**CSC**

# CSC Deutschland GmbH

# Deploying a Data Stack
# for Predictive Maintenance

**An instruction of HDP 2.4 deployment with Python & TensorFlow**
**for vibration data analysis in Jupyter Notebook**

Autor: Zilong Zhao
✉: zzhao3@csc.com
June 7, 2016

# Contents

# 1 Introduction

Hortonworks Data Platform[1] is an open source Apache Hadoop distribution. Officially Hortonworks offers a Sandbox on their website, which is an interactive Hadoop ecosystem, running in virtual machine. However, the integrated system is CentOS 6, which causes error during the installation of TensorFlow 0.8.0 due to the missing **GNU C Library** `glibc 2.16`.

Hence, in this introduction, we deploy the HDP 2.4 on the latest CentOS Linux distribution, namely Version 7. Furthermore, the Jupyter Notebook will be used as our data science IDE.

*Keywords:* CentOS 7, HDP 2.4, Spark, TensorFlow, Python, Jupyter Notebook

# 2 Deploying the HDP

We start our deployment from installation of Cent OS 7 on.

## 2.1 Create a CentOS 7 VM

The installation of Cent OS 7 should be done in couple of minutes. Meanwhile, we set two user accounts. One is `hdp` as administrator with password `csc`, and `root` with password `hearts4csc`. Applying the CSC proxy settings `web-gate.de.emea.csc.com:80` in Applications > System Tools > Settings > Network > Network proxy, we get some text output as above by typing

```
1   curl www.google.de
```

Also, in Applications > System Tools > Settings > Details > Device name, we set it as `wiesbaden`. First of all, we will install the Anaconda 2.5.0 from `https://repo.continuum.io/archive/index.html`
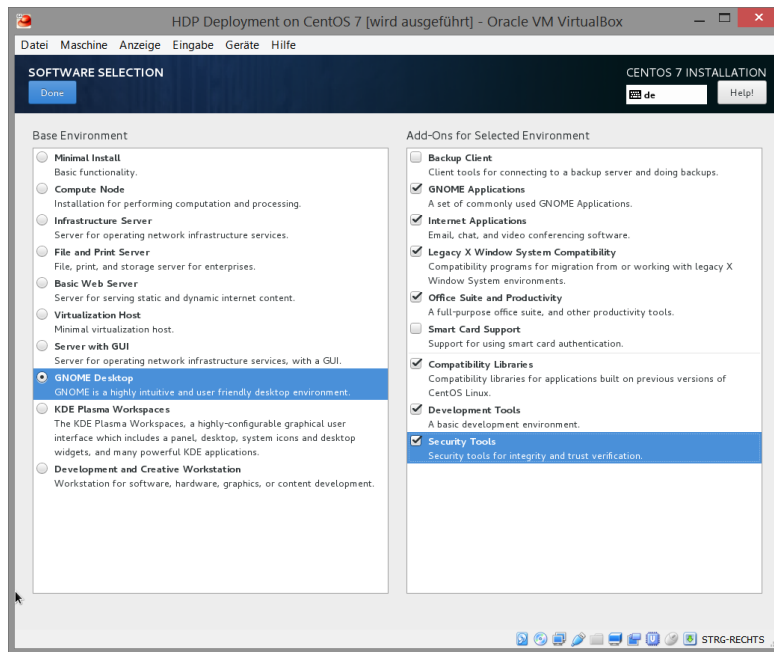
---

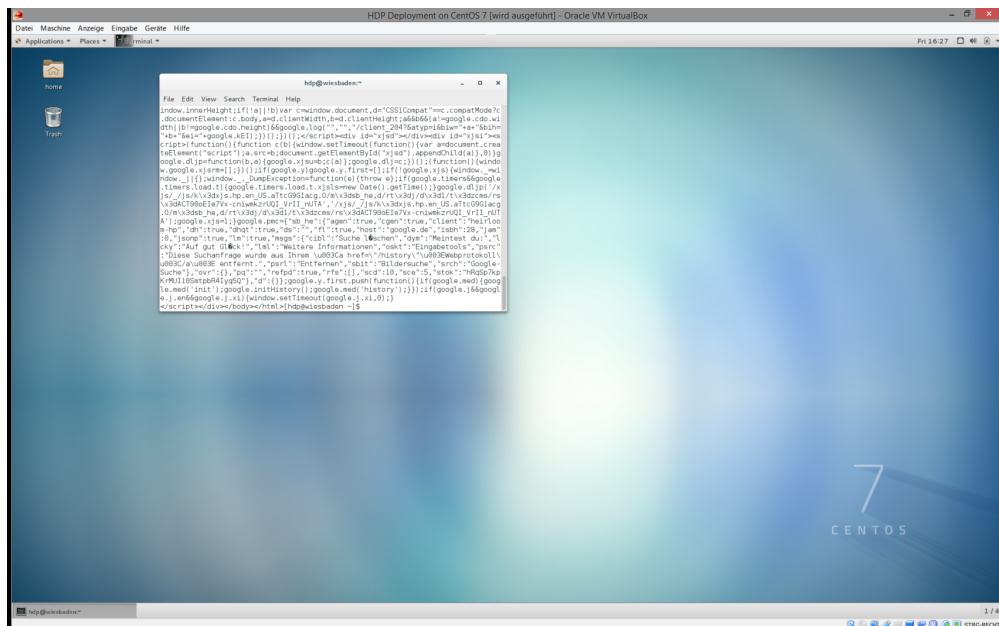[1]abbr. HDP

Figure 1: The Cent OS 7 package selection



Figure 2: The Cent OS 7 with gnome interface

## 2.2 Anaconda Python & TensorFlow

Anaconda is an open source distribution of the Python programming languages for large-scale data processing, predictive analytics, and scientific computing, that aims to simplify package management and deployment. The advantage is we have not only the `pip` package management pipeline, it gives us also a `conda` package management system to connect with the Anaconda cloud service.

```
1  [hdp@wiesbaden ~]$ cd ~/Downloads/
2  [hdp@wiesbaden Downloads]$ wget https://repo.continuum.io/archive/Anaconda2-2.5.0-Linux-x86_64.sh
```

To makes the file executable by everyone, we type

```
1  [hdp@wiesbaden Downloads]$ chmod +x Anaconda2-2.5.0-Linux-x86_64.sh
```

Now by running

```
1  [hdp@wiesbaden Downloads]$ ./Anaconda2-2.5.0-Linux-x86_64.sh
```

we start to install Anaconda Python. During the installation we accept the default path setting for anaconda2 and for `bash` we have

```
1  Python 2.7.11 :: Continuum Analytics, Inc.
2  creating default environment...
3  installation finished.
4  Do you wish the installer to prepend the Anaconda2 install location
5  to PATH in your /home/hdp/.bashrc ? [yes|no]
6  [no] >>> yes
```

After reopening a new terminal, we install Google TensorFlow with

```
1  [root@wiesbaden hdp]# pip install -U --ignore-installed setuptools
2  [root@wiesbaden hdp]# pip install --upgrade
   ↪    https://storage.googleapis.com/tensorflow/linux/cpu/tensorflow-0.9.0rc0-cp27-none-linux_x86_64.whl
```

## 2.3 Configure networking

Now our system is ready for Hortonworks Data Platform deployment. However, in the deploying process, we need a proxy free network.

We create two network adapters, "Adapter 1" set to "Host-only Adapter" and "Adapter 2" using "NAT" (If you haven't already created a "Host-only Network", you should do that first, under the VirtualBox Preferences Network tab).

Log in as root. In a "Basic Server" install the nework adapters aren't enabled by default so we need to do this. Modify the /etc/sysconfig/network-scripts/ifcfg-enp0s3 and /etc/sysconfig/network-scripts/ifcfg-enp0s8 files. Change the boot setting to ONBOOT=yes for both files so the adapters are enabled when we boot. When this is done just enter the reboot command. Once you log back in you should see both adapters available when entering the ifconfig command. Make note of the IP address of the enp0s8 interface, we will need that to access the VM from the host machine and we will also use that network for our Ambari configuration.

```
1  [root@wiesbaden hdp]# hostname
2  wiesbaden
3  [root@wiesbaden hdp]# ifconfig
4  ...
5  enp0s8: flags=4163<UP,BROADCAST,RUNNING,MULTICAST>  mtu 1500
6          inet 192.168.56.102  netmask 255.255.255.0  broadcast 192.168.56.255
7          inet6 fe80::a00:27ff:fe9a:7044  prefixlen 64  scopeid 0x20<link>
8  ...
```

We'll add enp0s8 IP address and hostname to /etc/hosts:

```
1  127.0.0.1    localhost localhost.localdomain localhost4 localhost4.localdomain4
2  ::1          localhost localhost.localdomain localhost6 localhost6.localdomain6
3
4  192.168.56.102   wiesbaden
```

## 2.4 Verify SSH settings

Now we need to verify SSH settings to check that the sshd service is running and enabled on boot:

```
1  [root@wiesbaden ~]# service sshd status
2  [root@wiesbaden ~]# chkconfig --list sshd
```

If it is not then run these commands:

```
1  [root@wiesbaden ~]#  chkconfig sshd on
2  [root@wiesbaden ~]#  service sshd start
```

## 2.5 Connect with SSH

Now we can log out and SSH to the server from a terminal on the host machine.

```
1  [hdp@wiesbaden ~]$ ssh root@192.168.56.102
```

As we log in as root through SSH the first time, we should see

```
1  [hdp@wiesbaden ~]$ ssh root@192.168.56.102
2  The authenticity of host '192.168.56.102 (192.168.56.102)' can't be established.
3  ECDSA key fingerprint is e3:92:52:2a:54:5b:8e:0b:b1:61:9a:7b:6f:b8:72:a8.
4  Are you sure you want to continue connecting (yes/no)? yes
5  Warning: Permanently added '192.168.56.102' (ECDSA) to the list of known hosts.
6  root@192.168.56.102's password:
7  Last login: Tue Jun  7 10:45:49 2016
8  ABRT has detected 1 problem(s). For more info run: abrt-cli list --since 1465289149
9  [root@wiesbaden ~]#
```

## 2.6 Update OpenSSL

Next step we will prepare for Ambari installation. Ambari is a cluster management for Hortonworks Data Platform. We update the installed OpenSSL package using:

```
1  [root@wiesbaden ~]#  yum -y update openssl
```

## 2.7 Set up Password-less SSH

It is a good idea to creat a password-less SSH, in order to do this we need to generate and store our SSH certificates in the .ssh/autorized_keys file on the local system. We can use ssh-keygen for this. We creat a password free SSH for the root account:

```
1   [root@wiesbaden hdp]# ssh-keygen
2   Generating public/private rsa key pair.
3   Enter file in which to save the key (/root/.ssh/id\_rsa):
4   Enter passphrase (empty for no passphrase):
5   Enter same passphrase again:
6   Your identification has been saved in /root/.ssh/id\_rsa.
7   Your public key has been saved in /root/.ssh/id\_rsa.pub.
8   The key fingerprint is:
9   28:f6:d9:2b:d7:37:3b:b9:7b:30:68:a4:77:e0:bb:1f root@wiesbaden
10  The key's randomart image is:
11  +--[ RSA 2048]----+
12  |                 |
13  |                 |
14  |                 |
15  |        .  o     |
16  |     o . S+ o    |
17  |   . o o. = +    |
18  |      o .+ oE+   |
19  |       . ..o =.. |
20  |        o. .+=B  |
21  +-----------------+
22  [root@wiesbaden hdp]# cd /root/.ssh/
23  [root@wiesbaden .ssh]# cat /root/.ssh/id_rsa.pub >> authorized_keys
```

We show the private key value with

```
1   [root@wiesbaden hdp]# cat /root/.ssh/id_rsa
```

and use the key as an option during the installation.

Next, we do the same process for hdp account.

```
1   [hdp@wiesbaden ~]$ ssh-keygen
2   Generating public/private rsa key pair.
3   Enter file in which to save the key (/home/hdp/.ssh/id_rsa):
4   Enter passphrase (empty for no passphrase):
5   Enter same passphrase again:
6   Your identification has been saved in /home/hdp/.ssh/id_rsa.
7   Your public key has been saved in /home/hdp/.ssh/id_rsa.pub.
8   The key fingerprint is:
9   1f:2f:1b:d4:de:26:f0:83:4f:13:82:f5:be:83:af:9c hdp@wiesbaden
10  The key's randomart image is:
11  +--[ RSA 2048]----+
12  |                 |
13  |                 |
14  |           .     |
15  |          o o    |
16  |         S = +   |
17  |        o X o    |
18  |         =.X o   |
19  |         ..B.*   |
20  |          Eo+.   |
21  +-----------------+
22  [hdp@wiesbaden ~]$ su
23  Password:
24  [root@wiesbaden hdp]# cat /home/hdp/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
25  [root@wiesbaden hdp]# chmod 600 ~/.ssh/authorized_keys
```

We creat a password free SSH for the root account:

```
1  [root@wiesbaden hdp]# ssh-keygen
2  Generating public/private rsa key pair.
3  Enter file in which to save the key (/root/.ssh/id\_rsa):
4  Enter passphrase (empty for no passphrase):
5  Enter same passphrase again:
6  Your identification has been saved in /root/.ssh/id\_rsa.
7  Your public key has been saved in /root/.ssh/id\_rsa.pub.
8  The key fingerprint is:
9  28:f6:d9:2b:d7:37:3b:b9:7b:30:68:a4:77:e0:bb:1f root@wiesbaden
10 The key's randomart image is:
11 +--[ RSA 2048]----+
12 |                 |
13 |                 |
14 |                 |
15 |       .  o      |
16 |    o . S+ o     |
17 |   . o o. = +    |
18 |      o .+ oE+    |
19 |      . ..o =..  |
20 |        o. .+=B   |
21 +-----------------+
22 [root@wiesbaden hdp]# cd /root/.ssh/
23 [root@wiesbaden .ssh]# cat /root/.ssh/id_rsa.pub >> authorized_keys
```

We show the private key value with

```
1  [root@wiesbaden hdp]# cat /root/.ssh/id_rsa
```

and use the key as an option during the installation.

We should now be able to create an SSH connection to wiesbaden without being prompted for a password (if you get a prompt for adding the host to the known hosts file, just enter yes):

```
1  [hdp@wiesbaden ~]$ ssh root@192.168.56.102
2  Last login: Mon Jun  6 15:41:18 2016 from wiesbaden
3  [root@wiesbaden ~]# exit
4  logout
5  Connection to 192.168.56.102 closed.
```

## 2.8 Enable ntpd

Next, Ambari needs time to be synchronized so we need to enable ntpd.

```
1  [root@wiesbaden hdp]# chkconfig ntpd on
2  Note: Forwarding request to 'systemctl enable ntpd.service'.
3  Created symlink from /etc/systemd/system/multi-user.target.wants/ntpd.service to
        ↪    /usr/lib/systemd/system/ntpd.service.
4  [root@wiesbaden hdp]# ntpdate pool.ntp.org
5   6 Jun 15:44:25 ntpdate[7998]: adjust time server 131.188.3.221 offset 0.208360 sec
6  [root@wiesbaden hdp]# service ntpd start
7  Redirecting to /bin/systemctl start  ntpd.service
```

## 2.9  Disable iptables

Then we disable iptables. A number of network ports need to be open on the VM for Ambari during setup. The easiest way to open the ports is to disable the iptables process:

```
1    # service iptables stop
2    # chkconfig iptables off
```

## 2.10  Disable SELinux and Check umask Value

Now we disable SELinux and Check umask Value. To disable SELinux during the Ambari setup:

```
1    # setenforce 0
```

To permanently disable SELinux for the VM modify /etc/sysconfig/selinux and change the config to SELINUX=disabled. Also, make sure umask is set to 0022 (it should be for a new install)

```
1    # umask
2    0022
```

## 2.11  Disable IPv6

Then we need to disable IPv6, log in as root and cut-and-paste the following commands into your terminal window to disable IPv6:

```
1    mkdir -p /etc/sysctl.d
2    ( cat > /etc/sysctl.d/99-hadoop-ipv6.conf <<-'EOF'
3    ## Disabled ipv6
4    ## Provided by Ambari Bootstrap
5    net.ipv6.conf.all.disable_ipv6 = 1
6    net.ipv6.conf.default.disable_ipv6 = 1
7    net.ipv6.conf.lo.disable_ipv6 = 1
8    EOF
9        )
10   sysctl -e -p /etc/sysctl.d/99-hadoop-ipv6.conf
```

## 2.12 Disable Transparent Huge Pages (THP)

When installing Ambari, one or more host checks may fail if you have not disabled Transparent Huge Pages on all hosts.

To disable THP log in as root and add the following commands to your /etc/rc.local file:

```
1  if test -f /sys/kernel/mm/redhat_transparent_hugepage/defrag;
2    then echo never > /sys/kernel/mm/redhat_transparent_hugepage/defrag
3  fi
4  if test -f /sys/kernel/mm/redhat_transparent_hugepage/enabled;
5    then echo never > /sys/kernel/mm/redhat_transparent_hugepage/enabled
6  fi
```

To confirm, reboot the host and then run the command:

```
1  [root@wiesbaden hdp]# cat /sys/kernel/mm/transparent_hugepage/enabled
2  [always] madvise never
```

## 2.13 Install httpd

We need to have the web server running so log in as root and install it with the following commands:

```
1  [root@wiesbaden hdp]# yum -y install httpd
```

We'll set the ServerName to be ${hostname}:80 in /etc/httpd/conf/httpd.conf

```
1  # ServerName gives the name and port that the server uses to identify itself.
2  # This can often be determined automatically, but we recommend you specify
3  # it explicitly to prevent problems during startup.
4  #
5  # If your host doesn't have a registered DNS name, enter its IP address here.
6  #
7  ServerName ${hostname}:80
```

Now we can start the httpd server.

```
1  chkconfig httpd on
2  service httpd start
```

Congratulations! We finally move to step for Ambari installation.

## 2.14 Set up Ambari `yum` repo

Download the Ambari repository file to a directory on your installation host.

```
1  [root@wiesbaden hdp]# wget -nv
     ↪    http://public-repo-1.hortonworks.com/ambari/centos7/2.x/updates/2.2.2.0/ambari.repo -O
     ↪    /etc/yum.repos.d/ambari.repo
```

Confirm that the repository is configured by checking the repo list.

```
1   [root@wiesbaden hdp]# yum repolist
2   Loaded plugins: fastestmirror, langpacks
3   Updates-ambari-2.2.2.0                                | 2.9 kB    00:00
4   Updates-ambari-2.2.2.0/primary_db                    | 6.3 kB    00:00
5   Loading mirror speeds from cached hostfile
6    * base: wftp.tu-chemnitz.de
7    * extras: ftp.plusline.de
8    * updates: centos.schlundtech.de
9   repo id                          repo name                          status
10  Updates-ambari-2.2.2.0           ambari-2.2.2.0 - Updates                8
11  base/7/x86_64                    CentOS-7 - Base                     9,007
12  extras/7/x86_64                  CentOS-7 - Extras                     305
13  updates/7/x86_64                 CentOS-7 - Updates                  1,683
14  repolist: 11,003
```

## 2.15 Install Ambari server packages

Install the Ambari packages. This also installs the default PostgreSQL Ambari database. Enter y when prompted to to confirm transaction and dependency checks.

```
1   [root@wiesbaden hdp]# yum install ambari-server
2   ...
3     Installing : postgresql-libs-9.2.15-1.el7_2.x86_64                   1/4
4     Installing : postgresql-9.2.15-1.el7_2.x86_64                        2/4
5     Installing : postgresql-server-9.2.15-1.el7_2.x86_64                 3/4
6     Installing : ambari-server-2.2.2.0-460.x86_64                        4/4
7     Verifying  : postgresql-libs-9.2.15-1.el7_2.x86_64                   1/4
8     Verifying  : postgresql-server-9.2.15-1.el7_2.x86_64                 2/4
9     Verifying  : ambari-server-2.2.2.0-460.x86_64                        3/4
10    Verifying  : postgresql-9.2.15-1.el7_2.x86_64                        4/4
11
12  Installed:
13    ambari-server.x86_64 0:2.2.2.0-460
14
15  Dependency Installed:
16    postgresql.x86_64 0:9.2.15-1.el7_2
17    postgresql-libs.x86_64 0:9.2.15-1.el7_2
18    postgresql-server.x86_64 0:9.2.15-1.el7_2
19
20  Complete!
```

Note: Accept the warning about trusting the Hortonworks GPG Key. That key will be automatically downloaded and used to validate packages from Hortonworks. You will see the following message:

```
1   Retrieving key from http://public-repo-1.hortonworks.com/ambari/centos7/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
2   Importing GPG key 0x07513CAD:
3    Userid     : "Jenkins (HDP Builds) <jenkin@hortonworks.com>"
4    Fingerprint: df52 ed4f 7a3a 5882 c099 4c66 b973 3a7a 0751 3cad
5    From       : http://public-repo-1.hortonworks.com/ambari/centos7/RPM-GPG-KEY/RPM-GPG-KEY-Jenkins
6   Is this ok [y/N]: y
```

## 2.16 Configure and start the Ambari server

Before starting the Ambari Server, you must set it up. Setup configures Ambari to talk to the Ambari database, installs the JDK and allows you to customize the user account the Ambari Server daemon will run as. The ambari-server setup command manages the setup process. Run the following command on the Ambari server host to start the setup process. You may also append Setup Options to the command.

```
1   [root@wiesbaden hdp]# ambari-server setup
2   Using python  /usr/bin/python
3   Setup ambari-server
4   Checking SELinux...
5   SELinux status is 'enabled'
6   SELinux mode is 'permissive'
7   WARNING: SELinux is set to 'permissive' mode and temporarily disabled.
8   OK to continue [y/n] (y)? y
9   Customize user account for ambari-server daemon [y/n] (n)? n
10  Adjusting ambari-server permissions and ownership...
11  Checking firewall status...
12  Redirecting to /bin/systemctl status  iptables.service
13
14  Checking JDK...
15  [1] Oracle JDK 1.8 + Java Cryptography Extension (JCE) Policy Files 8
16  [2] Oracle JDK 1.7 + Java Cryptography Extension (JCE) Policy Files 7
17  [3] Custom JDK
18  ==============================================================================
19  Enter choice (1): 1
20  To download the Oracle JDK and the Java Cryptography Extension (JCE) Policy Files you must accept the
     ↪    license terms found at http://www.oracle.com/technetwork/java/javase/terms/license/index.html and
     ↪    not accepting will cancel the Ambari Server setup and you must install the JDK and JCE files
     ↪    manually.
21  Do you accept the Oracle Binary Code License Agreement [y/n] (y)? y
22  Downloading JDK from http://public-repo-1.hortonworks.com/ARTIFACTS/jdk-8u60-linux-x64.tar.gz to
     ↪    /var/lib/ambari-server/resources/jdk-8u60-linux-x64.tar.gz
23  jdk-8u60-linux-x64.tar.gz... 100% (172.8 MB of 172.8 MB)
24  Successfully downloaded JDK distribution to /var/lib/ambari-server/resources/jdk-8u60-linux-x64.tar.gz
25  Installing JDK to /usr/jdk64
26  Successfully installed JDK to /usr/jdk64/
27  Downloading JCE Policy archive from http://public-repo-1.hortonworks.com/ARTIFACTS/jce_policy-8.zip to
     ↪    /var/lib/ambari-server/resources/jce_policy-8.zip
28
29  Successfully downloaded JCE Policy archive to /var/lib/ambari-server/resources/jce_policy-8.zip
30  Installing JCE policy...
31  Completing setup...
32  Configuring database...
33  Enter advanced database configuration [y/n] (n)? n
34  Configuring database...
35  Default properties detected. Using built-in database.
36  Configuring ambari database...
37  Checking PostgreSQL...
38  Running initdb: This may take upto a minute.
39  Initializing database ... OK
40
41
42  About to start PostgreSQL
43  Configuring local database...
44  Connecting to local database...done.
45  Configuring PostgreSQL...
46  Restarting PostgreSQL
47  Extracting system views...
48  ambari-admin-2.2.2.0.460.jar
49  ......
50  Adjusting ambari-server permissions and ownership...
51  Ambari Server 'setup' completed successfully.
```

Now, start the Ambari server:

```
1   [root@wiesbaden hdp]# ambari-server start
```

We can now connect to Ambari server on `localhost:8080`, username is `admin` and password is `admin`. Now it's the time to create our own Hadoop Cluster. From the Ambari Dashboard click the Launch Install Wizard link. We

configured our single-node VMs hostname as the only host. We name our cluster
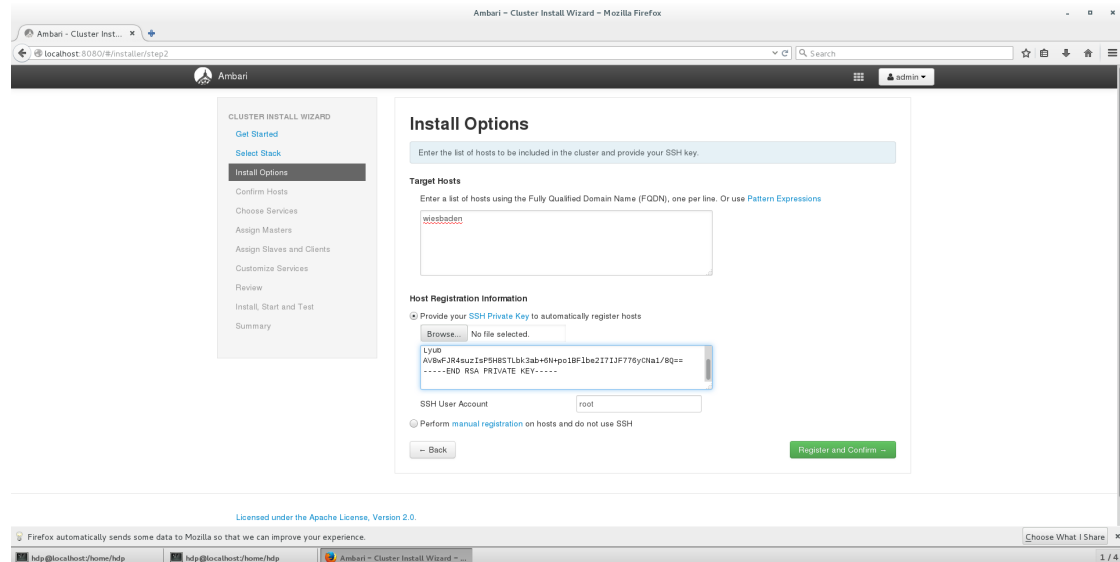as `xip`



Figure 3: Install Options

With the correct setting in Fig. 3, we get the Fig. 4

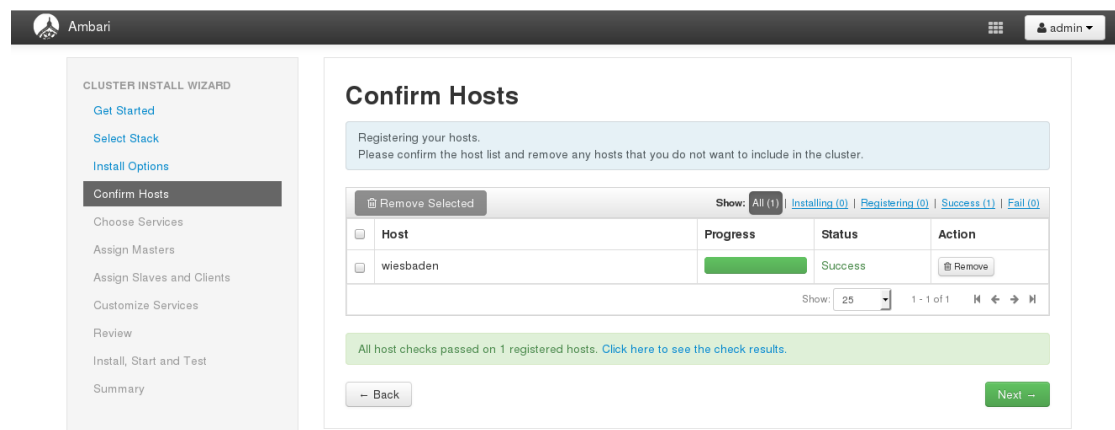Figure 4: Confirm Hosts

After choosing the services we want to install on our cluster, we assign the masters all with wiesbaden due to the single node. With several times Next clicks, we are finally able to deploy our single node cluster.
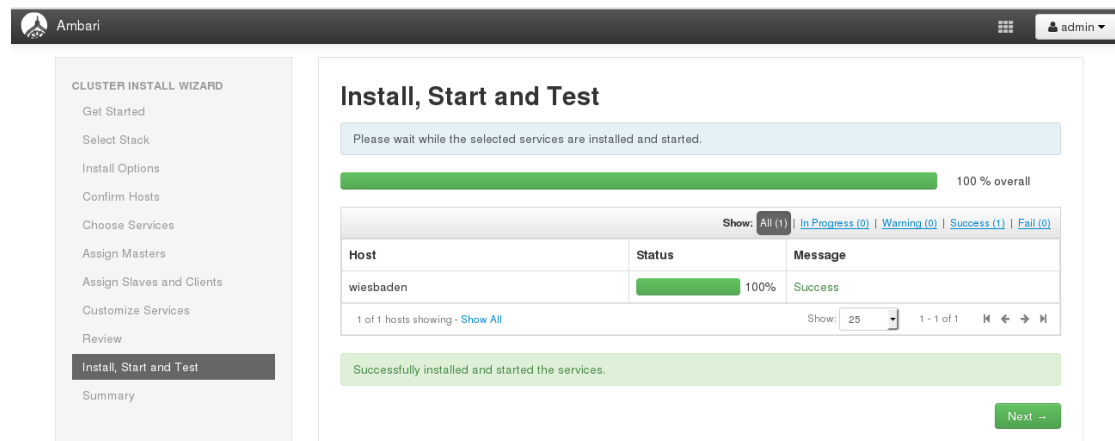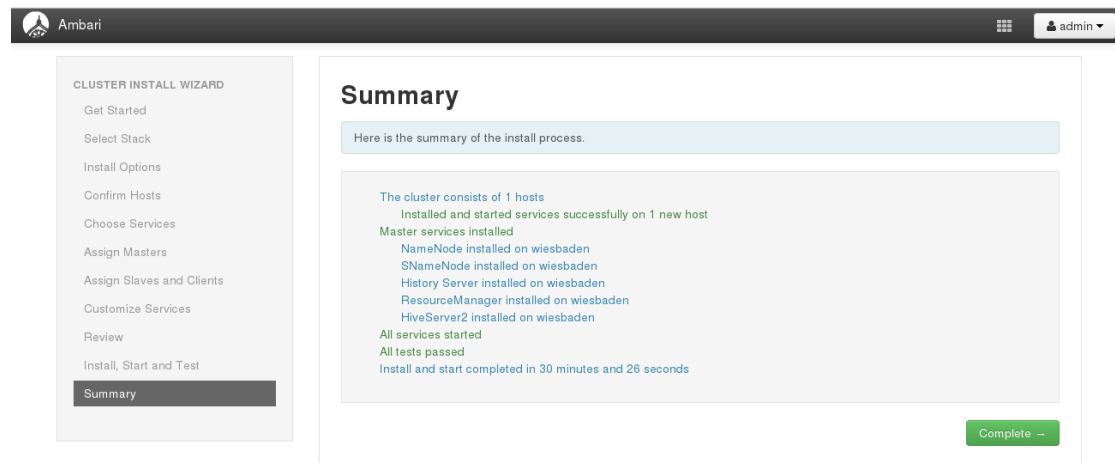


Figure 5: Install, Start and Test

Figure 6: Summary

With a successful install process we get a summary: all services started and all tests passed. The following screenshot shows how the Ambari dashboard looks like:
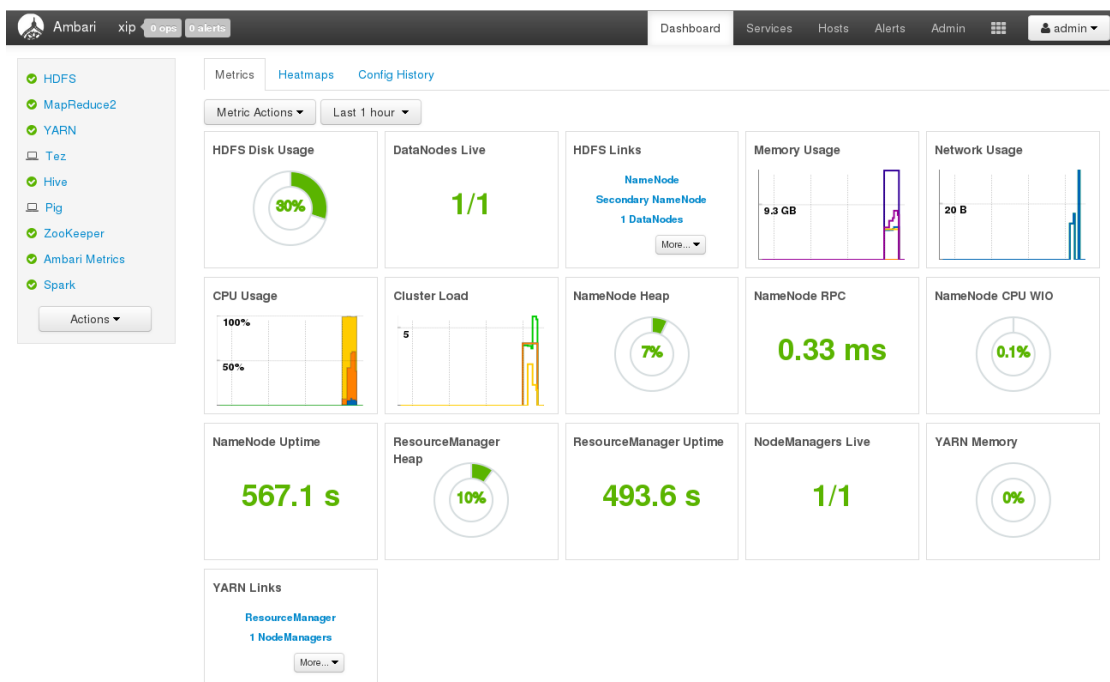
Figure 7: The Ambari Dashboard

## 2.17 Configuring Jupyter Notebook

Since we want to use IPython with Apache Spark we have to use the Python interpreter which is built with Apache Spark, `pyspark`, instead of the default Python interpreter. We begin with upgrading the Jupyter modules

```
1  [root@wiesbaden hdp]# pip install --upgrade jupyter-client
2  [root@wiesbaden hdp]# pip install --upgrade jupyter-console
3  [root@wiesbaden hdp]# pip install --upgrade jupyter-core
```

Let's create a Jupyter Notebook profile for pyspark

```
1  [root@wiesbaden hdp]# ipython profile create pyspark
```

Next generate a jupyter config file:

```
1  [root@wiesbaden hdp]# jupyter notebook --generate-config
```

You should see the following output:

```
1  Writing default config to: /root/.jupyter/jupyter_notebook_config.py
```

Now open your preferred editor. I'm using vi to edit jupyter_notebook_config.py. Let's change some settings as following:

```
1  # The directory to use for notebooks and kernels.
2  c.NotebookApp.notebook_dir = u'/usr/hdp/current/spark-client/'
3
4  # The port the notebook server will listen on.
5  c.NotebookApp.port = 8889
6
7  # The IP address the notebook server will listen on.
8  c.NotebookApp.ip = '127.0.0.1'
```

Next we are going to create a shell script to set the appropriate values every time we want to start Jupyter Notebook with pyspark enviroment. We create a shell script with the following command:

```
1  [root@wiesbaden hdp]# vi ~/start_ipython_notebook.sh
```

Then copy the following lines into the file:

```
1  #!/bin/bash
2  IPYTHON_OPTS="notebook" pyspark
```

Running

```
1  [root@wiesbaden hdp]# chmod +x ~/start_ipython_notebook.sh
2  [root@wiesbaden hdp]# ~/start_ipython_notebook.sh
```

A browser will be automatically opened on your host machine and the URl http://127.0.0.1:8889 shows up and you should see the screen below:
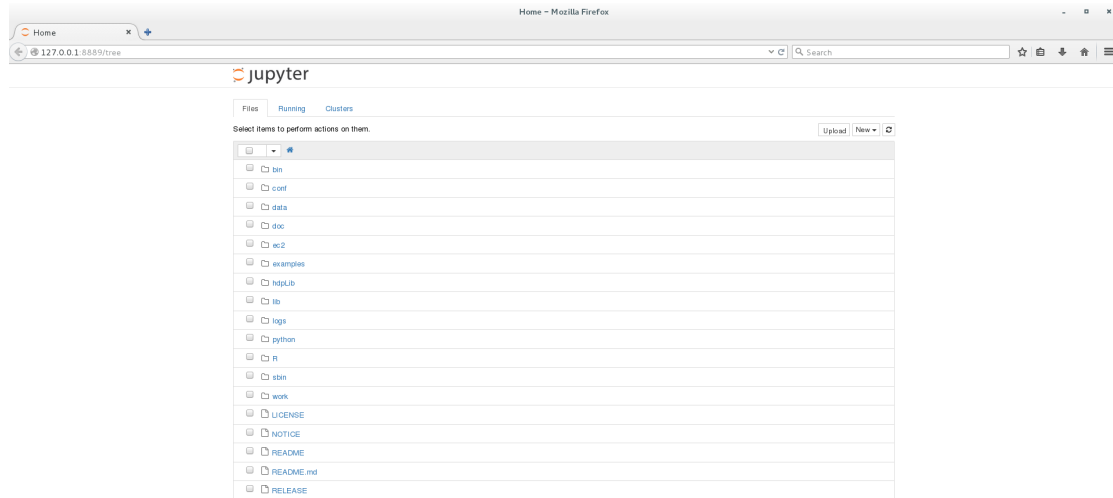
Figure 8: Jupyter Notebook

Voila! We have just configured jupyter notebook with Apache Spark on the Hortonworks Data Platform.

# 3 Summary

In this document, we have seen how to deploy the Hortonworks Data Platform on a single node cluster for vibration data analysis. In our future works, we would like to scale out the cluster for more intensive computing and storage tasks.