# Ambari Server

Using the Horton Works Virtual Machine, we downloaded [HortonWorks Sandbox](#) from the official web site, but I couldn't work in the Virtual Box in our OS (not enough Required Memory).

From this problem, we started using the Google Cloud by deploying the ambary server on a Virtual Machine with Centos6 (8 vCPU, 52 Go RAM), so to do that we created an account that give us free 300$ credits that we can use for the compute engine.

First, we started by setting up password-less SSH connections between the Ambari Server host and all other hosts in the cluster. To do that we generated a SSH public and private Keys with the RSA encryption, we copy the keys generated to the target host, we added the public key to the authorized keys, we set permission for the authorized keys.

```
ssh-keygen -t rsa
.ssh/id_rsa
.ssh/id_rsa.pub
cat id_rsa.pub >> authorized_keys
chmod 700 ~/.ssh
chmod 600 ~/.ssh/authorized_keys
```

Second, we will synchronize the clock of the cluster and the machine that runs the browser that we will see later to access to the ambari server.

```
service ntpd start
chkconfig ntpd on
```

After that, we add the host of the server by adding the local IP address and the name of domain to configure the DNS protocol, to do that we edited the **vi /etc/hosts** by adding the domain that is defined for our virtual machine.

```
10.156.0.2 hortonworks.c.cloud-service-190416.internal hortonworks
```

- **hortonworks.c.cloud-service-190416.internal :** name of the project in the google cloud
- **hortonworks :** Name of the chosen domain

By adding this domain, we need to configure the network of this host, we edit the **vi /etc/sysconfig/network** by modify the HOSTNAME to the domain chosen.

So, to make the communication between the hosts and the Ambari Server, we must open the ports that we use to get to the browser, so we disabled the **iptables** and **ip6tables**

```
service iptables stop

service ip6tables stop

chkconfig iptables off

chkconfig ip6tables off
```

We disabled the SELinux and Umask,

```
vi /etc/selinux/config

SELINUX=disabled

umask 0022
```

We used MySQL database for the Hive operation in the Ambari Server, so we check for the connector.

After the configuration that we made below, we downloaded the public repository Ambari Server and install it.

```
wget -nv http://public-
repo1.hortonworks.com/ambari/centos6/2.x/updates/2.6.1.0/ambari.repo -O
/etc/yum.repos.d/ambari.repo

yum install ambari-server

ambari-server setup

ambari-server start
```
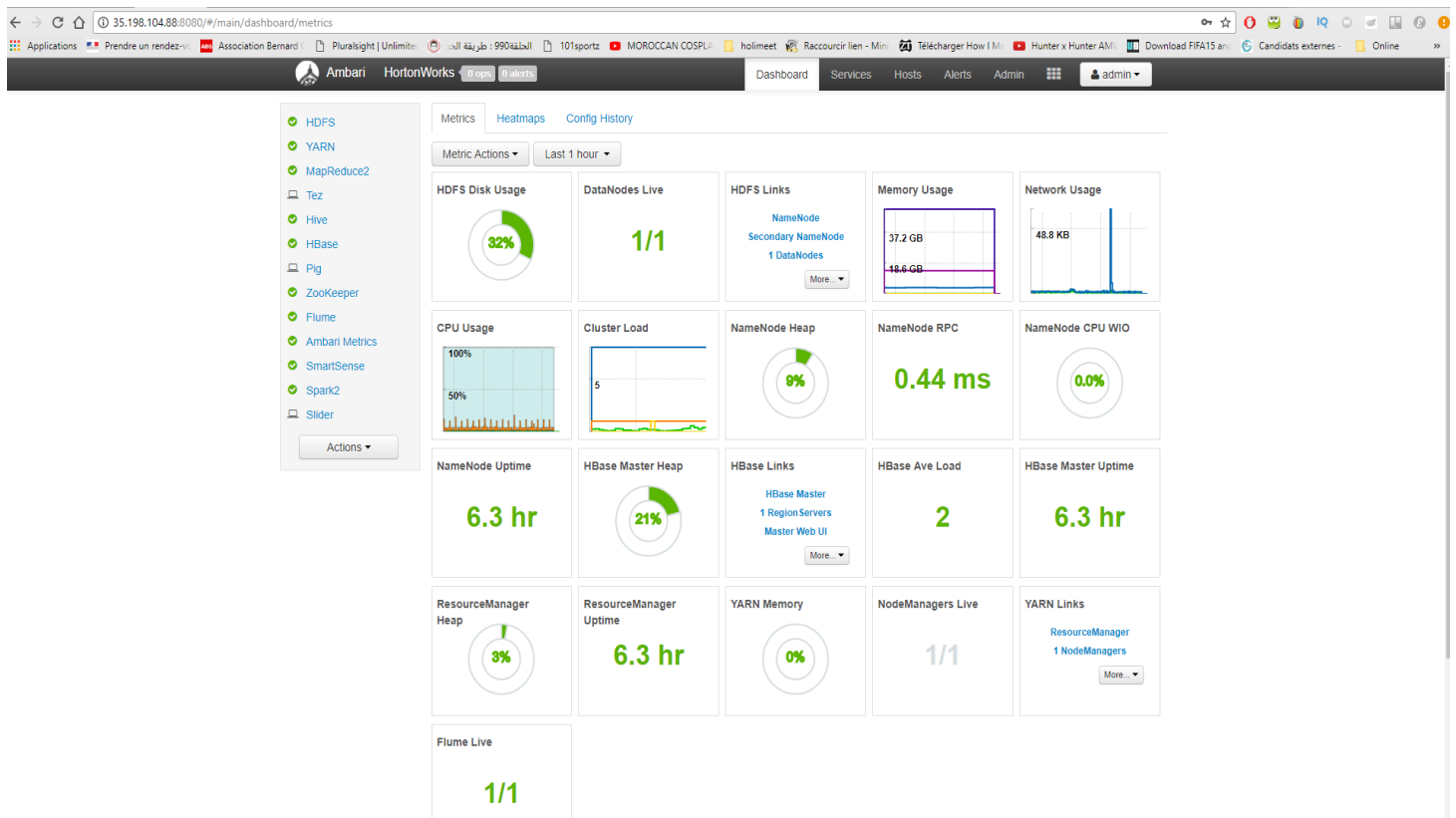


We press the SSH key.

After the installation and the running, the server is running in the external address (it changes while we restart the virtual machine) in the 8080 port,

We can log in with the **"admin" "admin"**, and right after, we configured the server by using the Private key that we generated. After, we added the services that we may use, (HDFS, MapReduce, Spark2, ZooKeeper, Hive, etc. …)



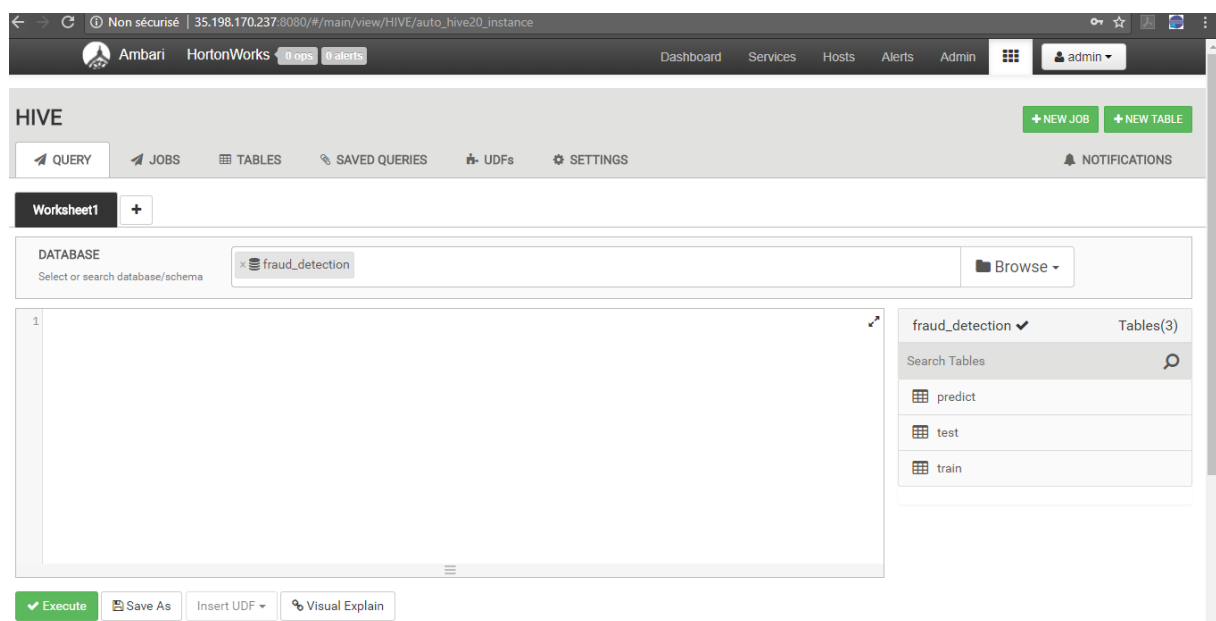So, the dashboard looks like this:

We will let the virtual machine running in the cloud for the next day, you can access to it by this link

http:// 35.198.170.237 :8080

Coming to Hive. It provides us data warehousing facilities on top of an existing Hadoop cluster. Along with that it provides an SQL like interface which makes your work easier, in case you are coming from an SQL background. You can create tables in Hive and store data there. Along with that you can even map your existing HBase tables to Hive and operate on them.

We click on New Table:



We drag our file .csv



Here we find our tables.



Selection request :

Results:



| predict.operationtype | predict.amount | predict.sourcename | predict.beforebalancesource | predict.afterbalancesource | predict.destinationname |
|---|---|---|---|---|---|
| TRANSFER | 1397689.48 | C1224276273 | 1397689.48 | 0.0 | C298592695 |
| TRANSFER | 290211.26 | C276602564 | 290211.26 | 0.0 | C58143708 |
| TRANSFER | 178344.25 | C1280852753 | 178344.25 | 0.0 | C461887553 |
| CASH_OUT | 1023189.9 | C1057733607 | 1023189.9 | 0.0 | C385233047 |
| TRANSFER | 592615.53 | C366677925 | 592615.53 | 0.0 | C1433310823 |
| PAYMENT | 2632.63 | C183500200 | 0.0 | 0.0 | M1765700212 |
| CASH_IN | 60402.09 | C113413524 | 3130888.26 | 3191290.35 | C105834049 |
| CASH_OUT | 152097.84 | C1978501596 | 0.0 | 0.0 | C1937322225 |
| CASH_IN | 251107.86 | C24724881 | 398483.0 | 649590.86 | C964775275 |
| CASH_OUT | 75734.99 | C109260257 | 75734.99 | 0.0 | C2020337583 |
| TRANSFER | 493.28 | C1723063411 | 493.28 | 0.0 | C1565674952 |
| PAYMENT | 34185.61 | C1512170440 | 7253.0 | 0.0 | M420143361 |

We also create a VM in Google cloud to execute python script.

```
saidou_khadija@python1:~$ sudo pip3 install --upgrade ipython numpy pandas matplotlib scipy \
>                          scikit-learn jupyter
The directory '/home/saidou_khadija/.cache/pip/http' or its parent directory is not owned by the curr
he cache has been disabled. Please check the permissions and owner of that directory. If executing pi
ou may want sudo's -H flag.
The directory '/home/saidou_khadija/.cache/pip' or its parent directory is not owned by the current u
g wheels has been disabled. check the permissions and owner of that directory. If executing pip with
want sudo's -H flag.
Collecting ipython
  Downloading ipython-6.2.1-py3-none-any.whl (745kB)
    100% |                                | 747kB 1.6MB/s
Requirement already up-to-date: numpy in /usr/local/lib/python3.5/dist-packages
Collecting pandas
  Downloading pandas-0.22.0-cp35-cp35m-manylinux1_x86_64.whl (25.7MB)
    100% |                                | 25.7MB 51kB/s
Collecting matplotlib
  Downloading matplotlib-2.1.2-cp35-cp35m-manylinux1_x86_64.whl (15.0MB)
    100% |                                | 15.0MB 97kB/s
Collecting scipy
  Downloading scipy-1.0.0-cp35-cp35m-manylinux1_x86_64.whl (49.6MB)
    90% |                           | 44.9MB 87.3MB/s eta 0:00:01
```

To send the predict.csv , test.csv, and train csv, we used FileZela:

```
saidou_khadija@python1:~$ python3 humrqApG.py
Training accuracy: 1.0
Test accuracy: 0.9341167645955745
16.630875804
saidou_khadija@python1:~$ ls
humrqApG.py  outScript.csv  predict.csv  test.csv  train.csv
saidou_khadija@python1:~$
```

To connect to the machine with SSH we generated a key with puttygen then we copied this key to th ssh configuration machine:

**Clés SSH**

☐ Bloquer les clés SSH à l'échelle du projet

**rsa-key-20180216**

ssh-rsa
AAAAB3NzaC1yc2EAAAABJQAAAQEAu3Bf.../s9GC2z2GzihpPe3Hn4tw== rsa-
key-20180216

Then we importe the ppk in putty: