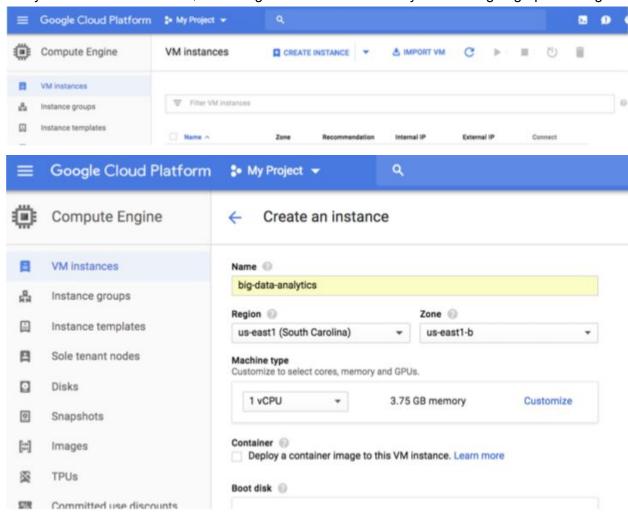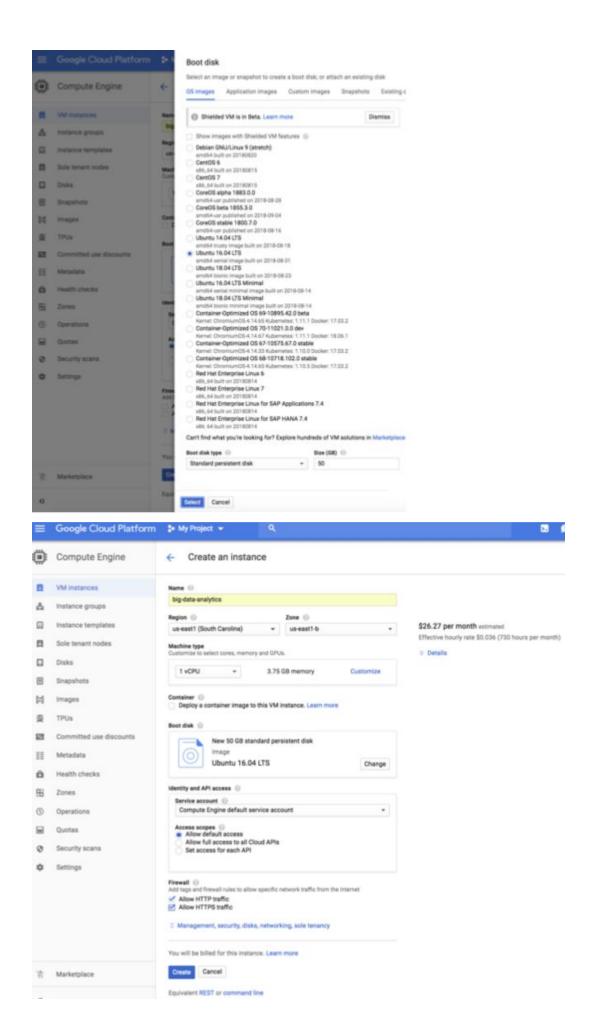# Tutorial for Beginners

## 1. Google Cloud Setup

Head over to https://console.cloud.google.com/
Everyone could use the free $300 Google credits valid for one year from signing up on Google Cloud.

# 2. Installations

## 2.1 Install JDK, JRE

```
sudo apt-get update
sudo apt-get install openjdk-8-jdk-headless default-jre ssh rsync readlink -f
/usr/bin/java | sed "s:bin/java::"
```

OUTPUT for last command: /usr/lib/jvm/java-8-openjdk-amd64/jre/

## 2.2 Install Hadoop

Firstly, you need to find out your username using

```
whoami
```

Then, download it using

```
wget http://www.eu.apache.org/dist/hadoop/common/stable/hadoop-2.9.2.tar.gz tar -xvzf
hadoop-2.9.2.tar.gz
mv hadoop-2.9.2 hadoop
```

(There may be more advanced version, like 2.9.3, so you would need to verify it by yourself.)
Then, edit the configurations using

```
vim ./.bashrc
```

Put these at the end of the file (To insert, type "i" for insert.)

```
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64 export PATH=$PATH:$HADOOP_PREFIX/bin
```

Now save the file with <esc> key followed by typing ":wq" for write-quit

```
vim ./hadoop/etc/hadoop/hadoop-env.sh
```

Then, comment the line "export JAVA_HOME=${JAVA_HOME}" and replace with these lines:

```
#export JAVA_HOME=${JAVA_HOME}
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
```

```
vim ./hadoop/etc/hadoop/core-site.xml
```

Put these lines inside

```xml
<property>
<name>fs.default.name</name>
<value>hdfs://localhost:1234</value>
</property>
```

```
vim ./hadoop/etc/hadoop/hdfs-site.xml
```

Put these lines inside

```xml
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
```

```
cp ./hadoop/etc/hadoop/mapred-site.xml.template ./hadoop/etc/hadoop/mapred-site.xml
vim ./hadoop/etc/hadoop/mapred-site.xml
```

Now Put these lines inside

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

```
vim ./hadoop/etc/hadoop/yarn-site.xml
```

Put these lines inside

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

Next line has two single quotes (take caution on typesetting).

```
ssh-keygen -t rsa -P '' -f ./.ssh/id_rsa
cat ./.ssh/id_rsa.pub >> ./.ssh/authorized_keys
```

Execute "ssh localhost" and enter "yes" when prompted. And then enter "exit" Followed by this, execute:

```
source ./.bashrc
```

Start Hadoop (first time execution)
Remember to execute these commands in order (even if they may seem duplicated commands):

```
cd ./hadoop
./sbin/start-dfs.sh
./bin/hdfs namenode -format
./sbin/stop-dfs.sh
./sbin/start-dfs.sh
./bin/hdfs dfs -mkdir /user
./bin/hdfs dfs -mkdir /user/your_username
./sbin/start-yarn.sh
jps
cd ../
```

After jps, it should show like this:



Sample Hadoop Examples

```
./hadoop/bin/hadoop jar
./hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.9.1.jar pi 10 100
```

Note, above is one line code.


## 2.3 Setup HBase

Downloads and Edit Configuration

```
wget http://www-eu.apache.org/dist/hbase/stable/hbase-1.4.7-bin.tar.gz tar -zxvf
hbase-1.4.7-bin.tar.gz
mv hbase-1.4.7 hbase
vim ./.bashrc
```

Add to bashrc file:

```
export HBASE_HOME=/home/your_username/hbase
export PATH=$PATH:$HBASE_HOME/bin
```

Save the file and execute this:

```
source ./.bashrc
```

```
vim ./hbase/conf/hbase-env.sh
```

Add to file:

```
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
```

```
vim ./hbase/conf/hbase-site.xml
```

```xml
<property>
<name>hbase.rootdir</name>
<value>hdfs://localhost:1234/hbase</value>
</property>
<property>
<name>hbase.zookeeper.property.dataDir</name>
<value>/home/your_username/zookeeper</value>
</property>
<property>
<name>hbase.cluster.distributed</name>
<value>true</value>
</property>
```

Start HBase

```
start-hbase.sh ./hadoop/bin/hdfs dfs -ls /hbase
```

If HBase setup is successful, a folder /hbase would be created and will be listed with the last command.

## 2.4 Install Hive

Downloads and Edit Configuration

```
wget https://archive.apache.org/dist/hive/stable/apache-hive-1.2.2-bin.tar.gz
tar -xvzf apache-hive-1.2.2-bin.tar.gz
mv apache-hive-1.2.2-bin hive
```

```
vim ./.bashrc
```

Add to bashrc file:

```
export HIVE_HOME=/home/your_username/hive export PATH=$PATH:$HIVE_HOME/bin
```

Save the file and execute this:

```
source ./.bashrc
```

```
cp ./hive/conf/hive-env.sh.template ./hive/conf/hive-env.sh
vim ./hive/conf/hive-env.sh
```

Add to file:

```
HADOOP_HOME=/home/your_username/hadoop
```

Start Hive

```
hive
exit;
```

## 2.5 Install spark

Downloads and Edit Configuration

```
wget https://archive.apache.org/dist/spark/spark-2.3.1/spark-2.3.1-bin-hadoop2.7.tgz
tar -xvzf spark-2.3.1-bin-hadoop2.7.tgz
mv spark-2.3.1-bin-hadoop2.7 spark
vim ./.bashrc
```

```
export SPARK_HOME=/home/your_username/spark
export PATH=$SPARK_HOME/bin:$PATH
```

```
source ./.bashrc
```

Start Spark

```
spark-shell
:quit
```

## 2.6 Install Jupyter notebook

Downloads and Execution

```
wget https://repo.continuum.io/archive/Anaconda3-5.1.0-Linux-x86_64.sh
bash Anaconda3-5.1.0-Linux-x86_64.sh
```

Accept terms: yes
Install folder: <ENTER>
PATH in .bashrc: yes
VSCode: no

```
source ./.bashrc
cd ./anaconda3/bin/
jupyter notebook --generate-config
cd ..
vim ./.jupyter/jupyter_notebook_config.py
```

Add to the beginning of file:

```
c = get_config()
c.NotebookApp.ip = '*'
c.NotebookApp.open_browser = False
c.NotebookApp.port = 5000
```

```
source ./.bashrc
tmux new -s jupyter
jupyter notebook
```

Before launching jupyter notebook, you should open the port on google cloud using step 3.
Then, copy the url shown and replace "localhost" with the IP on google cloud landing page and open this in browser.
To come out of a "tmux" session, use: Ctrl and B together, and then leaving the two keys, press D,
i.e. Ctrl+B, D.
To go back into the session, use: "tmux a -t jupyter"

# 3. Open Port in Google Cloud

🔲 VPC network                    ← Create a firewall rule

⊞ VPC networks

▢ External IP addresses

⊞ Firewall rules

✕ Routes

◇ VPC network peering

✉ Shared VPC

**Name** ⓘ

big-data-analytics

**Description** (Optional)

**Network** ⓘ

default ▾

**Priority** ⓘ
Priority can be 0 - 65535 Check priority of other firewall rules

1000

**Direction of traffic** ⓘ
● Ingress
○ Egress

**Action on match** ⓘ
● Allow
○ Deny

**Targets** ⓘ

All instances in the network ▾

**Source filter** ⓘ

IP ranges ▾

**Source IP ranges** ⓘ

0.0.0.0/0 ⊗

**Second source filter** ⓘ

None ▾

**Protocols and ports** ⓘ
○ Allow all
● Specified protocols and ports
  ✓ tcp :   5000
  ☐ udp :   all
  ☐ Other protocols
      protocols, comma separated, e.g. ah, sctp

⟳ Disable rule

Create    Cancel

---

🔲 VPC network      **Firewall rules**      ➕ CREATE FIREWALL RULE      ⟳ REFRESH      🗑 DELETE

⊞ VPC networks

▢ External IP addresses

⊞ Firewall rules

✕ Routes

◇ VPC network peering

✉ Shared VPC

Firewall rules control incoming or outgoing traffic to an instance. By default, incoming traffic from outside your network is blocked. Learn more
Note: App Engine firewalls are managed here.

▼ Filter resources                                            ⓘ  Columns ▾

| Name | Type | Targets | Filters | Protocols / ports | Action | Priority | Network ▲ |
|---|---|---|---|---|---|---|---|
| big-data-analytics | Ingress | Apply to all | IP ranges: 0.0.0.0/0 | tcp:5000 | Allow | 1000 | default |
| default-allow-http | Ingress | http-server | IP ranges: 0.0.0.0/0 | tcp:80 | Allow | 1000 | default |
| default-allow-https | Ingress | https-server | IP ranges: 0.0.0.0/0 | tcp:443 | Allow | 1000 | default |