

EEL 6761: Cloud Computing
Homework 2
Spring 2022

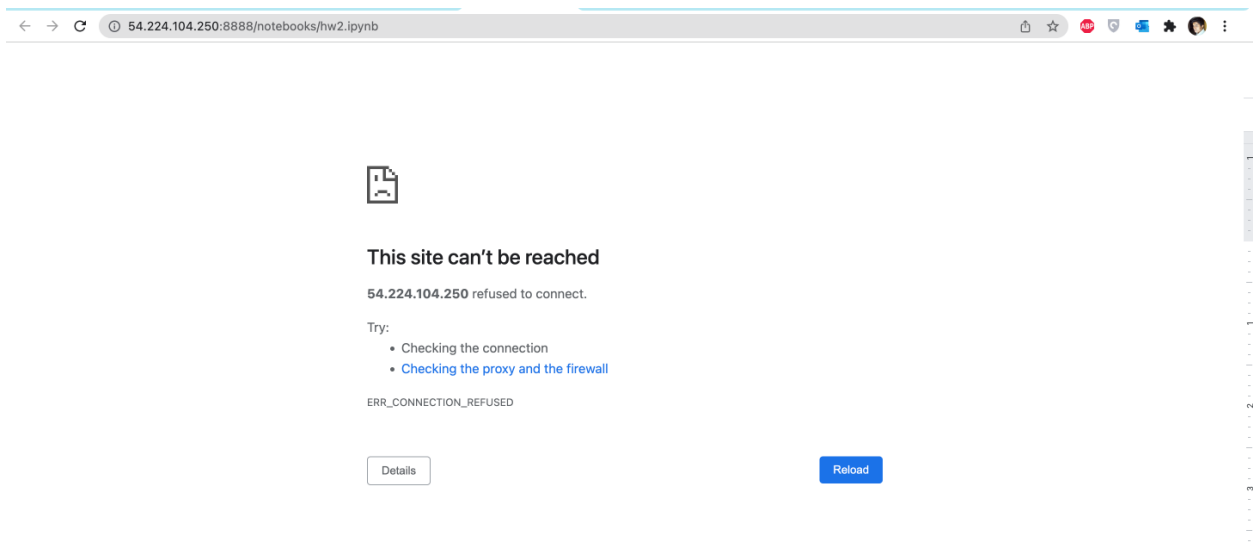
Deploying an Elastic Cloud Service

Anmol Lingan Gouda Patil
UFID: 1967 3150

A. When configuring PySpark to use Jupyter as its driver, you were asked to set its execution options using the `PYSPARK_DRIVER_PYTHON_OPTS` environment variable. Why is the `--ip=*` option needed? What happens if it is omitted?

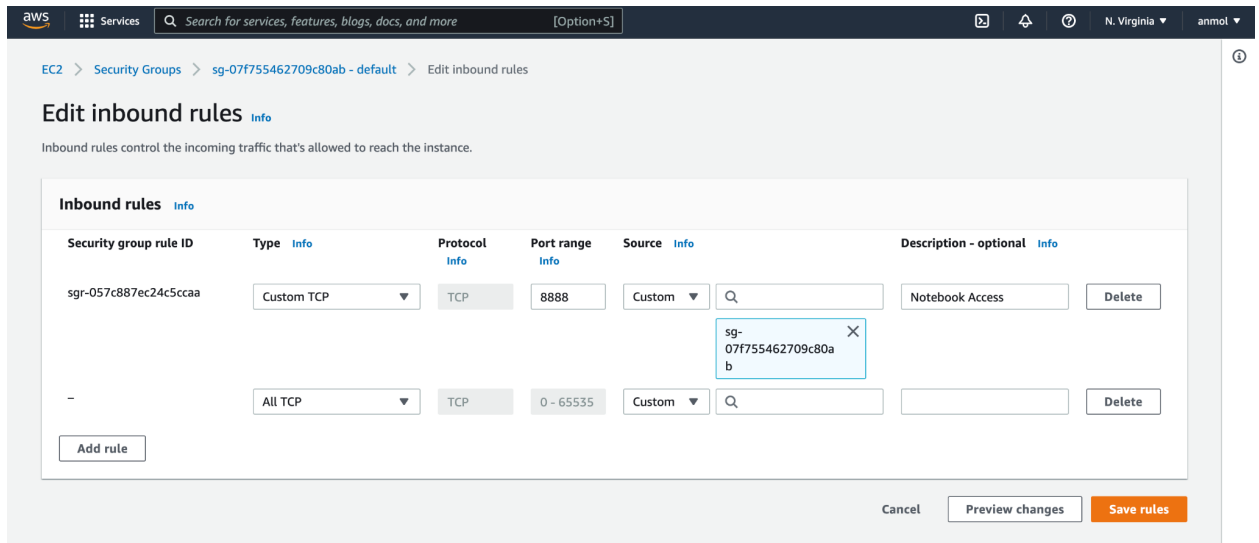
Ans: As per the Jupyter Notebook user documentation, by default a notebook server runs locally on the localhost i.e `ip = 127.0.0.1` and `port = 8888`. And the access to the notebook server via a browser is limited only to <http://127.0.0.1:8888>. But when we are trying to use the jupyter notebook from an external source outside of the local machine(launched EC2), we need the jupyter server to allow accepting requests from different/external interfaces(ips). Thus the `--ip=*` option allows the jupyter server to bind on all interfaces for the public server(launched EC2) thereby allowing us to access the notebook outside of the local instance.

When it is omitted we will not be able to access the jupyter notebook server from outside the EC2 instance, and a connection refused error like the one below is thrown.



B. What is the difference between “Custom TCP” and “All TCP” types in the EC2 Security Groups? How can you protect your jupyter notebook apart from the credentialed access? Provide the steps.

Ans: As seen in the image below, when the **Custom TCP** option is selected we must specify a port or a range of ports to limit the incoming tcp traffic to only such port or range ports. However when the **All TCP** option is chosen, AWS auto populates the port range to all the ports from 0 to 65535, stating that when you choose All TCP, incoming tcp traffic shall bind on any ports from 0 to 65535.



Besides protecting the jupyter notebook using a password here are some more steps to achieve additional security:

1. **Preparing a hashed password:** we can prepare a hashed password manually by using the function: `notebook.auth.security.passwd()`. This function when called will ask you for your password and verify it by retyping and then outputs a sha1 hash.
2. **Using hashed password in notebook configuration file/**
PYSPARK_DRIVER_PYTHON_OPTS:
 we can then add the hashed password from previous step to `jupyter_notebook_config.py` file or set the `–NotebookApp.password =` the hashed string in unicode, if we are using the **PYSPARK_DRIVER_PYTHON_OPTS**.
3. **Encrypted communication using secure socket layer(SSL):** a hashed password needs to be encrypted before it is sent via the browser, thus it is prudent to use SSL with a self signed web certificate. Mention the certificate file path using the `–certfile=` option along with the key using `–keyfile yourkey.key`. A self signed certificate can be generated using the following command using the openssl:

```
openssl req -x509 -nodes -days 365 -newkey rsa:2048 -keyout your_key.key
-out your_cert.pem
```

C. When configuring the pyspark service used by systemd, you were asked to include the following lines in your `pyspark.service` file:

```
After=network.target
After=systemd-user-sessions.service
After=network-online.target
```

Why are these lines needed? What happens if they are omitted?

Ans:

network.target:

This unit is supposed to indicate when network functionality is available.

systemd-user-sessions.service:

is a service that controls user logins through pam_nologin(Prevent non-root users from login).

network-online.target:

This target unit is intended to pull in a service that delays further execution until the network is sufficiently set up.

So when we append these options to **After=** , we are essentially stating that we do not start our service until the system on which it is going to run has network functionality up and running and the root user has logged into the system. And if the internet connectivity is not yet established, delay it until it is sufficiently set up.

If these lines are omitted, we will not be able to connect to the jupyter notebook server from outside the local instance in which it is running.

D. Does increasing the number of instances in your cloud decrease the execution time of Spark jobs? Why or why not?

Ans:

In the current scope of the assignment, since the calculation of the value of **Pi** is an infinite process, hence we cannot measure the execution time.

However, since we are running the spark jobs in a **standalone mode** and not **cluster mode** there is no distribution of work across instances. All instances are running their own version of the process. If we would have run the spark job in a cluster mode, then actual distribution of work could be done using parallelization that spark takes care of. Thus increasing the number of instances will not decrease the execution time of the spark jobs

E. What happens if the minimum number of instances in the autoscaling group is set to zero?

Ans:

If the minimum number of instances in the autoscaling group is set to 0 along with the **desired capacity** option = 0. The autoscaling group will not launch any instances. But even if your minimum capacity is zero, the autoscaling group launches the maximum of (min capacity, desired capacity) # of instances. Suppose we set minimum capacity to 0 and desired capacity to 0 there would be no process started/ instance launched. However if the desired capacity is greater than 0 , say 1 then the autoscaling group will launch an instance and then later follows the target policy set.

SCREENSHOTS

← → ↺ console.aws.amazon.com/ec2/v2/home?region=us-east-1#instanceDetails:instanceId=i-034adde1b789402ff

aws

Services

Search for services, features, blogs, docs, and more

[Option+S]

N. Virginia

anmol

☰

EC2 > Instances > i-034adde1b789402ff

Instance summary for i-034adde1b789402ff

Info

Updated less than a minute ago

Refresh

Connect

Instance state

Actions

| | | |
|--|---|---|
| Instance ID i-034adde1b789402ff | Public IPv4 address 35.171.7.18 open address | Private IPv4 addresses 172.31.88.97 |
| IPv6 address - | Instance state Running | Public IPv4 DNS ec2-35-171-7-18.compute-1.amazonaws.com open address |
| Hostname type IP name: ip-172-31-88-97.ec2.internal | Private IP DNS name (IPv4 only) ip-172-31-88-97.ec2.internal | Answer private resource DNS name IPv4 (A) |
| Instance type t2.micro | Elastic IP addresses - | VPC ID vpc-08de330c45d57b992 |
| AWS Compute Optimizer finding Opt-in to AWS Compute Optimizer for recommendations. Learn more | IAM Role - | Subnet ID subnet-08773bb0e7439d66d |

Details

Security

Networking

Storage

Status checks

Monitoring

Tags

▼ Instance details Info

| | | |
|--|---|------------------------------------|
| Platform Ubuntu (Inferred) | AMI ID ami-04505e74c0741db8d | Monitoring disabled |
| Platform details Linux/UNIX | AMI name ubuntu/images/hvm-ssd/ubuntu-focal-20.04-amd64-server-20211129 | Termination protection Disabled |
| Launch time Mon Feb 28 2022 16:46:52 GMT-0500 (Eastern Standard Time) (2 minutes) | AMI location 099720109477/ubuntu/images/hvm-ssd/ubuntu-focal-20.04-amd64-server-20211129 | Lifecycle normal |

A

```
12 updates can be applied immediately.
8 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable

The list of available updates is more than a week old.
To check for new updates run: sudo apt update

Last login: Sat Feb 19 01:35:18 2022 from 184.188.101.162
ubuntu@ip-172-31-88-97:~$
```

B

C

← → ↻ Not Secure | 35.171.7.18:8888/notebooks/hw2.ipynb

jupyter hw2 Last Checkpoint: 02/13/2022 (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (pykernel)

```
In [1]: import math
        from pyspark.sql import SparkSession

In [2]: spark = SparkSession.builder.getOrCreate()
        sc = spark.sparkContext

In [*]: import datetime
        x = 0
        i = 0

        while True:
            n = 10000
            k_values = [i + x for x in range(1, n + 1)]
            num_partitions = 100
            dat = sc.parallelize(k_values, num_partitions)
            sqrs = dat.map(lambda k: (1.0/k) ** 2)
            x += sqrs.reduce(lambda a,b: a + b)
            pi = math.sqrt(x * 6)
            i += n
            print("Value of pi = ", pi, "Estimated at ", datetime.datetime.now().time())

Value of pi = 3.1414971639472076 Estimated at 22:19:50.111838
```

D-1

← → ↻ Not Secure | 35.171.7.18:8888/notebooks/hw2.ipynb

jupyter hw2 Last Checkpoint: 02/13/2022 (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

```
In [1]: import math
        from pyspark.sql import SparkSession

In [2]: spark = SparkSession.builder.getOrCreate()
        sc = spark.sparkContext

In [*]: import datetime
        x = 0
        i = 0

        while True:
            n = 10000
            k_values = [i + x for x in range(1, n + 1)]
            num_partitions = 100
            dat = sc.parallelize(k_values, num_partitions)
            sqrs = dat.map(lambda k: (1.0/k) ** 2)
            x += sqrs.reduce(lambda a,b: a + b)
            pi = math.sqrt(x * 6)
            i += n
            print("Value of pi = ", pi, "Estimated at ", datetime.datetime.now().time())

Value of pi = 3.141589367328518 Estimated at 22:20:47.324613

Value of pi = 3.141589470494622 Estimated at 22:20:49.368158

Value of pi = 3.141589573174996 Estimated at 22:20:51.387827

Value of pi = 3.141589669437854 Estimated at 22:20:53.351080

Value of pi = 3.141589759866605 Estimated at 22:20:55.292653

Value of pi = 3.141589844976023 Estimated at 22:20:57.277114
```

D-2

← → ↻ Not Secure | 35.171.7.18:8888/notebooks/hw2.ipynb

jupyter hw2 Last Checkpoint: 02/13/2022 (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

```
In [1]: import math
        from pyspark.sql import SparkSession

In [2]: spark = SparkSession.builder.getOrCreate()
        sc = spark.sparkContext

In [*]: import datetime
        x = 0
        i = 0

        while True:
            n = 10000
            k_values = [i + x for x in range(1, n + 1)]
            num_partitions = 100
            dat = sc.parallelize(k_values, num_partitions)
            sqrs = dat.map(lambda k: (1.0/k) ** 2)
            x += sqrs.reduce(lambda a,b: a + b)
            pi = math.sqrt(x * 6)
            i += n
            print("Value of pi = ", pi, "Estimated at ", datetime.datetime.now().time())
```

Value of pi = 3.1415909782756235 Estimated at 22:21:43.174962

Value of pi = 3.1415910071603332 Estimated at 22:21:45.155353

Value of pi = 3.141591035065901 Estimated at 22:21:47.084930

Value of pi = 3.1415910620412832 Estimated at 22:21:49.032800

Value of pi = 3.1415910881322278 Estimated at 22:21:51.012221

D-3