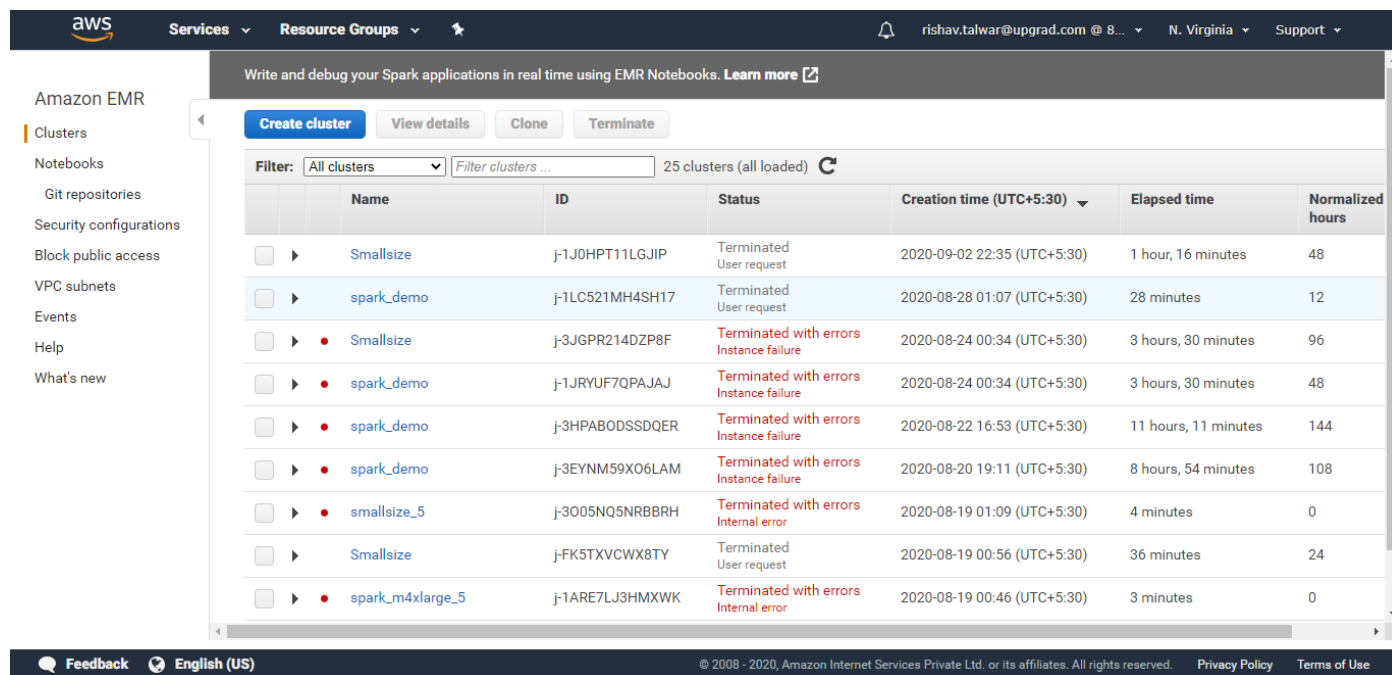# Creating an EMR Cluster on AWS

Before starting to create an EMR cluster, you will have to go to the EMR home page on your AWS account.

Once you arrive at the EMR home page, it should look something like this:



From here, follow the steps below to create an EMR Cluster:

**Step 1:** First, you need to click on the blue '**Create cluster**' button near the top left of the screen. Once you have clicked on this button, the following page should open on your screen:

### General Configuration

Cluster name    My cluster

☑ Logging ⓘ

S3 folder    s3://aws-logs-864328032829-us-east-1/elasticmaprec

Launch mode    ● Cluster ⓘ    ○ Step execution ⓘ

### Software configuration

Release    emr-5.30.1

Applications    ● Core Hadoop: Hadoop 2.8.5, Hive 2.3.6, Hue 4.6.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.2

○ HBase: HBase 1.4.13, Hadoop 2.8.5, Hive 2.3.6, Hue 4.6.0, Phoenix 4.14.3, and ZooKeeper 3.4.14

○ Presto: Presto 0.232 with Hadoop 2.8.5 HDFS and Hive 2.3.6 Metastore

○ Spark: Spark 2.4.5 on Hadoop 2.8.5 YARN and Zeppelin 0.8.2

☐ Use AWS Glue Data Catalog for table metadata ⓘ

### Hardware configuration

Instance type    m5.xlarge    The selected instance type adds 64 GiB of GP2 EBS storage per instance by default. Learn more ⎘

Number of instances    3    (1 master and 2 core nodes)

Cluster scaling    ☐ scale cluster nodes based on workload

### Security and access

EC2 key pair    Choose an option    ⓘ Learn how to create an EC2 key pair.

Permissions    ● Default    ○ Custom

Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

**Step 2:** Here, under 'General Configuration', write an appropriate Cluster name. Disable 'Logging' for this cluster, as this will lead to additional costs, and we do not need Logging services for the purposes of this module. Also, select 'Cluster' as the launch mode.

## General Configuration

Cluster name    Cluster_1

☐ Logging ⓘ

Launch mode    ● Cluster ⓘ    ○ Step execution ⓘ

**Step 3:** Under Software configuration', select **"emr-5.30.1"** as the Release. After this, you need to choose the Applications for your EMR cluster. Under 'Applications', since we will need Spark in our EMR cluster, choose the fourth option, which is for a Spark cluster.

## Software configuration

| | |
|---|---|
| Release | emr-5.30.1 ⌄ ⓘ |
| Applications | ◯ Core Hadoop: Hadoop 2.8.5, Hive 2.3.6, Hue 4.6.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.2 |
| | ◯ HBase: HBase 1.4.13, Hadoop 2.8.5, Hive 2.3.6, Hue 4.6.0, Phoenix 4.14.3, and ZooKeeper 3.4.14 |
| | ◯ Presto: Presto 0.232 with Hadoop 2.8.5 HDFS and Hive 2.3.6 Metastore |
| | ⦿ Spark: Spark 2.4.5 on Hadoop 2.8.5 YARN and Zeppelin 0.8.2 |
| | ☐ Use AWS Glue Data Catalog for table metadata ⓘ |

**Step 4:** Under 'Hardware configuration', we need to choose the instance type and the number of instances for our EMR Spark cluster. Here, please note that the type of instance and the number of instances that you choose will significantly affect the costs that you will incur when using the EMR cluster. For the purposes of this module, we will be using a **three-machine cluster** of instance type '**m4.large**'. Uncheck the checkbox for 'Cluster scaling'. Choose the configuration as shown below:

## Hardware configuration

| | | |
|---|---|---|
| Instance type | m4.large ⌄ | The selected instance type adds 32 GiB of GP2 EBS storage per instance by default. Learn more ↗ |
| Number of instances | 3  (1 master and 2 core nodes) | |
| Cluster scaling | ☐ scale cluster nodes based on workload | |

**Step 5:** Finally, under 'Security and access', you need to select the EC2 key pair that you have used until now in this program for practising on the EC2 machine. This will be used if you want to SSH to the Master node of your EMR cluster. For permissions, EMR role and EC2 instance profile, the default settings will be used.

## Security and access

**EC2 key pair** spark123    ⓘ Learn how to create an EC2 key pair.

**Permissions** ● Default ○ Custom

Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

**EMR role** EMR_DefaultRole ⧉ ⓘ

**EC2 instance profile** EMR_EC2_DefaultRole ⧉ ⓘ

**Step 6:** You can now proceed ahead and click on the 'Create cluster' button and your cluster will be created. As soon as you do that, you will be taken to the screen below, which shows that your Spark EMR cluster is starting.

| Clone | Terminate | AWS CLI export |
|---|---|---|

## Cluster: Cluster_1   Starting

**Summary** | Application user interfaces | Monitoring | Hardware | Configurations | Events | Steps | Bootstrap actions

### Summary

**ID:** j-1VZFOCUXPC11Y
**Creation date:** 2020-09-03 23:01 (UTC+5:30)
**Elapsed time:** 0 seconds
**After last step completes:** Cluster waits
**Termination protection:** Off Change
**Tags:** -- View All / Edit
**Master public DNS:** --

### Configuration details

**Release label:** emr-5.30.1
**Hadoop distribution:** Amazon
**Applications:** Spark 2.4.5, Zeppelin 0.8.2
**Log URI:** --
**EMRFS consistent view:** Disabled
**Custom AMI ID:** --

### Application user interfaces

**Persistent user interfaces** ⧉: --
**On-cluster user interfaces** ⧉: --

### Network and hardware

**Availability zone:** --
**Subnet ID:** subnet-1a4cbb57 ⧉
**Master:** Provisioning 1 m4.large
**Core:** Provisioning 2 m4.large
**Task:** --
**Cluster scaling:** Not enabled

### Security and access

**Key name:** spark123
**EC2 instance profile:** EMR_EC2_DefaultRole
**EMR role:** EMR_DefaultRole
**Visible to all users:** All Change
**Security groups for Master:**
**Security groups for Core & Task:**

**Note:** If you need to access the Spark history server for your EMR cluster, you can easily get it by clicking on the 'Application user interface' tab and then clicking on the 'Spark history server' link. This will take you to the Spark history server. If you get any pop-up block error, ignore it and click on the link again.

## Persistent application user interfaces

Applications installed on the Amazon EMR cluster publish u cluster.

| Application user interface ⬀ |
|---|
| Spark history server |
| YARN timeline server |

**Spark** 2.4.5-amzn-0 **History Server**

**Event log directory:** s3a://prod.us-east-1.appinfo.src/j-1VZFOCUXPC11Y/sparklogs

Last updated: 2020-09-03 23:35:30

Client local time zone: Asia/Calcutta

**No completed applications found!**

Did you specify the correct logging directory? Please verify your setting of *spark.history.fs.logDirectory* listed above and whether you have the permissions to access it. It is also possible that your application did not run to completion or did not stop the SparkContext.

Show incomplete applications

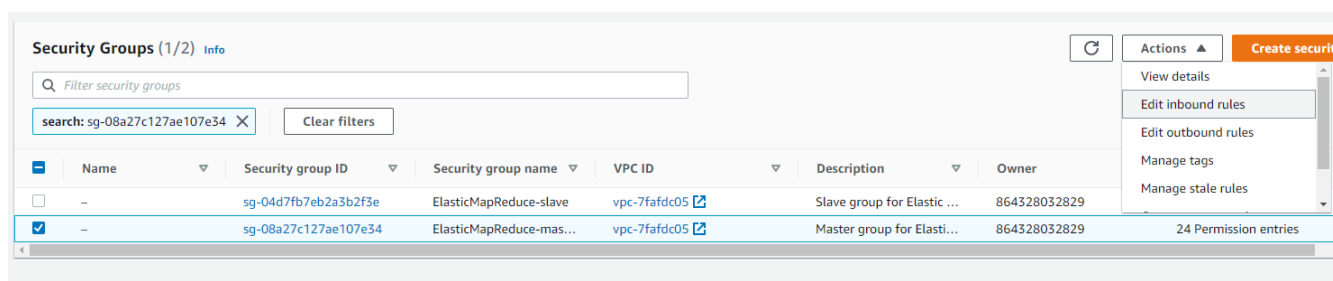# Setting Up the Security Group for the Master Node

You need to configure the security group for the Master node if you want to SSH to it. You will be doing this in the third session wherein you will be running the spark-submit commands on the Master node of the EMR cluster. The steps for setting up the security group of the Master node are the same as that for normal EC2 instances.

**Step 1:** First, click on the security group link that is displayed after the text 'Security groups for Master'. This will take you to the security group configuration screen.



**Step 2:** Next, click the checkbox for the Master group for EMR and then click on Actions, followed by clicking on Edit inbound rules, as shown below:

**Step 3:** Here, you need to add another Rule. Select the type of rule as 'ALL TCP' and source as 'My IP'. After this, save the rules and this step complete.

| All TCP ▼ | TCP | 0 - 65535 | My IP ▼ | 🔍 |
| | | | | 122.162.36.11/32 ✕ |

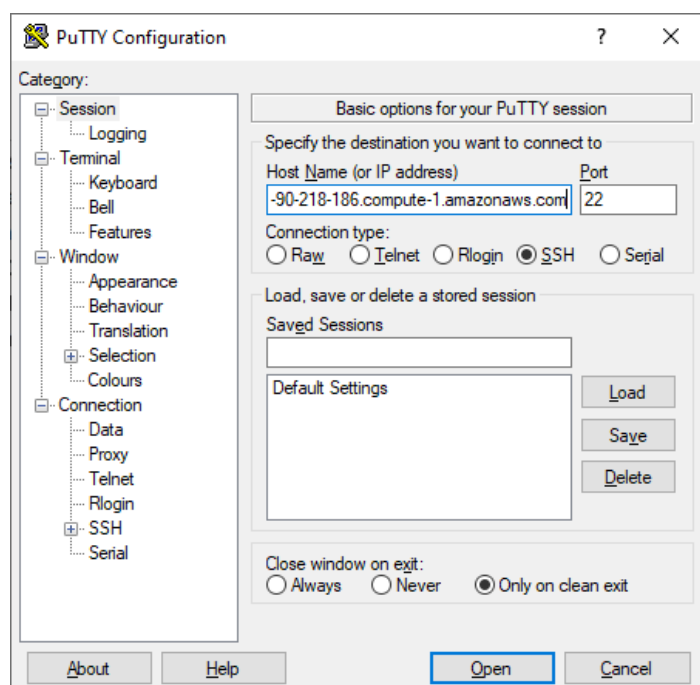# Connecting to the Master Node of Your Spark EMR Cluster

Once your cluster is up and running, you can connect to the Spark EMR cluster using SSH. The steps for this are as follows:
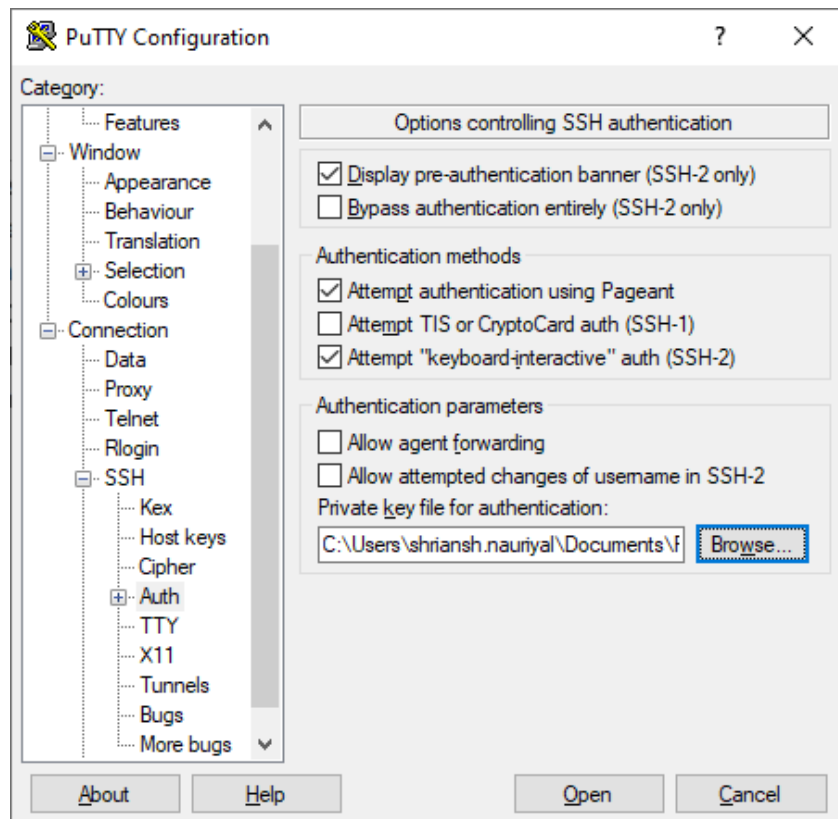
**For Windows Users**

**Step 1:** First, you need to copy the Master public DNS from the cluster Summary page. You can do this by clicking on the blue square icon adjacent to the DNS address to the right.

**Master public DNS:** ec2-54-90-218-186.compute-1.amazonaws.com 🔲
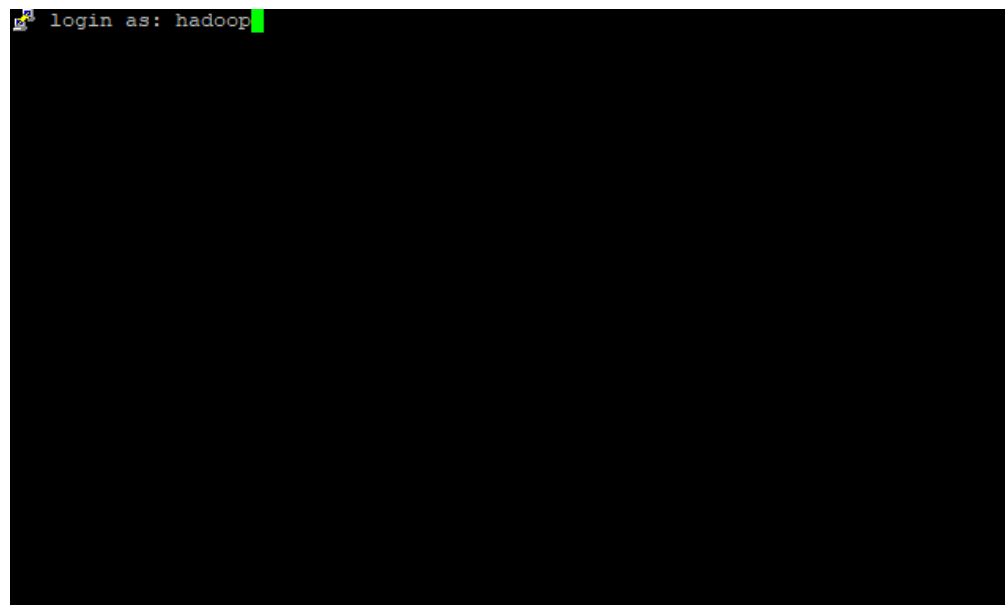Connect to the Master Node Using SSH

**Step 2:** After this, the steps to connect to the Master node are the same as the steps for connecting to an EC2 instance. First, open PuTTY and then paste the DNS address in the 'Host Name' field.

**Step 3:** Next, proceed to opening the SSH category, and under that, click on 'Auth' and then browse for the ppk file that you have configured to be used for this cluster.



**Step 4:** Next, click on 'Open'. Click 'Yes' if any pop-ups appear and then in 'login as', you need to type '**hadoop**', as you did for accessing an EC2 instance, and then click on Enter.

With this, you have gained access to the Master node of your EMR cluster.



**For Mac/Linux Users**

**Step 1:** First, you need to open the terminal on your machines.

**Step 2:** Next, you need to type the following command:

```
ssh -i <path to pem file> hadoop@<Master Public DNS>
```

For example, in my case, I can execute the following command to log in to the Master node:

```
ssh -i ~/spark123.pem hadoop@ec2-54-226-100-62.compute-1.amazonaws.com
```

**Step 3:** Finally, type 'Yes' to dismiss any warnings that may be displayed:

# Termination of an EMR Cluster

Once you have completed all your coding exercises on your EMR cluster, it is very important that you terminate the cluster. EMR is quite a costly service, and it must be kept activated only for the time that you are using it to run Spark jobs.
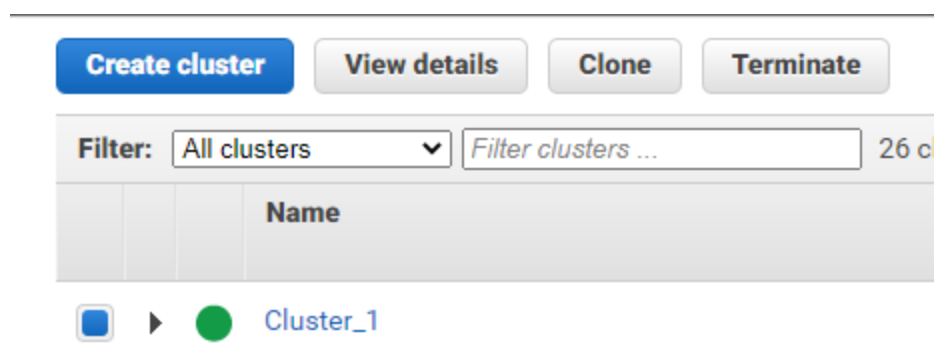
You can easily terminate your EMR cluster by following the steps below:

**Step 1:** You need to go to the cluster home page, which contains an entire list of all the EMR clusters in your account. You can do this by clicking on the 'Clusters' button to the top left under Amazon EMR.



**Step 2:** Next, you need to click the checkbox for your Spark EMR cluster and then click on the 'Terminate' button.

**Step 3:** A pop-up will appear, asking you to confirm whether you want to terminate the EMR cluster. Click on the 'Terminate' button in the pop-up. Once you click this button, your Spark EMR cluster will be terminated.