# 1. Getting the Environment Ready doing R Analytics on Hadoop

In the section, students shall be introduced to the basic platforms and setup which should be installed to be able to do R analytics using Big Data. The three main setups include:

1. The Virtual Machine (VM) (VMWare Workstation player)
2. Hadoop Environment (Hortonworks)
3. R/RStudio on Hadoop

# 2. Installing the Virtual Machine

VMWare Workstation player (non-commercial) edition has been used. Follow the following steps to download and install the software.
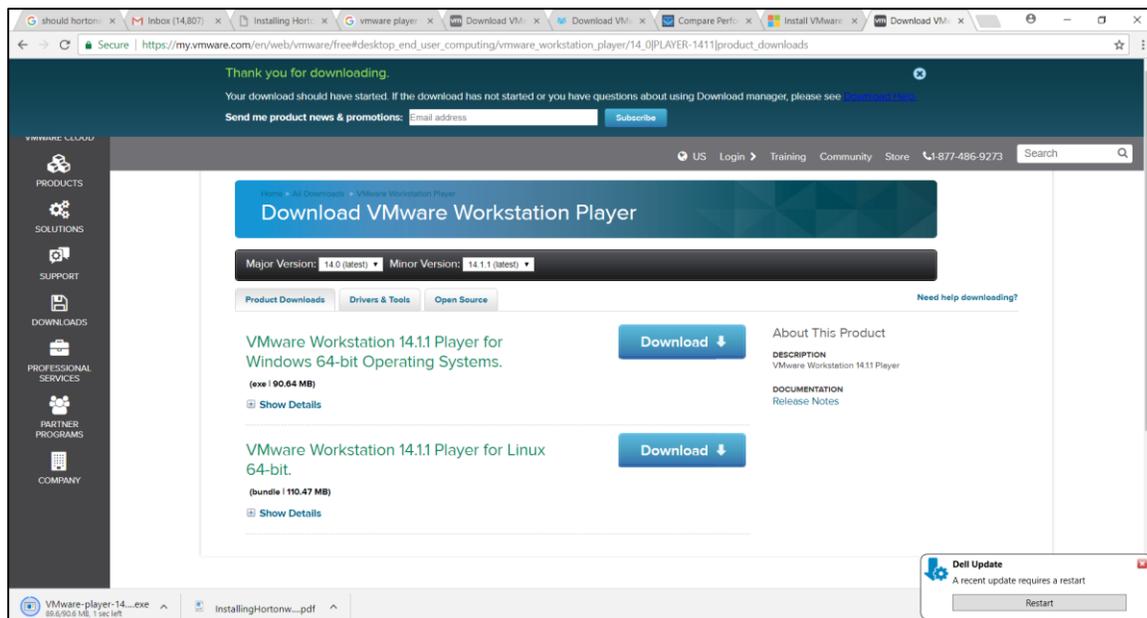
- *There might be slight variations when you download and install depending on the version that you are using and the OS that you have.*

**Note it!**

## Steps

1. Download VMWare Workstation player by searching on the web.
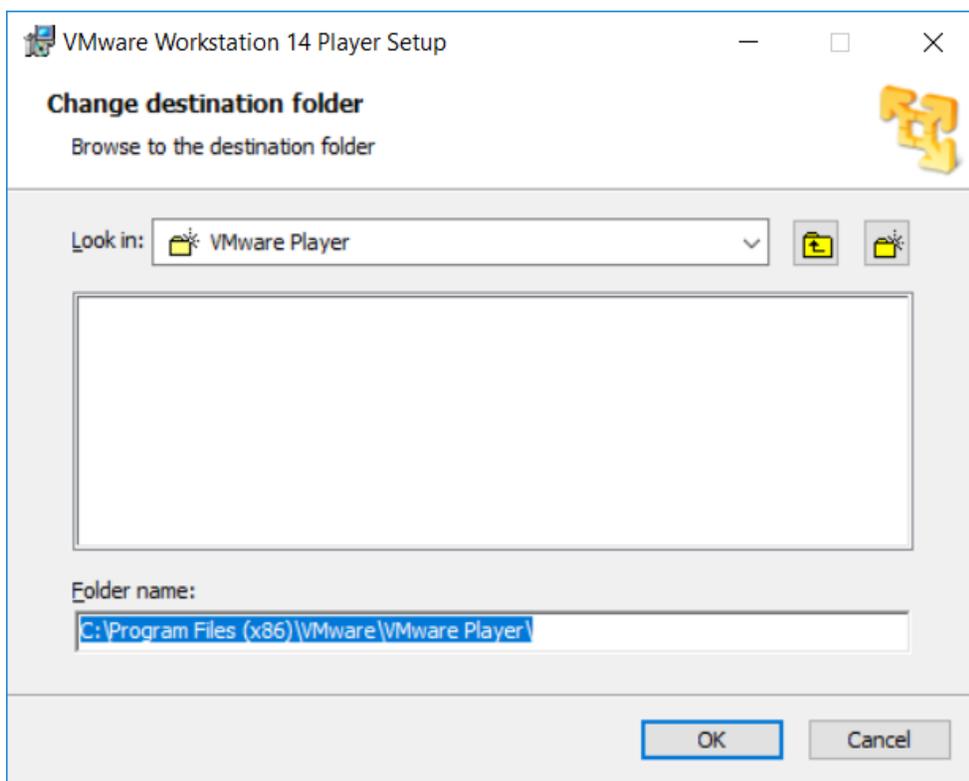
VMware Player on Windows

2. Install once download is complete. Click on Next and Accept the recommended settings.
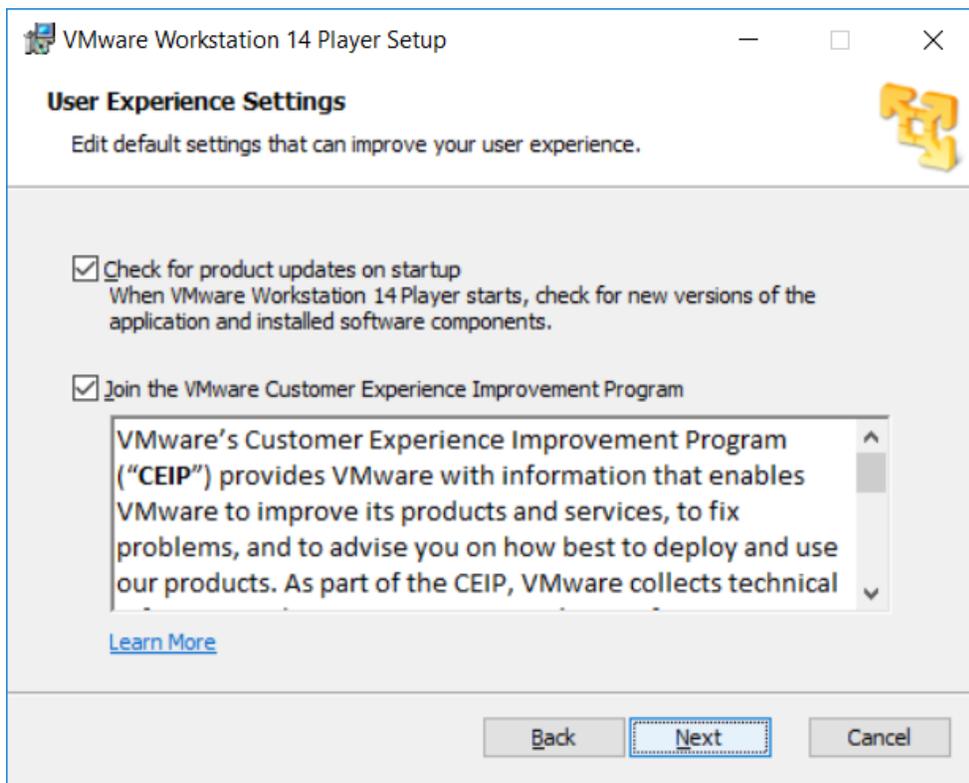

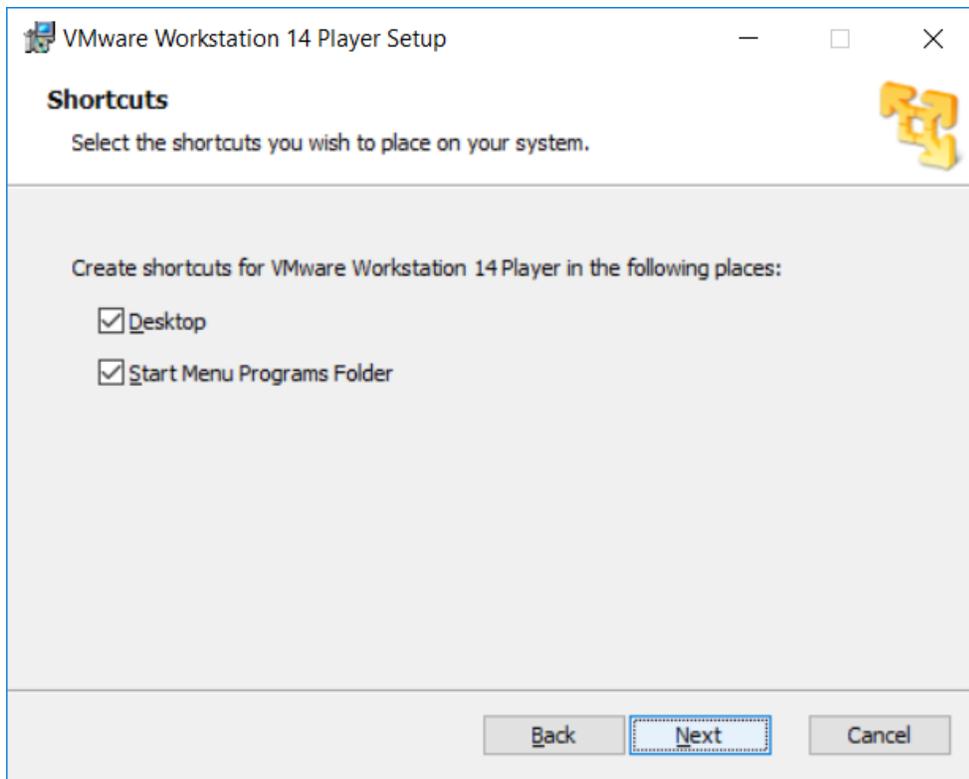
3. Accept the License agreement by clicking on the option.

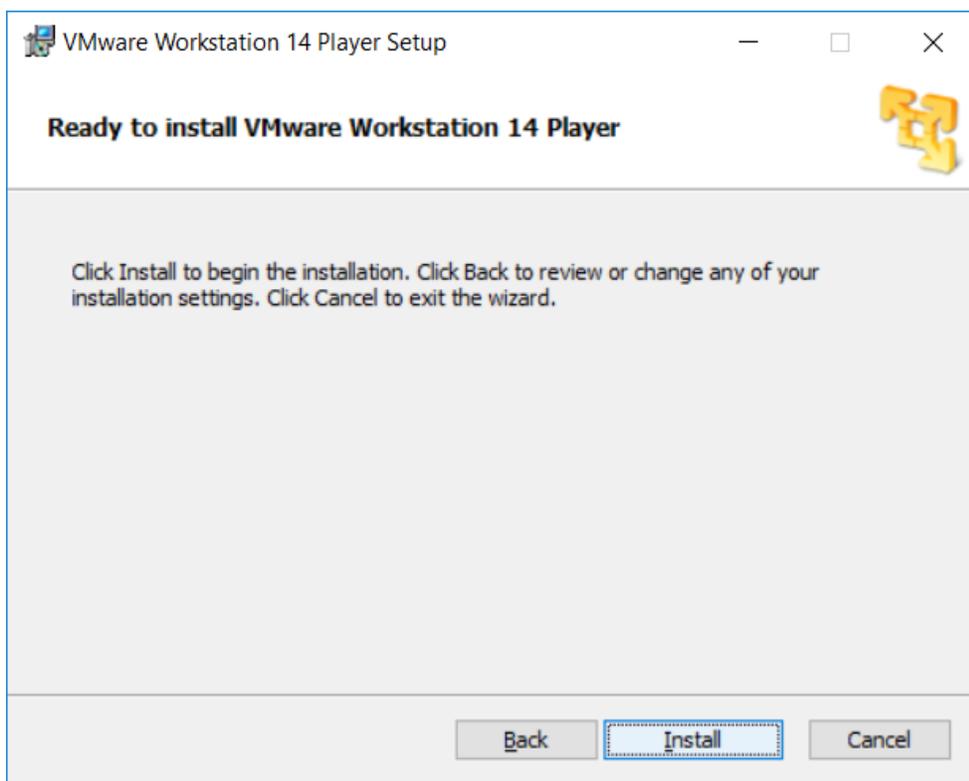4. Select folder to install and Click on "Ok"

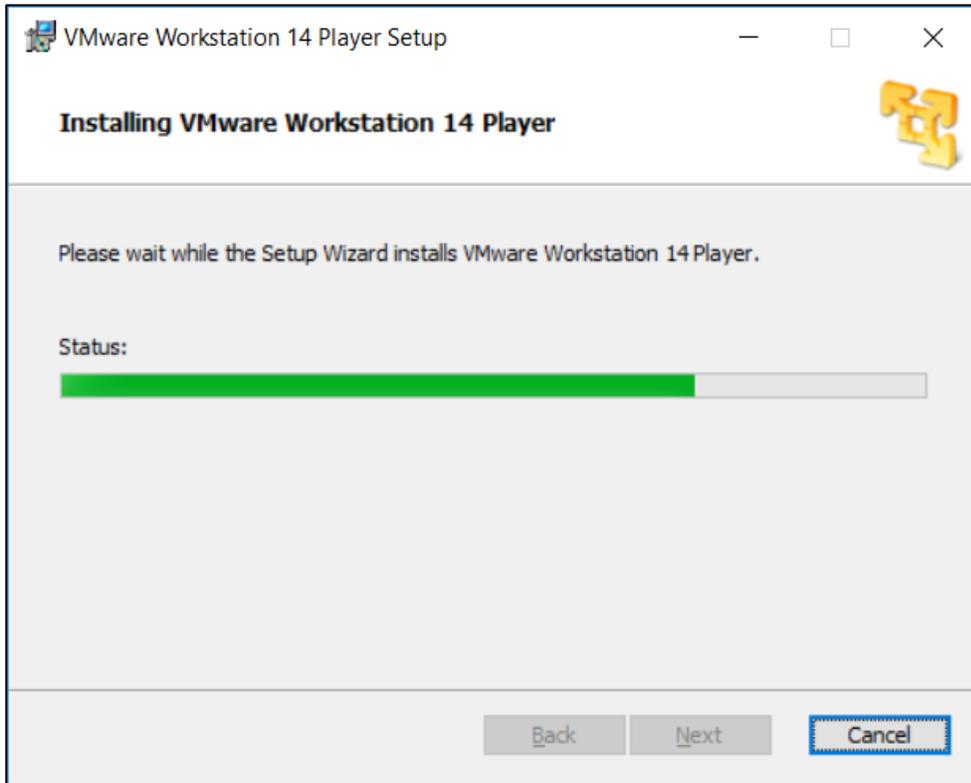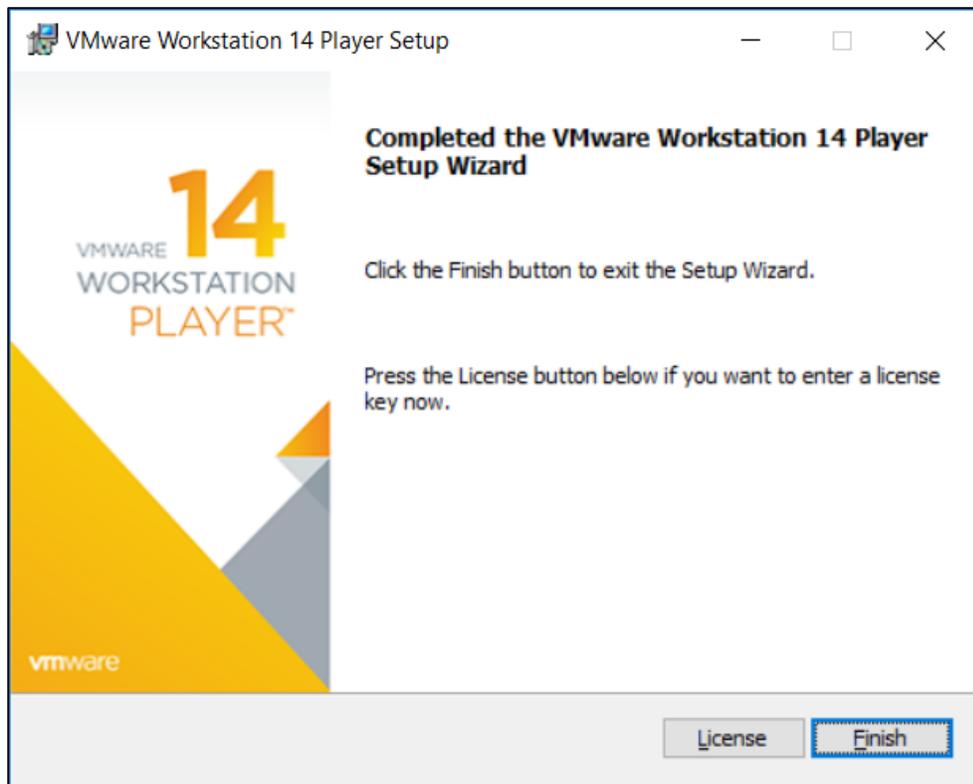5. Continue to follow the instructions.



5. Select "Next".

6. Click "Install".

7. Wait for the progress window.
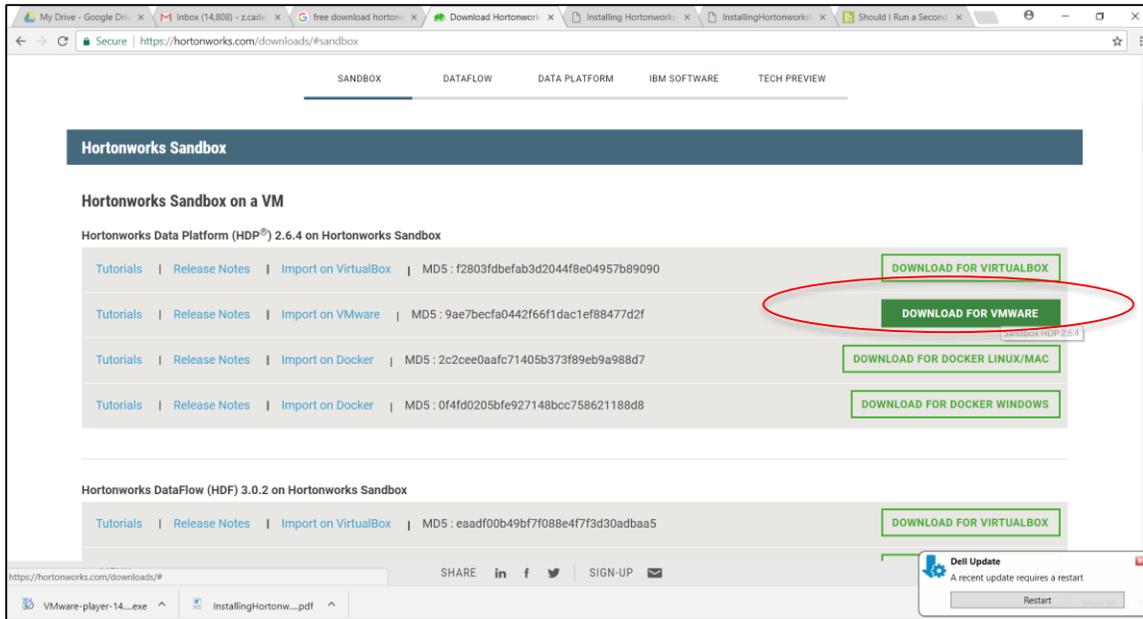


8. Click "Finish".

Your VM should now be ready and you can proceed with the next setup to install your Hadoop environment.

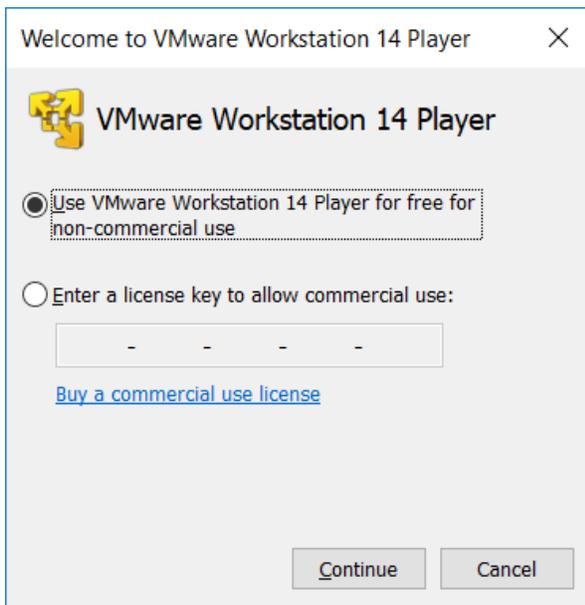## 3. Installing the Hadoop Environment

Hortonworks Sandbox has been chosen as the Big Data Platform. Follow the following steps to get Hortonworks ready for VMWare.
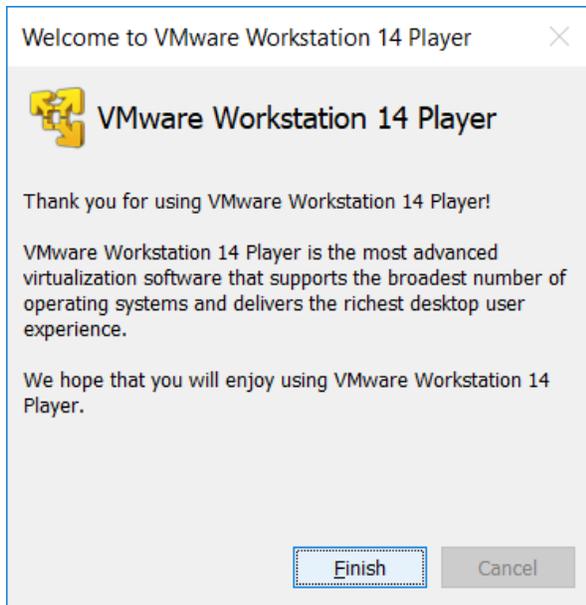
**Steps**

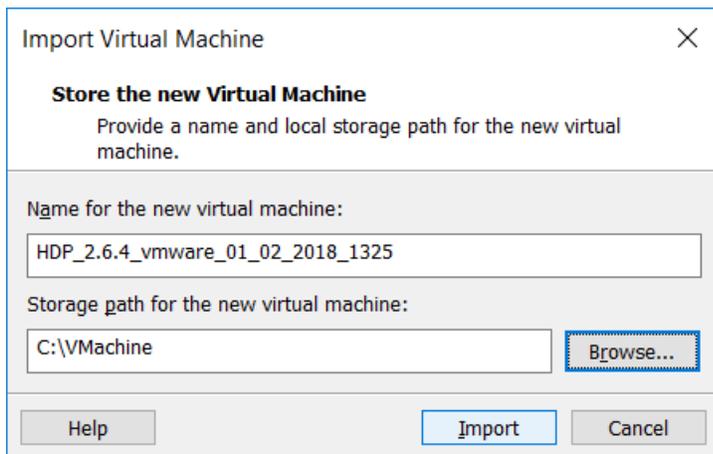1. Download Hortonworks Sandbox for VMWare.

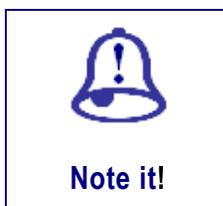2. Choose for free and click on "Continue".



3. Click "Finish".

4. Choose the storage location for the VM and leave the defaults settings (unless you have specific needs). Click on "Import".
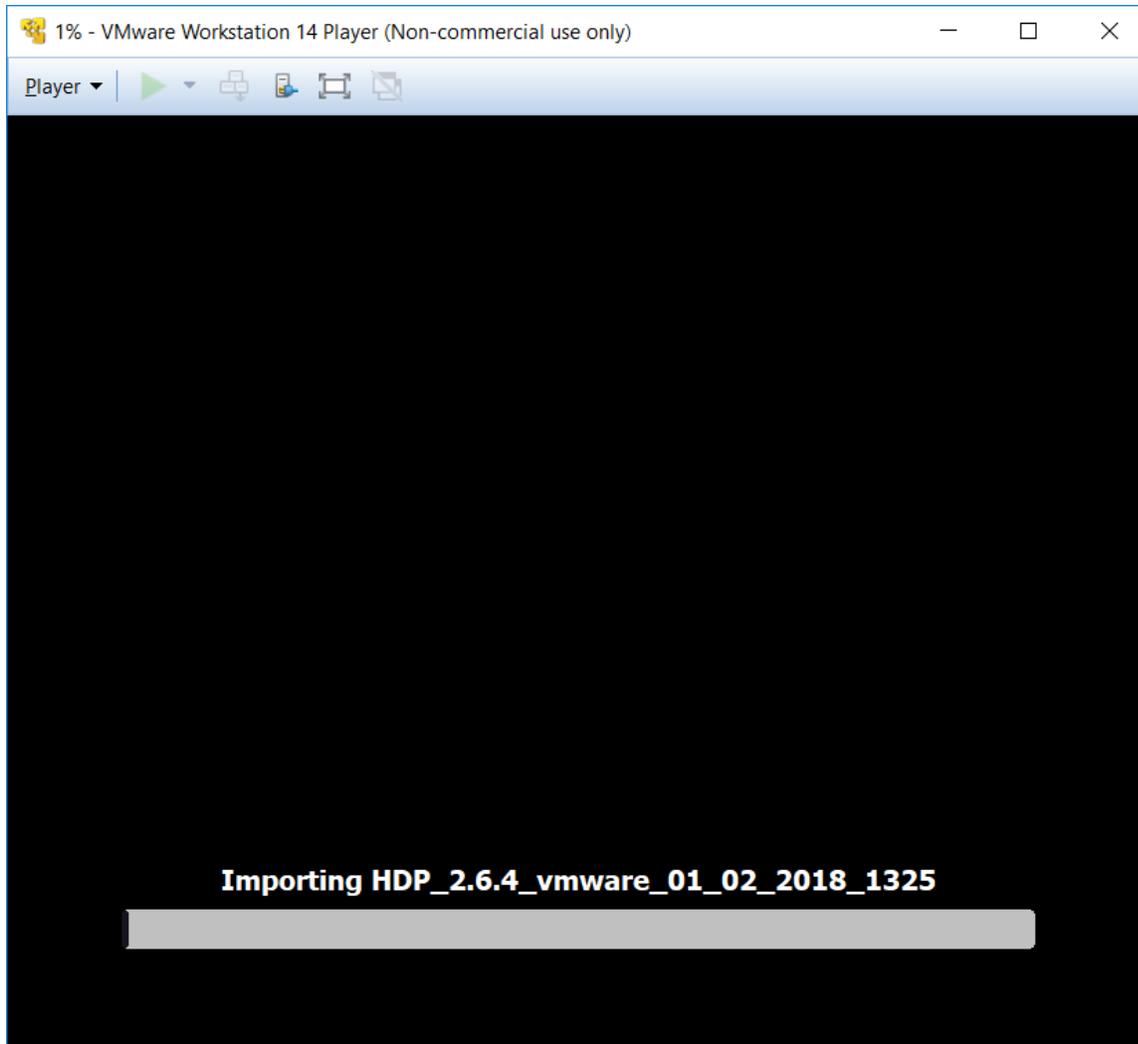


5. You will see that Hortonworks being imported automatically in the virtual machine.



**Note it!**

- *In this lab, HDP2..6.4 for VMWare has been used.*

6. You will see the progress being shown and the packages being installed

6. This might take some time. Wait for the progress

7. The following screen will appear after installation. Note the IP which is displayed which will have to be opened from your host computer.

**http://192.168.223.108:8888**



Note it!

- *You shall be accessing the VM from your host computer by using this IP.*

8. From the VM, press Alt+F5 to log into your virtual machine. The user name and password is root and hadoop respectively

9. From the VM, go to player > manage > virtual machine setting and press the tab "Hardware".

10. The following screen will appear. Click on Network adapter. Check the Network connection option. You can keep to NAT or host only depending on your requirements. For this setup, NAT was selected.

11. You can check the other settings. Select Memory as shown below. You will see that 8 GB memory has automatically been allocated for this VM as well as a portion of the hard disk.

12. Go to your host computer and type the IP saved from Step 7 above.

**http://192.168.223.108:8888**

13. The following screen should appear and your Big Data platform should be ready to use.

*You might not have the same IP and port number. In case you are not able*

11. Click on launch dashboard and select "Ambari". The following user name and password have to be used.
raj_ops (same username and password)

# 4. Installing R/RStudio on Hadoop

To be able to run R on the Big Data platform, the related R packages and IDE have to be installed. The following should be setup in your Big Data platform:

1. R
2. R studio

## 4.1 Installing R and RStudio on the VM



- *For each installation/download, you should check whether the*

<table>
<tr><td>**Note it!**</td><td>*installation or download is successful. Otherwise, you have to look for other versions depending on the OS/version/virtual machine that you are using.*<br>• *For this VM, linux centos version 5 has been used.*<br>• *When installing select "y" whenever prompted to do so to continue with the installation*</td></tr>
</table>

## Steps

1. Start by installing R. Login to your virtual machine (if not already login). We should first enable the epel repository using the following command. Type the following command and press enter.

```
yum install epel-release -y
```

2. Install R by typing the following and press enter. R version 3.4.4 has been used in this setup.

```
yum install R -y
```

3. Install RServer IDE for RStudio. For this, you should first download RStudio and then install.

4. Use command wget to download RStudio.

<table>
<tr><td><br>**Note it!**</td><td>• *If you use wget and you get the error message "command not found", you should first install wget using yum install wget*</td></tr>
</table>

```
yum install wget
```

5. Download RStudio using

```
wget https://download2.rstudio.org/rstudio-server-rhel-1.0.44-x86_64.rpm
```

6. Install RStudio after successful download using

```
yum install --nogpgcheck rstudio-server-rhel-1.0.44-x86_64.rpm
```

7. Once the package has been get and installed, it will be automatically starting the service. To check this use
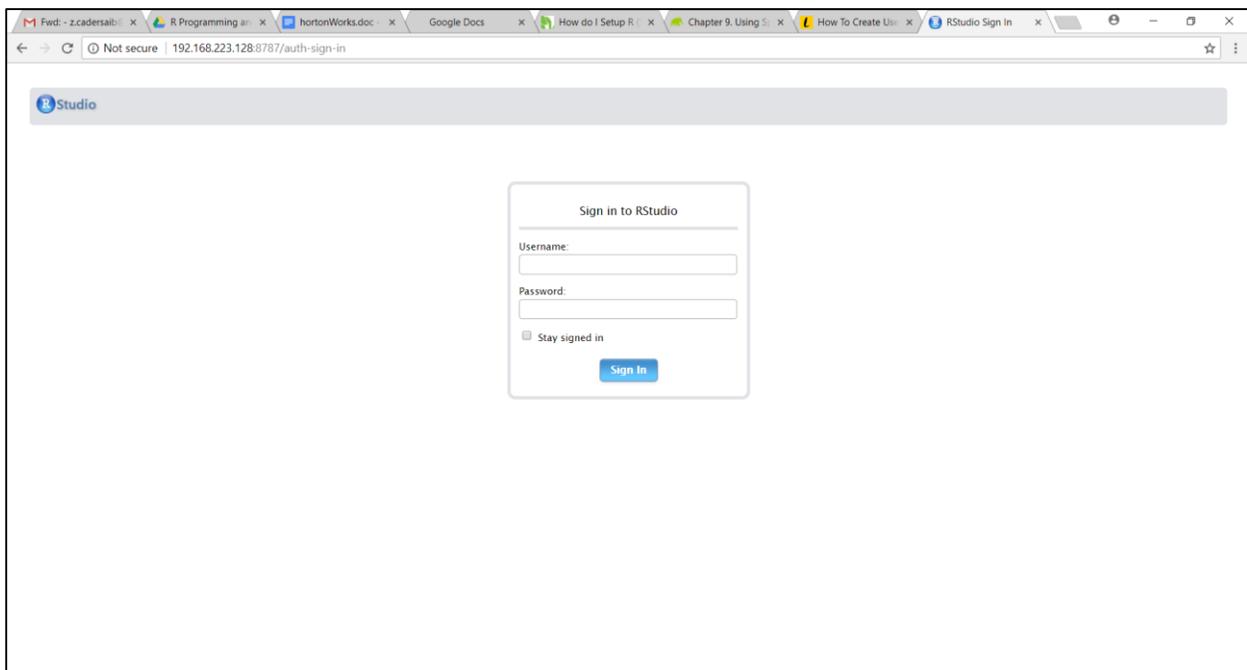
```
systemctl status rstudio-server.service
```



- *You should have the message Active (running).*

**Note it!**

8. Now you can open Rstudio from the browser of your host computer using port 8787 by typing the following IP

**http://192.168.223.128:8787**

9. The following window should be opened prompting you to enter the user name and password.
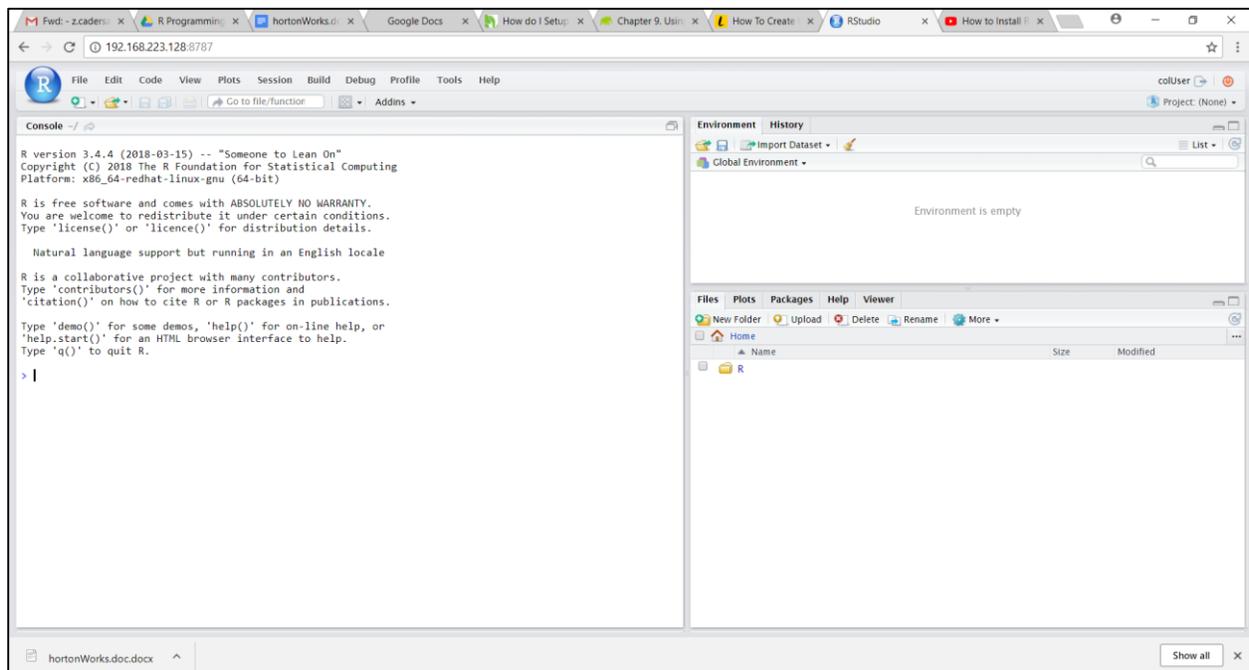


10. You should create **another user than root**. See instructions below for creating a new user from the VM command line. Go back to the VM to create a new user. Creating a new user, **colUser (or any other user)** using the command

```
useradd -m colUser
```

11. Set the password of the new user created (in this case colUser) using

```
passwd colUser
```

12. Back in the browser of your host computer, use the new user (colUser) created to login and RStudio should be opened successfully as below



13. Go back to the terminal and install the developers' tool using

```
yum groupinstall "Development Tools"
```

14. After successful installation of the development packages, install packages in RStudio by going back to browser where RStudio is opened and install the following packages.

```
>install.packages('quantmod')

>install.packages('txtplot')
```

|   Note it! | • *Refer to chapter 1 to know how to install packages in RStudio.* |
|---|---|