



João Torres

Follow

Feb 4, 2020 · 7 min read · Listen



Save



How to Install and Set Up an Apache Spark Cluster on Hadoop 18.04

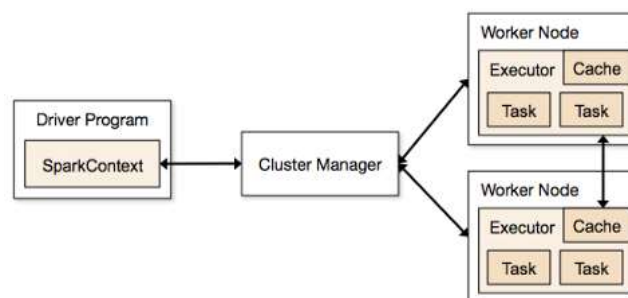


In this article I will explain how to install Apache Spark on a multi-node cluster, providing step by step instructions.

Spark Architecture

Apache Spark follows a master/slave architecture with two main daemons and a cluster manager.

- Master Daemon — (Master/Driver Process)
- Worker Daemon — (Slave Process)
- Cluster Manager



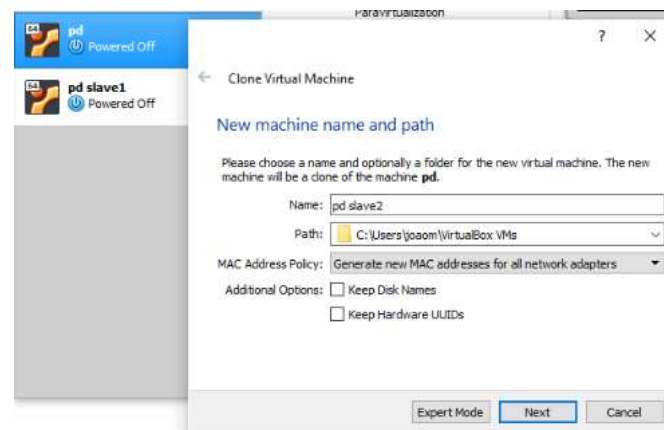
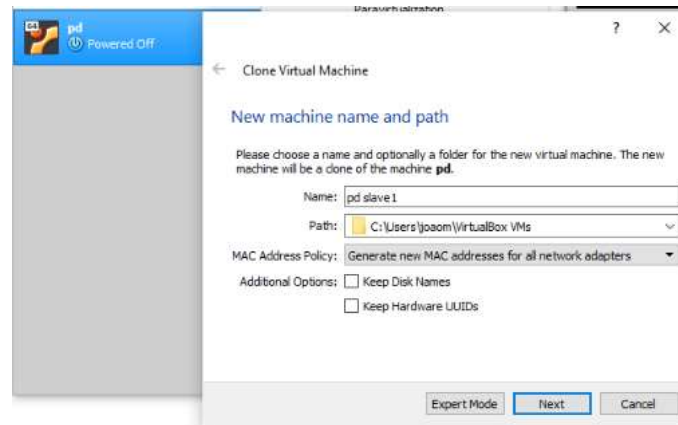
A spark cluster has a single Master and any number of Slaves/Workers. The driver and the executors run their individual Java processes and users can run them on the same horizontal spark cluster or on separate machines.

Pre-requirements

- Ubuntu 18.04 installed on a virtual machine.

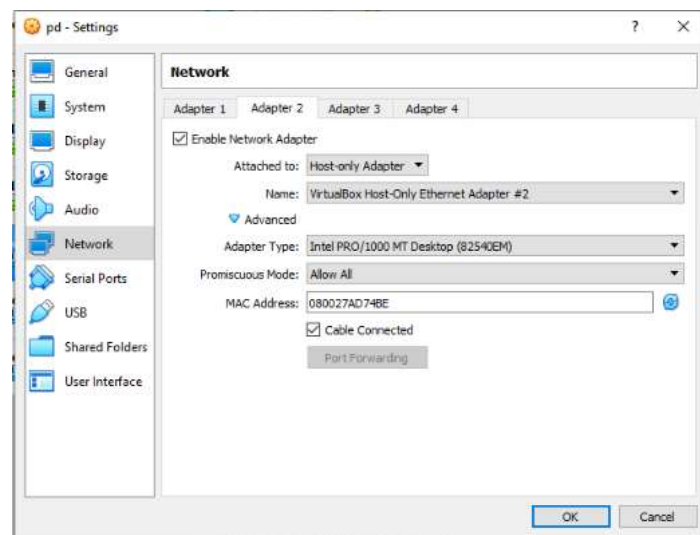
1st Step:





2nd Step:

Make sure all the VM's have the following network configuration on Adapter 2:



3rd Step:

Let's change the hostname on each virtual machine. Open the file and type the name of the machina. Use this command:

```
sudo nano /hostname
```



[Get unlimited access](#)[Open in app](#)

pd-master

GNU nano 2.9.3 /hostname Modified

pd-slave1

^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify
^X Exit ^R Read File ^\ Replace ^U Uncut Text ^T To Spell

torres@torres-VirtualBox: ~

File Edit View Search Terminal Help

GNU nano 2.9.3 /hostname Modified

pd-slave1

GNU nano 2.9.3 /hostname Modified

pd-slave2

^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify
^X Exit ^R Read File ^\ Replace ^U Uncut Text ^T To Spell

torres@torres-VirtualBox: ~

File Edit View Search Terminal Help

GNU nano 2.9.3 /hostname Modified

pd-slave2

GNU nano 2.9.3 /hostname Modified

pd-slave2

^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify
^X Exit ^R Read File ^\ Replace ^U Uncut Text ^T To Spell

4th Step:





Get unlimited access

Open in app

```

torres@torres-VirtualBox: ~
File Edit View Search Terminal Help
torres@torres-VirtualBox:~$ ip addr
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN group default qlen 1000
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
    inet 127.0.0.1/8 scope host lo
        valid_lft forever preferred_lft forever
    inet6 ::1/128 scope host
        valid_lft forever preferred_lft forever
2: enp0s3: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc fq_codel state UP group default qlen 1000
    link/ether 08:00:27:2a:eb:6b brd ff:ff:ff:ff:ff:ff
    inet 10.0.2.15/24 brd 10.0.2.255 scope global dynamic noprefixroute enp0s3
        valid_lft 86315sec preferred_lft 86315sec
    inet6 fe80::c66c:d5fa:7179:6628/64 scope link noprefixroute
        valid_lft forever preferred_lft forever
3: enp0s8: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc fq_codel state UP group default qlen 1000
    link/ether 08:00:27:ad:74:be brd ff:ff:ff:ff:ff:ff
    inet 192.168.205.10/24 brd 192.168.205.255 scope global dynamic noprefixroute enp0s8
        valid_lft 1116sec preferred_lft 1116sec
    inet6 fe80::9f4c:576b:e0dc:f4da/64 scope link noprefixroute
        valid_lft forever preferred_lft forever
torres@torres-VirtualBox:~$

```

This is on the master VM, as you can see our IP is 192.168.205.10. For you this will be different.

This means that our IP's are:

master: 192.168.205.10

slave1: 192.168.205.11

slave2: 192.168.205.12

5th Step:

We need to edit the **hosts** file. Use the following command:

```
sudo nano /etc/hosts
```

and add your network information:

```

torres@torres-VirtualBox: ~
File Edit View Search Terminal Help
GNU nano 2.9.3 /etc/hosts Modified
127.0.0.1    localhost
127.0.1.1    torres-VirtualBox

pd-master    192.168.205.10
pd-slave1    192.168.205.11
pd-slave2    192.168.205.12

# The following lines are desirable for IPv6 capable hosts
::1          ip6-localhost ip6-loopback
fe00::0      ip6-localnet
ff00::0      ip6-mcastprefix
ff02::1      ip6-allnodes
ff02::2      ip6-allrouters

```





```
sudo reboot
```

7th Step:

Do this step on all the Machines, master and slaves.

Now, in order to install Java we need to do some things. Follow these commands and give permission when needed:

```
$ sudo apt-get install software-properties-common
$ sudo add-apt-repository ppa:webupd8team/java
$ sudo apt-get update
$ sudo apt-get install openjdk-11-jdk
```

To check if java is installed, run the following command.

```
$ java -version
```

```
torres@pd-master:~$ java -version
openjdk version "11.0.6" 2020-01-14
OpenJDK Runtime Environment (build 11.0.6+10-post-Ubuntu-1ubuntu118.04.1)
OpenJDK 64-Bit Server VM (build 11.0.6+10-post-Ubuntu-1ubuntu118.04.1, mixed mode)
torres@pd-master:~$
```

8th Step:

Now let's install Scala on the master and the slaves. Use this command:

```
$ sudo apt-get install scala
```

```
torres@pd-master:~$ sudo apt-get install scala
Reading package lists... Done
Building dependency tree
Reading state information... Done
```

To check if Scala was correctly installed run this command:

```
$ scala -version
```

```
torres@pd-master:~$ scala -version
Scala code runner version 2.11.12 -- Copyright 2002-2017, LAMP/EPFL
torres@pd-master:~$
```

As you can see, Scala version 2.11.12 is now installed on my machine.





Get unlimited access

Open in app

```
ssh-copy-id user@pd-master
ssh-copy-id user@pd-slave1
ssh-copy-id user@pd-slave2
```

```
torres@pd-master:~$ ssh-copy-id torres@pd-slave1
The authenticity of host 'pd-slave1 (192.168.205.11)' can't be established.
ECDSA key fingerprint is SHA256:MG56/CnMiqav/9cS7HP8hXNN0EYzzol/7JLKq4cTLbE.
Are you sure you want to continue connecting (yes/no)? yes
/usr/bin/ssh-copy-id: INFO: attempting to log in with the new key(s), to filter
out any that are already installed
/usr/bin/ssh-copy-id: INFO: 1 key(s) remain to be installed -- if you are promp
ted now it is to install the new keys
torres@pd-slave1's password:

Number of key(s) added: 1

Now try logging into the machine, with: "ssh 'torres@pd-slave1'"
and check to make sure that only the key(s) you wanted were added.

torres@pd-master:~$
```

```
torres@pd-master:~$ ssh-copy-id torres@pd-slave2
The authenticity of host 'pd-slave2 (192.168.205.12)' can't be established.
ECDSA key fingerprint is SHA256:vchan4ALBXytJSnv+3L+nChJlQ1GmER3bWpXnMPy4k0.
Are you sure you want to continue connecting (yes/no)? yes
/usr/bin/ssh-copy-id: INFO: attempting to log in with the new key(s), to filter
out any that are already installed
/usr/bin/ssh-copy-id: INFO: 1 key(s) remain to be installed -- if you are promp
ted now it is to install the new keys
torres@pd-slave2's password:

Number of key(s) added: 1

Now try logging into the machine, with: "ssh 'torres@pd-slave2'"
and check to make sure that only the key(s) you wanted were added.

torres@pd-master:~$
```

```
torres@pd-master:~$ ssh-copy-id torres@pd-slave2
The authenticity of host 'pd-slave2 (192.168.205.12)' can't be established.
ECDSA key fingerprint is SHA256:vchan4ALBXytJSnv+3L+nChJlQ1GmER3bWpXnMPy4k0.
Are you sure you want to continue connecting (yes/no)? yes
/usr/bin/ssh-copy-id: INFO: attempting to log in with the new key(s), to filter
out any that are already installed
/usr/bin/ssh-copy-id: INFO: 1 key(s) remain to be installed -- if you are promp
ted now it is to install the new keys
torres@pd-slave2's password:

Number of key(s) added: 1

Now try logging into the machine, with: "ssh 'torres@pd-slave2'"
and check to make sure that only the key(s) you wanted were added.
```

Let's check if everything went well, try to connect to the slaves:

```
$ ssh slave01
$ ssh slave02
```

```
torres@pd-master:~$ ssh pd-slave1
Welcome to Ubuntu 18.04.2 LTS (GNU/Linux 4.18.0-15-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/advantage

 * Canonical Livepatch is available for installation.
   - Reduce system reboots and improve kernel security. Activate at:
     https://ubuntu.com/livepatch

463 packages can be updated.
250 updates are security updates.

Your Hardware Enablement Stack (HWE) is supported until April 2023.
torres@pd-slave1:~$
```





Get unlimited access

Open in app

```
* Management: https://landscape.canonical.com
* Support: https://ubuntu.com/advantage

* Canonical Livepatch is available for installation.
- Reduce system reboots and improve kernel security. Activate at:
https://ubuntu.com/livepatch

463 packages can be updated.
250 updates are security updates.

Your Hardware Enablement Stack (HWE) is supported until April 2023.
torres@pd-slave2:~$
```

As you can see everything went well, to exit just type the command:

```
exit
```

```
torres@pd-slave2:~$ exit
logout
Connection to pd-slave2 closed.
```

10th Step:

Now we Download the latest version of Apache Spark.

NOTE: Everything inside this step must be done on all the virtual machines.

Use the following command :

```
$ wget http://www-us.apache.org/dist/spark/spark-2.4.4/spark-2.4.4-bin-hadoop2.7.tgz
```

This is the most recent version as of the writing of this article, it might have changed if you try it later. Anyway, I think you'll still be good using this one.

```
torres@pd-master:~$ wget http://www-us.apache.org/dist/spark/spark-2.4.4/spark-
2.4.4-bin-hadoop2.7.tgz
--2020-02-03 19:41:41-- http://www-us.apache.org/dist/spark/spark-2.4.4/spark-
2.4.4-bin-hadoop2.7.tgz
Resolving www-us.apache.org (www-us.apache.org)... 40.79.78.1
Connecting to www-us.apache.org (www-us.apache.org)|40.79.78.1|:80... connected
.
HTTP request sent, awaiting response... 200 OK
Length: 230091034 (219M) [application/x-gzip]
Saving to: 'spark-2.4.4-bin-hadoop2.7.tgz'

.4.4-bin-hadoop2.7.  0%[                ] 856,50K  231KB/s  eta 16m 45s
```

Extract the Apache Spark file you just downloaded

Use the following command to extract the Spark tar file:

```
$ tar xvf spark-2.4.4-bin-hadoop2.7.tgz
```





Get unlimited access

Open in app

```
spark-2.4.4-bin-hadoop2.7/R/lib/sparkr.zip
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/INDEX
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/html/
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/html/R.css
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/html/00Index.html
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/help/
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/help/aliases.rds
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/help/AnIndex
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/help/SparkR.rdx
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/help/SparkR.rdb
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/help/paths.rds
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/worker/
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/worker/worker.R
spark-2.4.4-bin-hadoop2.7/R/lib/SparkR/worker/daemon.R
```

Move Apache Spark software files

Use the following command to move the spark software files to respective directory (*/usr/local/bin*)

```
$ sudo mv spark-2.4.4-bin-hadoop2.7 /usr/local/spark
```

```
torres@pd-master:~$ sudo mv spark-2.4.4-bin-hadoop2.7 /usr/local/spark
[sudo] password for torres:
torres@pd-master:~$
```

Set up the environment for Apache Spark

Edit the *bashrc* file using this command:

```
$ sudo gedit ~/.bashrc
```

```
torres@pd-master:~$ sudo gedit ~/.bashrc
```

Add the following line to the file. This adds the location where the spark software file are located to the PATH variable.

```
export PATH = $PATH:/usr/local/spark/bin
```

```
Open ▾  *.*bashrc  Save  ⋮  🔍  🌐  🗑️  📄  📁  📂  📅  📆  📇  📈  📉  📊  📋  📌  📍  📎  📏  📐  📑  📒  📓  📔  📕  📖  📗  📘  📙  📚  📛  📜  📝  📞  📟  📠  📡  📢  📣  📤  📥  📦  📧  📨  📩  📪  📫  📬  📭  📮  📯  📰  📱  📲  📳  📴  📵  📶  📷  📸  📹  📺  📻  📼  📽  📾  📿  📺  📻  📼  📽  📾  📿
```

```
# Add an "alert" alias for long running commands. Use like so:
# sleep 10; alert
alias alert='notify-send --urgency=low -i "${?} = 0" && echo terminal ||
echo error)' "${history|tail -n1|sed -e '\s/^s*[0-9]+\s+//;s/[:&|]
\s*alerts/\'\'\'')"
```

```
# Alias definitions.
# You may want to put all your additions into a separate file like
# ~/.bash_aliases, instead of adding them here directly.
# See /usr/share/doc/bash-doc/examples in the bash-doc package.

if [ -f ~/.bash_aliases ]; then
. ~/.bash_aliases
fi

# enable programmable completion features (you don't need to enable
# this, if it's already enabled in /etc/bash.bashrc and /etc/profile
# sources /etc/bash.bashrc).
if ! shopt -oq posix; then
if [ -f /usr/share/bash-completion/bash_completion ]; then
. /usr/share/bash-completion/bash_completion
elif [ -f /etc/bash_completion ]; then
. /etc/bash_completion
fi
fi

export PATH = $PATH:/usr/local/spark/bin
```

Note: this screenshot has a mistake, when you're doing this don't leave a space like I did. Just write "PATH=\$PATH".



[Get unlimited access](#)[Open in app](#)

```
torres@pd-master:~$ source ~/.bashrc
torres@pd-master:~$
```

11th Step:

Apache Spark Master Configuration (do this step on the Master VM only)

Edit spark-env.sh

Move to spark *conf* folder and create a copy of the template of *spark-env.sh* and rename it.

```
$ cd /usr/local/spark/conf
$ cp spark-env.sh.template spark-env.sh
```

```
torres@pd-master:/usr/local/spark/conf$ cp spark-env.sh.template spark-env.sh
torres@pd-master:/usr/local/spark/conf$
```

Now edit the configuration file *spark-env.sh*.

```
$ sudo vim spark-env.sh
```

And add the following parameters:

```
export SPARK_MASTER_HOST='<MASTER-IP>'export JAVA_HOME=<Path_of_JAVA_installation>
```

Add Workers

Edit the configuration file *slaves* in (*/usr/local/spark/conf*).

```
$ sudo nano slaves
```

```
torres@pd-master:/usr/local/spark/conf$ sudo nano slaves
```

And add the following entries.

```
pd-master
pd-slave01
pd-slave02
```





Get unlimited access

Open in app

```
pd-master
pd-slave1
pd-slave2
```

^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify
 ^X Exit ^R Read File ^\ Replace ^U Uncut Text ^T To Spell

12th Step:

Let's try to start our Apache Spark Cluster, hopefully everything is ok!

To start the spark cluster, run the following command on master.:

```
$ cd /usr/local/spark
$ ./sbin/start-all.sh
```

```
torres@pd-master:/usr/local/spark/conf$ cd /usr/local/spark
torres@pd-master:/usr/local/spark$ ./sbin/start-all.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/log
s/spark-torres-org.apache.spark.deploy.master.Master-1-pd-master.out
pd-slave1: starting org.apache.spark.deploy.worker.Worker, logging to /usr/loca
l/spark/logs/spark-torres-org.apache.spark.deploy.worker.Worker-1-pd-slave1.out
pd-master: starting org.apache.spark.deploy.worker.Worker, logging to /usr/loca
l/spark/logs/spark-torres-org.apache.spark.deploy.worker.Worker-1-pd-master.out
```

I won't stop it, but in case you want to stop the cluster, this is the command:

```
$ ./sbin/stop-all.sh
```

13th Step:

To check if the services started we use the command:


```
$ jps
```

```
torres@pd-master:/usr/local/spark$ jps
9280 Jps
9065 Master
9212 Worker
```

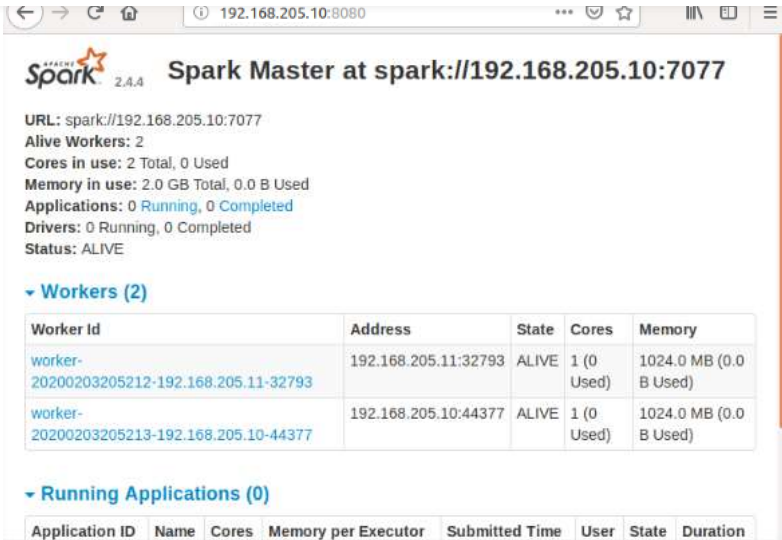
14th Step:

Browse the Spark UI to know about your cluster. To do this, go to your browser and type:





Get unlimited access [Open in app](#)



The screenshot shows the Apache Spark Master web interface at `spark://192.168.205.10:7077`. The interface displays the following information:

- URL:** `spark://192.168.205.10:7077`
- Alive Workers:** 2
- Cores in use:** 2 Total, 0 Used
- Memory in use:** 2.0 GB Total, 0.0 B Used
- Applications:** 0 Running, 0 Completed
- Drivers:** 0 Running, 0 Completed
- Status:** ALIVE

Below this information, there is a section for **Workers (2)** with the following table:

Worker Id	Address	State	Cores	Memory
worker-20200203205212-192.168.205.11-32793	192.168.205.11:32793	ALIVE	1 (0 Used)	1024.0 MB (0.0 B Used)
worker-20200203205213-192.168.205.10-44377	192.168.205.10:44377	ALIVE	1 (0 Used)	1024.0 MB (0.0 B Used)

At the bottom, there is a section for **Running Applications (0)** with a table header:

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------

As you can see we have 2 Alive Workers, our slaves, which means it's all done!

Final considerations:

Hopefully you managed to successfully follow this tutorial and have a perfectly working Apache Spark Cluster.

Any doubts feel free to ask me.

See you later!



Get unlimited access

Open in app





Get unlimited access

Open in app

