## Installing Spark Pre-built for Hadoop version 2.6 and later

You still need to install Ubuntu on Virtual Machine and Java using same commands provided earlier by Ashwin(Step 1-5 remains same).

1. Download VirtualBox from Oracle Website

https://www.virtualbox.org/wiki/Downloads

2. Install Oracle VirtualBox like any other software.

3. Download Ubuntu Desktop Image from the following URL

http://www.ubuntu.com/download

4. Use the instructions in the below link to setup Ubuntu in VirtuaBox

http://www.wikihow.com/Install-Ubuntu-on-VirtualBox

5. After Installation try to installing Ubuntu, open Terminal and execute the following operations

a. Login as root user using the command – "su root"

b. # sudo apt-get update

c. Install Java in Ubuntu

       # sudo add-apt-repository ppa:webupd8team/java

       # sudo apt-get update

       # sudo apt-get install oracle-java7-installer

d. Set Java environment variable

       # vi /etc/environment

       Paste the following line in that document

       export JAVA_HOME=/usr/lib/jvm/java-7-oracle

 Press ESC key and type the following command to save the file in VI Editor

:wq

6. Download the latest Spark version 1.3.0 from http://spark.apache.org/downloads.html

       a. Select the Pre-built for Hadoop 2.6 and later

Download the .tgz file.



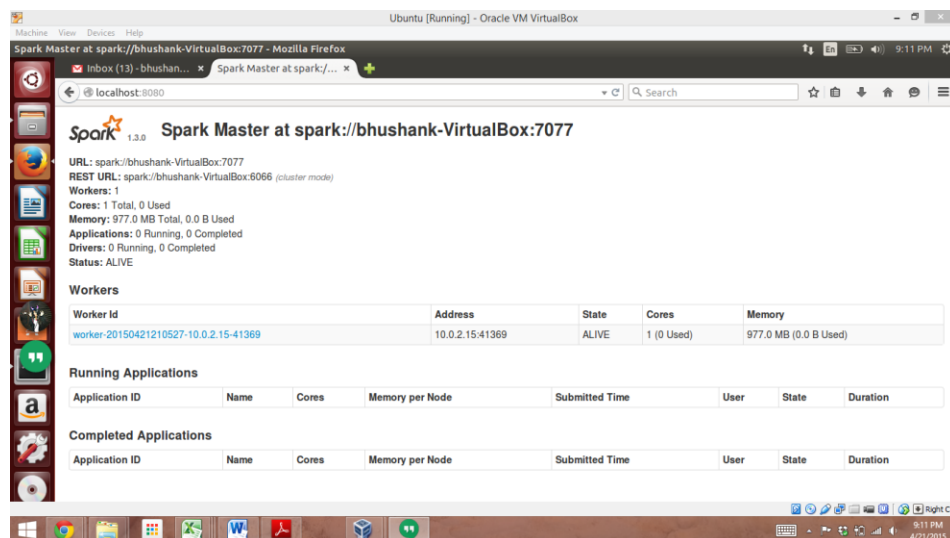b. extract the file using the following command
# tar zvfx spark-1.3.0.tgz

c. get inside the extracted folder using the following command
# cd spark-1.3.0-bin-hadoop2.6

d. Goto "sbin" folder and execute the following commands and check for **localhost:8080**
# ./stop-all.sh
# ./start-all.sh

d. Get back to spark-1.3.0-bin-hadoop2.6 folder and execute following command
# ./bin/pyspark

7. To remove logging in pyspark, execute the following command.
a. Get into "conf" folder
# cp log4j.properties.template log4j.properties
# vi log4j.properties
Change
log4j.rootCategory = INFO, console
to
log4j.rootCategory = ERROR, console

Happy Coding !!