
Lab Distributed Big Data Analytics

Worksheet-1: Setting up environment and getting started with Scala

[Dr. Hajira Jabeen](#), [Gezim Sejdiu](#), [Prof. Dr. Jens Lehmann](#)

April 18, 2017

Over the lab of the semester, we are going to use a variety of tools, including Apache Hadoop (HDFS)¹, Apache Spark², Docker³ and many more. Installing these tools and getting started can often be a hassle, mainly because of the dependencies. You are allowed to choose which option : (i) on your own machine; (ii) using Virtualbox⁴ and Docker; (iii) or via VMs, to use when you install these packages.

Tasks:

1. **Virtualbox**: - it is a virtualization software package similar to VMWare or other virtual tools. We will make use of it to setup and configure our working environment in order to complete assignments.

Here are the steps to be followed:

- a) Download the latest Ubuntu ISO from <http://www.ubuntu.com/download/desktop> (use 64 bit).
- b) Create a new virtual machine with options: Type = Linux, Version = Ubuntu (64 bit).
- c) Recommended memory size: 2GB
- d) Select: "Create a Virtual Hard Drive Now".
 - i) Leave the setting for Hard Drive File Type unchanged (i.e., VDI).
 - ii) Set the hard drive to be "*Dynamically Allocated*".
 - iii) Size: ~10GB
- e) The virtual machine is now created.
- f) Press "**Start**"
 - i) Navigate to the Ubuntu ISO that you have downloaded, and Press Start.
 - ii) On the Boot Screen: "Install Ubuntu"
 - iii) Deselect both of "Download Updates while Installing" and "Install Third-Party Software"
 - iv) Press "Continue"

¹ <http://hadoop.apache.org/>

² <http://spark.apache.org/>

³ <https://www.docker.com/>

⁴ <https://www.virtualbox.org/>

- v) Select "Erase disk and install Ubuntu"
 - vi) Add your account informations:
 - vii) Name = "yourname"; username = "username"; password = "****";
 - viii) Select "Log In Automatically"
 - ix) Press "Restart Now"
 - x)
- g) Log in to the machine.
- i) Open the terminal (Ctrl + Alt + T) and execute these commands:
 - 1) Download and upgrade the packages lists


```
sudo apt-get update
sudo apt-get upgrade
```
 - 2) Installing JAVA


```
sudo add-apt-repository ppa:webupd8team/java
sudo apt-get update
sudo apt-get install oracle-java8-installer
```

- Setting the JAVA_HOME Environment Variable

```
sudo update-alternatives --config java
sudo nano /etc/environment
```

- At the end of this file, add the following line, making sure to replace the highlighted path with your own copied path.

```
JAVA_HOME="/usr/lib/jvm/java-8-oracle"
source /etc/environment
echo $JAVA_HOME
```
 - 3) Install Maven


```
sudo apt-get update
sudo apt-get install maven
```
 - 4) Install SBT


```
echo "deb https://dl.bintray.com/sbt/debian /" | sudo tee -a
/etc/apt/sources.list.d/sbt.list
sudo apt-key adv --keyserver hkp://keyserver.ubuntu.com:80
--recv 642AC823
sudo apt-get update sudo apt-get install sbt
```
 - 5) Install Hadoop (Single Node Setup)

- Creating a Hadoop user for accessing HDFS

```
sudo addgroup hadoop
sudo adduser --ingroup hadoop hduser
```

- Installing SSH

```
sudo apt-get install openssh-server
```

Configuring SSH

#First login with hduser (and from now use only hduser account for further steps).

```
sudo su hduser
```

Generate ssh key for hduser account

```
ssh-keygen -t rsa -P ""
```

```

## Copy id_rsa.pub to authorized keys from hduser
cat $HOME/.ssh/id_rsa.pub >> $HOME/.ssh/authorized_keys
- Installations Steps
# Download latest Apache Hadoop source from Apache mirrors
wget
http://mirror.nohup.it/apache/hadoop/common/hadoop-2.7.1/hadoop-2.7.1.tar.gz
## Extract Hadoop source
tar xzf hadoop-2.7.1.tar.gz
rm hadoop-2.7.1.tar.gz
## Move hadoop-2.7.1 to hadoop folder
sudo mv hadoop-2.7.1 /usr/local
sudo ln -sf /usr/local/hadoop-2.7.1/ /usr/local/hadoop
## Assign ownership of this folder to Hadoop user./
sudo chown -R hduser:hadoop /usr/local/hadoop-2.7.1/
## Create Hadoop temp directories for Namenode and Datanode
sudo mkdir -p /usr/local/hadoop/yarn_data/hdfs/namenode
sudo mkdir -p /usr/local/hadoop_tmp/hdfs/datanode
## Again assign ownership of this Hadoop temp folder to Hadoop user
sudo chown hduser:hadoop -R /usr/local/hadoop/yarn_data/
- Update Hadoop configuration files
## User profile : Update $HOME/.bashrc
nano ~/.bashrc
# Set Hadoop-related environment variables
export HADOOP_PREFIX=/usr/local/hadoop
export HADOOP_HOME=/usr/local/hadoop
export HADOOP_MAPRED_HOME=${HADOOP_HOME}
export HADOOP_COMMON_HOME=${HADOOP_HOME}
export HADOOP_HDFS_HOME=${HADOOP_HOME}
export YARN_HOME=${HADOOP_HOME}
export HADOOP_CONF_DIR=${HADOOP_HOME}/etc/hadoop
# Native path
export
HADOOP_COMMON_LIB_NATIVE_DIR=${HADOOP_PREFIX}/lib/native
export
HADOOP_OPTS="-Djava.library.path=${HADOOP_PREFIX}/lib/native"
# Java path
export JAVA_HOME="/usr/lib/jvm/java-8-oracle"
# Add Hadoop bin/ directory to PATH
export
PATH=$PATH:$HADOOP_HOME/bin:$JAVA_PATH/bin:$HADOOP_HOME/sbin

```

In order to have the new environment variables in place, reload .bashrc
source ~/.bashrc

- Configure Hadoop

```
cd /usr/local/hadoop/etc/hadoop
```

```
nano yarn-site.xml
```

```
<configuration>
  <property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
  </property>
</configuration>
```

```
nano core-site.xml:
```

```
<configuration>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://localhost:9000</value>
  </property>
</configuration>
```

```
cp mapred-site.xml.template mapred-site.xml
```

```
nano mapred-site.xml
```

```
<configuration>
  <property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
  </property>
</configuration>
```

```
nano hdfs-site.xml:
```

```
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:/usr/local/hadoop/yarn_data/hdfs/namenode</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>file:/usr/local/hadoop/yarn_data/hdfs/datanode</value>
  </property>
</configuration>
```

- Finally, set to "/usr/lib/jvm/java-8-oracle" the JAVA_HOME variable in /usr/local/hadoop/etc/hadoop/hadoop-env.sh.

- Starting Hadoop

```
sudo su hduser
hdfs namenode -format
start-dfs.sh
start-yarn.sh
```

- Create a directory on HDFS.

```
hdfs dfs -mkdir /user
hdfs dfs -mkdir /user/hduser
```

- Run a MapReduce job.

```
hadoop jar
/usr/local/hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.
7.1.jar pi 10 50
```

- Track/Monitor/Verify

```
jps
```

For ResourceManager – [Http://localhost:8088](http://localhost:8088)

For NameNode – [Http://localhost:50070](http://localhost:50070)

Finally, to stop the hadoop daemons, simply invoke `stop-dfs.sh` and `stop-yarn.sh`.

6) Install Spark

```
mkdir $HOME/spark
cd $HOME/spark
wget
http://d3kbcqa49mib13.cloudfront.net/spark-2.1.0-bin-hadoop2.7.tgz
tar xvf spark-2.1.0-bin-hadoop2.7.tgz
nano ~/.bashrc
export SPARK_HOME=$HOME/spark/spark-2.1.0-bin-hadoop2.7/
export PATH=$PATH:$SPARK_HOME/bin
source ~/.bashrc
./sbin/start-master.sh
./sbin/start-slave.sh <master-spark-URL>
./bin/spark-shell <master-spark-URL>
```

Installing Scala

```
wget https://downloads.lightbend.com/scala/2.12.1/scala-2.12.1.tgz
sudo tar xvf scala-2.12.1.tgz
nano ~/.bashrc
```

```
export SCALA_HOME=$HOME/scala-2.12.1/
export PATH=$SCALA_HOME/bin:$PATH
source ~/.bashrc
scala -version
```

2. Configure IDE with Scala and Spark

Here are steps how to configure scala-ide for eclipse.

1. Go to <http://scala-ide.org/download/sdk.html> and download the version needed for your Linux 32 or 64-bit to your linux server/workstation. The latest version is 4.5.0.
2. Meet JDK requirements: JDK 8
3. Copy the archive to your preferred folder and decompress.
4. Find eclipse executable and execute it

After the scala-ide have been configured properly we could set up a spark-template project to get started with Scala and Spark.

```
git clone https://github.com/SANSA-Stack/SANSA-Template-Maven-Spark.git
```

5. Open eclipse.
6. Click File > Import.
7. Type Maven in the search box under Select an import source:
8. Select Existing Maven Projects.
9. Click Next.
10. Click Browse and select the folder that is the root of the Maven project (probably contains the pom.xml file)
11. Click OK.

Assignment:

IN CLASS

1. Data Store & Processing using HDFS
 - a. Start HDFS and verify its status.
 - b. Create a new directory /yourname on HDFS.
 - c. Download page_links_simple.nt.bz2, unzip it, on your local filesystem and upload it to HDFS under /yourname folder.
 - d. View the content and the size of /yourname directory.
 - e. Copy the file just created on /yourname into newfile.bck
 - f. Copy your file back to local filesystem and name it youfile.bck
 - g. Remove your file from HDFS.
 - h. Remove /yourname directory from HDFS.
2. Basics of Scala

- a. Define a class `Point` which describes an (x, y) coordinate.
- b. Create a companion object `Point` such will allow you to instantiate `Point` without using `new`.
- c. Create a singleton object `Origin` that represents the $(0, 0)$ coordinates of `Point`.
- d. Instantiate the two object of `Origin` and check if both refer to the same object in memory.
- e. Implement a function `distanceTo` which calculates the distance between two `Point` instances.

AT HOME

1. Read and explore
 - a. Functional Programming
 - b. Recursion, and Tail recursion
 - c. Anonymous functions
 - d. High Order Functions
 - e. Currying
2. Read a textfile and do a word count on that file. Hint: create and populate a `Map` with words as keys and counts of the number of occurrences of the word as values⁵.
 - a. Hello Hello World
 - b. (Hello, 2)
 - c. (World, 1)
3. Apply what you have read above by creating a function that returns sum of a range of integers by applying a user defined function to it.

$$\sum_{i=a}^b f(i)$$

e.g. `f=cube`, and `sum(2,4)` would result in $2^3+3^3+4^3$.
4. Further readings
 - a. <http://www.artima.com/scalazine/articles/steps.html>
 - b. <http://www.scala-lang.org/>
 - c. http://twitter.github.io/scala_school/

⁵ <http://ampcamp.berkeley.edu/big-data-mini-course/introduction-to-the-scala-shell.html>