

Configuración de Spark en AWS EC2 con Python y Jupyter Notebooks

Nota: La instancia de EC2 en AWS debe de tener mínimo 30 GB de espacio en disco.

Descarga e Instala anaconda:

```
$ wget https://repo.anaconda.com/archive/Anaconda3-2023.07-2-Linux-x86_64.sh
```

```
$ bash Anaconda3-2023.07-2-Linux-x86_64.sh
```

Ver que version de Python estás utilizando y cambiarlo:

```
$ which python3
```

```
$ source .bashrc
```

```
$ which python3
```

Configuración Jupyter Notebook:

```
$ jupyter notebook --generate-config
```

Crear certificados:

```
$ mkdir certs
```

```
$ cd certs
```

```
$ sudo openssl req -x509 -nodes -days 365 -newkey rsa:1024 -keyout mycert.pem -out mycert.pem
```

Editar archivo de configuración:

```
$ cd ~/.jupyter/
```

```
$ vi jupyter_notebook_config.py
```

```
c = get_config()

# Notebook config this is where you saved your pem cert
# c.NotebookApp.certfile = u'/home/ubuntu/certs/mycert.pem'

# Run on all IP addresses of your instance
c.NotebookApp.ip = '*'

# Don't open browser by default
c.NotebookApp.open_browser = False

# Fix port to 8888
c.NotebookApp.port = 8888
```

Agregar el Puerto de Jupyter como Inbound rules.
Seleccionar la instancia en el listado, ir a Security.

Instance: i-0230ed333d4bd7238 (aws spark)

Details

Security

Networking

Storage

Status checks

Monitoring

Tags

▼ Security details

IAM Role

-

Owner ID

330189987704

Launch time

Mon Oct 02 2023 10:53:59 GMT-0500 (Central Daylight Time)

Seleccionar el Security group

Security groups

 **sg-0e73fb7ee53828781 (launch-wizard-4)**

Clic en Edit inbound rules

Inbound rules

Outbound rules

Tags

Inbound rules (5)

Manage tags

Edit inbound rules

☐

Name

▼

☐

Security group rule...

▼

☐

IP version

▼

☐

Type

▼

☐

Protocol

▼

☐

Port range

▼

<input type="checkbox"/>	-	sg-0edc75ff58f507eb8	IPv4	HTTP	TCP	80
<input type="checkbox"/>	-	sg-0a5a9819219738...	IPv4	SSH	TCP	22
<input type="checkbox"/>	-	sg-0c518626c916a8b...	IPv4	Custom TCP	TCP	8890
<input type="checkbox"/>	-	sg-0ad648d3462338...	IPv4	HTTPS	TCP	443
<input type="checkbox"/>	-	sg-0115e7e51aac6d50b	IPv4	Custom TCP	TCP	8888

Agregar el Puerto 8888 o según corresponda.

0.0.0.0/0 ✕

sg-0c518626c916a8b7b

Custom TCP ▼

TCP

8890

Cust... ▼

Delete

sg-0ad648d346233896a

HTTPS ▼

TCP

443

Cust... ▼

Delete

sg-0115e7e51aac6d50b

Custom TCP ▼

TCP

8888

Cust... ▼

Delete

Add rule

Lanzar Jupyter Notebook:

```
$ jupyter notebook
```

Instalar Java:

```
$ sudo apt-get update
```

```
$ sudo apt-get install default-jre
```

```
$ java -version
```

Instalar Scala:

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

```
$ scala -version
```

Instalar py4j:

```
$ export PATH=$PATH:$HOME/anaconda3/bin
```

```
$ conda install pip
```

```
$ which pip
```

```
$ pip install py4j
```

Instalar Spark y Hadoop:

*asegurate de estar en el directorio principal

```
$ cd
```

```
$ wget http://archive.apache.org/dist/spark/spark-3.5.0/spark-3.5.0-bin-hadoop3.tgz
```

```
$ sudo tar -zxvf spark-3.5.0-bin-hadoop3.tgz
```

Decirle a Python donde encontrar Spark:

```
$ export SPARK_HOME='/home/ubuntu/spark-3.5.0-bin-hadoop3'
```

```
$ export PATH=$SPARK_HOME:$PATH
```

```
$ export PYTHONPATH=$SPARK_HOME/python:$PYTHONPATH
```

También puedes añadir esas variables a tu profile:

```
$ vi ~/.bashrc  
$ source ~/.bashrc
```

Lanzar Jupyter Notebook:

```
$ jupyter notebook
```

Lanzar Spark:

```
from pyspark import SparkContext  
sc = SparkContext()
```

```
import pyspark
from pyspark.sql import SparkSession

# Create SparkSession
spark = SparkSession.builder.master("local[*]") \
    .appName('BigData-ETL.com') \
    .getOrCreate()

print(f'The PySpark {spark.version} version is running...')

The PySpark 3.5.0 version is running...
```