

EASY BUTTON

Note: THIS WILL DELETE ALL DATA ON YOUR EBS VOLUME.

This document is put together from various sources from the w205 course materials

1.0 Install Environment

Preliminaries

1. If you have not already, create an EBS volume. Note the region in which the volume was created.
2. Start an instance of [ucbw205_complete_plus_postgres_PY2.7 \(ami-558fc730\)](#) in the **same region** as your EBS volume.

Use the below image to set up the security of the ports.

Inbound rules for sg-0abf456c (Selected security groups: sg-0abf456c)

Type ⓘ	Protocol ⓘ	Port Range ⓘ	Source ⓘ
Custom TCP Rule	TCP	4040	0.0.0.0/0
Custom TCP Rule	TCP	8080	0.0.0.0/0
Custom TCP Rule	TCP	50070	0.0.0.0/0
SSH	TCP	22	0.0.0.0/0
Custom TCP Rule	TCP	10000	0.0.0.0/0
Custom TCP Rule	TCP	8088	0.0.0.0/0

3. Attach the volume to the instance.
4. Login to the instance via ssh. At login, you will be the root user.

Setup and install

1. Determine which device is your EBS drive by running `fdisk -l`
2. The last entry is typically your EBS volume. For example `/dev/xvdf`
3. Copy the device path (`/dev/xvdf`) to your clipboard
4. Download the setup script like this:

`wget https://s3.amazonaws.com/ucbdatasciencew205/setup_ucb_complete_plus_postgres.sh`

5. Run the script like this:

```
bash setup_ucb_complete_plus_postgres.sh <paste your device path>
```

6. Follow the onscreen instructions

After the script runs

1. Hadoop will be installed and started
2. Postgres will be installed and started
3. Hive and SparkSQL will use Postgres as a metastore
4. A w205 user will exist. This is the user you should work as
5. A script called setup_zeppelin.sh has been created
 - a. If you want to setup zeppelin, type: ./setup_zeppelin.sh
 - b. You can start zeppelin by typing /data/zeppelin/bin/zeppelin.sh

You only need to go through this process **ONCE**. After the install, you should interact with your instance like this:

2.0 Set up Spark 1.5

Open <https://spark.apache.org/downloads.html> in your browser

Select a release as follows:

Spark 1.5.0

Pre-built for Hadoop 2.6 or later

Direct download

Copy the URL to download spark

As your personal user,

wget <url for spark>

tar xvzf spark-1.5.0-bin-hadoop2.6.tgz

mv spark-1.5.0-bin-hadoop2.6 spark15

export SPARK_HOME=\$HOME/spark15

export HADOOP_CONF_DIR=/etc/hadoop/conf

You can start pyspark as follows:

\$SPARK_HOME/bin/pyspark --master yarn

To get rid of pesky "INFO" lines,

Go to /data/spark15/conf

Make a copy of the _____ file

Change INFO to WARNING and save.

3.0 Integrating SparkSQL and the Hive Metastore

Integrating SparkSQL with the Hive Metastore is straightforward. However, we need to make sure that SparkSQL knows where our Hive metadata is.

Place a Hive configuration in Spark

```
mv spark15 /data

ln -s /data/spark15 $HOME/spark15

cp /data/hadoop/hive/conf/hive-site.xml /data/spark15/conf

export SPARK_HOME=/data/spark15
```

Edit /data/spark15/conf/hive-site.xml and change the following:

```
<!-- <property>

    <name>hive.metastore.uris</name>

    <value>thrift://localhost:9083</value>

    <description>IP address (or fully-qualified domain name) and port of the metastore host</description>

</property>

-->
```

To

```
<property>

    <name>hive.metastore.uris</name>

    <value>thrift://localhost:9083</value>

    <description>IP address (or fully-qualified domain name) and port of the metastore host</description>

</property>
```

Set up a Hive Metastore Service Script

In a file called /data/start_metastore.sh place the following:

```
#!/bin/bash  
  
nohup hive --service metastore &
```

In a file called /data/stop_metastore.sh place the following:

```
#!/bin/bash  
  
ps aux|grep org.apache.hadoop.hive.metastore.HiveMetaStore|awk '{print $2}'|xargs kill -9
```

From now on, when you want to use Hive Data in Spark, you must do:

```
/data/start_postgres.sh
```

```
/data/start_metastore.sh
```

Then start spark using one of the following

```
/data/spark15/bin/pyspark
```

OR

```
/data/spark15/bin/spark-sql
```

OR

```
/data/spark15/bin/spark
```

OR

```
/data/spark15/bin/spark-submit
```

4.0 Zeppelin

Configuring Zeppelin as a primary interface for Pyspark, Spark and Spark-SQL (picks up Hive metadata)

As root, do the following:

- Copy the hive-site.xml from /data/spark15/conf to /data/zeppelin/conf
- Copy the Hadoop configurations (*-site.xml) from /etc/hadoop/conf to /data/zeppelin/conf

Before starting zeppelin, **make sure your metastore is started!**

If you are going to use zeppelin, run: /data/zeppelin/bin/zeppelin.sh

Open <the URL of your EC2 instance>:8080 in your web browser

You will NOT be able to access zeppelin unless you have configured your instance security group to have port 8080 open.

5.0 START AND STOP

To connect to instance (change amazonaws.com part):

```
sudo ssh -i Angela_Gunn.pem  
root@ec2-54-165-19-196.compute-1.amazonaws.com
```

To run all the setups (as root):

```
mount -t ext4 /dev/xvdf /data  
chown hadoop:hadoop /data  
chmod ug+wx /data
```

```
./start-hadoop.sh  
cd /data  
./start_postgres.sh  
./start_metastore.sh
```

```
su - w205  
cd /data
```

To shut down the setus (as root):

```
cd /data  
  
./stop_metastore.sh  
./stop_postgres.sh  
cd $home  
./stop-hadoop.sh
```

```
umount /data
```

SPARK Commands:

```
$SPARK_HOME/bin/spark-sql
```

```
$SPARK_HOME/bin/spark-shell
```

```
$SPARK_HOME/bin/spark-submit pyspark_test.py
```

Cloning repository:

```
git clone https://anggunn:j6uKY8iL0g@github.com/anggunn/MIDS-W205.git
```

Copy files from AWS to local machine (the period at the end is important):

```
scp -i Angela_Gunn.pem  
root@ec2-52-90-52-124.compute-1.amazonaws.com:/data/MIDS-W205/205proj  
ect/query_scripts/*.png .
```