Solution of the error*

As we can see in figure 14, installing "markupsafe" helped in removing the error.



*Figure 15: Successful installation of jupyter notebook*

Figure 15 shows the successful implementation of the jupyter notebook. Using any one of the three links mentioned in the output we can open our jupyter notebook and activate clusters thereafter.

Our system is set up with all the needed libraries and tools to run Spark. We would now use the Jupyter Notebook to activate the cluster and test it.

**Activate Spark**

To activate the cluster, we will use the following commands-

```
$cd $SPARK_HOME

$./sbin/start-master.sh

$./sbin/start-worker.sh spark://ubuntu1:7077
```

Figure 16 shows how to successfully start the clusters. We first enter Spark with Hadoop and start the mast and worker clusters. Apache spark helps us to create and stop clusters.

```
sanjana_mohile@Sanjanaaaa:~$ cd $SPARK_HOME
sanjana_mohile@Sanjanaaaa:~/spark-3.1.1-bin-hadoop3.2$ ./sbin/start-master.sh
starting org.apache.spark.deploy.master.Master, logging to /home/sanjana_mohile/spark-3.1.1-bin-hadoop3.2/logs/s
park-sanjana_mohile-org.apache.spark.deploy.master.Master-1-Sanjanaaaa.out
sanjana_mohile@Sanjanaaaa:~/spark-3.1.1-bin-hadoop3.2$ ./sbin/start-worker.sh spark://ubuntu1:7077
starting org.apache.spark.deploy.worker.Worker, logging to /home/sanjana_mohile/spark-3.1.1-bin-hadoop3.2/logs/s
park-sanjana_mohile-org.apache.spark.deploy.worker.Worker-1-Sanjanaaaa.out
```

*Figure 16: Starting the cluster*

**Testing**

```
sanjana_mohile@Sanjanaaaa:~/spark-3.1.1-bin-hadoop3.2$ spark-shell
22/11/11 23:11:12 WARN Utils: Your hostname, Sanjanaaaa resolves to a loopback address: 127.0.1.1; using 172.23.
171.1 instead (on interface eth0)
22/11/11 23:11:12 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.spark.unsafe.Platform (file:/home/sanjana_mohile/spark-3.1.1-bi
n-hadoop3.2/jars/spark-unsafe_2.12-3.1.1.jar) to constructor java.nio.DirectByteBuffer(long,int)
WARNING: Please consider reporting this to the maintainers of org.apache.spark.unsafe.Platform
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
22/11/11 23:11:13 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin
-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/11/11 23:11:19 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.
Spark context Web UI available at http://172.23.171.1:4041
Spark context available as 'sc' (master = local[*], app id = local-1668226279774).
Spark session available as 'spark'.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /___/ .__/\_,_/_/ /_/\_\   version 3.1.1
      /_/

Using Scala version 2.12.10 (OpenJDK 64-Bit Server VM, Java 11.0.17)
Type in expressions to have them evaluated.
Type :help for more information.
```

*Figure 17: Testing the Spark*

The $spark-shell command is used to check if the spark is installed and to know its version. Although it is done in the earlier parts of the installation, it is a good practice to check before we start the process.

We will now test our PySpark using Jupyter Notebook. The command $jupyter-notebook will be used to test our PySpark.
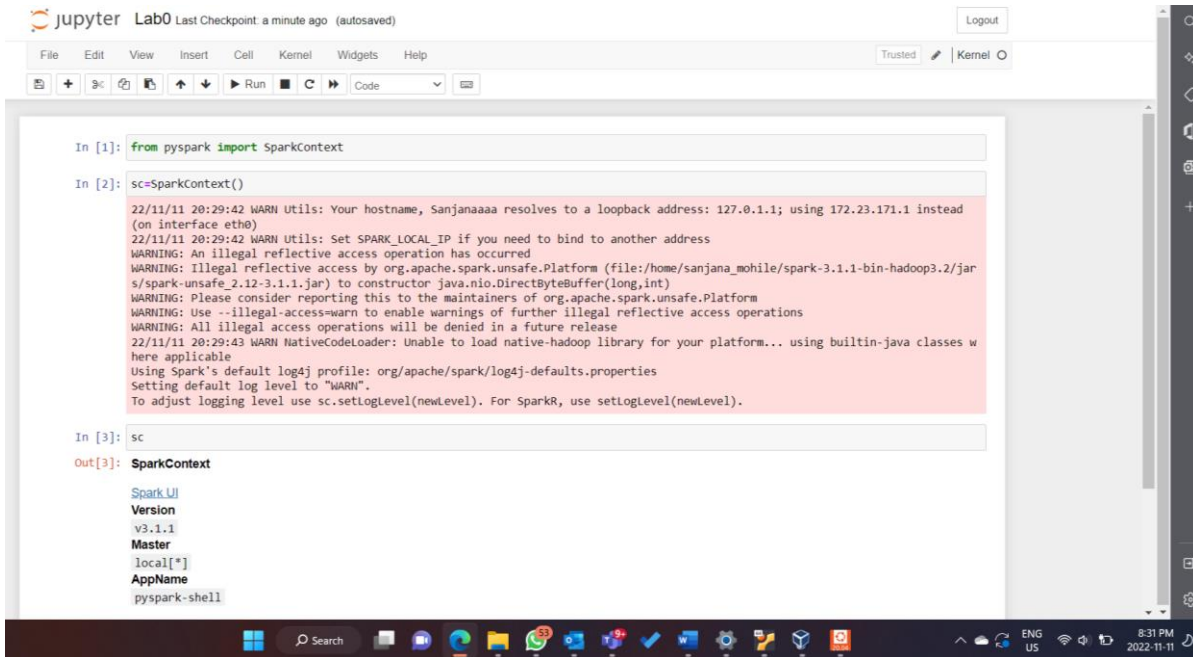
*Figure 18: Jupyter Notebook*

Figure 18 shows the successful installation of PySpark through Jupyter Notebook. We import the library SparkContext from PySpark first. Then we obtain the Spark UI version, Master, and AppName that would confirm the successful working of PySpark.

**Stop the cluster**

Our final step would be to stop the cluster after testing it. We can not leave the clustered unattended. Once we complete the testing of our cluster, we stop it to avoid heavy billings. We would use the following command to stop the cluster –

*$./sbin/stop-master.sh*

*$./sbin/stop-worker.sh*



*Figure 19: Checking the files*

Before we stop the clusters, it is important to know which and how many clusters are open. The "ls" command shows the list of files in that folder.



*Figure 20: Stopping the cluster*

Figure 20 shows the attempt of stopping our cluster. To be sure of it, we double-run the command of stopping the cluster.



*Figure 21: Checking if the clusters have stopped*

We now check if the clusters are still running or are stopped by using the commands shown in figure 21. "Sudo" means "superuser do". Hence, we are now sure that our clusters have successfully stopped.

## Conclusion

We have now successfully created an environment that will help us in practicing the various concepts of Big Data Analytics. The installation of different tools and libraries was interesting to learn and understand. We also were able to create, test, and close the clusters. We also learned about the basic steps of installing WSL2 in our systems.

## References

I. Banerjee, P. (June 17, 2022) JDK in Java. Geeks for geeks. Retrieved November 11, 2022, from https://www.geeksforgeeks.org/jdk-in-java

II. Fox, A. (January 23, 2018) What is Bashrc and why should you edit it? Make Tech Easier. Retrieved November 11, 2022, from https://www.maketecheasier.com/what-is-bashrc

III. Pedamkar, P. (September 22, 2021) Spark Shell Commands. Educba. Retrieved November 11, 2022, from https://www.educba.com/spark-shell-commands

IV. My Repository - https://github.com/SanjanaMohile/ALY-6110---Data-Management-and-Big-Data