# Reto 5-1

1. Creación y configuración del clúster en AWS.

## Name and applications Info

**Name**

cluster1

**Amazon EMR release** | Info
A release contains a set of applications which can be installed on your cluster.

emr-6.3.1 ▼

**Application bundle**

| Spark | Core Hadoop | HBase | Presto | PrestoSQL | Custom |
|---|---|---|---|---|---|
| Spark | hadoop | HBASE | presto | trino | aws |

▼ **Customize your application bundle**

**Applications included in bundle**

- ☐ Flink 1.12.1
- ☐ HBase 2.2.6
- ☑ Hadoop 3.2.1
- ☑ Hue 4.9.0
- ☑ JupyterHub 1.2.0
- ☐ MXNet 1.7.0
- ☐ Phoenix 5.0.0
- ☐ Presto 0.245.1
- ☑ Spark 3.1.1
- ☐ TensorFlow 2.4.1
- ☑ Zeppelin 0.9.0
- ☐ Ganglia 3.7.2
- ☑ HCatalog 3.1.2
- ☑ Hive 3.1.2
- ☑ JupyterEnterpriseGateway 2.1.0
- ☑ Livy 0.7.0
- ☐ Oozie 5.2.1
- ☐ Pig 0.17.0
- ☐ PrestoSQL 350
- ☑ Sqoop 1.4.7
- ☐ Tez 0.9.2
- ☐ ZooKeeper 3.4.14

**AWS Glue Data Catalog settings**
Use the AWS Glue Data Catalog to provide an external metastore for your application.
- ☑ Use for Hive table metadata
- ☑ Use for Spark table metadata

**Custom Amazon Machine Image (AMI)** | Info

🔍 Choose or enter an AMI ID

- ☑ Update all installed packages on reboot

## Cluster configuration Info
Choose a configuration method for the primary, core, and task node groups for your cluster.

- ⦿ **Instance groups**
  Choose one instance type per node group
- ○ **Instance fleets**
  Choose any combination of instance types within each node group

### Instance groups

**Primary**
Choose EC2 instance type

| m4.xlarge<br>4 vCore   16 GiB memory   EBS only storage<br>On-Demand price: -   Lowest Spot price: - ▼ | Actions ▼ |
|---|---|

- ☐ Use multiple primary nodes
  To improve cluster availability, use 3 primary nodes with the same configuration and bootstrap actions. You can not use multiple primary nodes with instance fleets.

**Core**

Choose EC2 instance type

| m4.xlarge<br>4 vCore   16 GiB memory   EBS only storage<br>On-Demand price: -   Lowest Spot price: - | ▼ |

| Actions ▼ |

▶ Node configuration - *optional*

---

**Task 1 of 1**

| Remove instance group |

Name

| Task - 1 |

Choose EC2 instance type

| m4.xlarge<br>4 vCore   16 GiB memory   EBS only storage<br>On-Demand price: -   Lowest Spot price: - | ▼ |

| Actions ▼ |

▶ Node configuration - *optional*

---

| Add task instance group |

You can add up to 47 more task instance groups.

---

### Cluster scaling and provisioning option Info

Amazon EMR console only supports EMR-managed scaling. To create a cluster with auto-scaling, use CLI or SDK.

| ● Set cluster size manually<br>Use this option if you know your workload patterns in advance. | ○ Use EMR-managed scaling<br>Monitor key workload metrics so that EMR can optimize the cluster size and resource utilization. |

| Name | Instance type | Size | | Use Spot purchasing option |
| --- | --- | --- | --- | --- |
| Core | m4.xlarge | 1 ⇕ | instance(s) | ☐ |
| Task - 1 | m4.xlarge | 1 ⇕ | instance(s) | ☐ |

---

### Networking Info

Virtual private cloud (VPC)   Info

| vpc-03a3ccc8477f9b4bd |

| Browse |   | Create VPC ⧉ |

Subnet   Info

| subnet-0075a8b1af477f0dd |

| Browse |   | Create subnet ⧉ |

## ▼ Software settings - *optional* Info

| ● Enter configuration | ○ Load JSON from Amazon S3 |
|---|---|

```json
1 ▼ [
2 ▼   {
3       "Classification": "jupyter-s3-conf",
4 ▼     "Properties": {
5         "s3.persistence.enabled": "true",
6         "s3.persistence.bucket": "st0263eatorresm"
7       }
8     }
9 ]
```

JSON    Ln 9, Col 2    ⚙

## Identity and Access Management (IAM) roles Info
Choose or create a service role and instance profile for the EC2 instances in your cluster.

### Amazon EMR service role Info
The service role is an IAM role that Amazon EMR assumes to provision resources and perform service-level actions with other AWS services.

| ● Choose an existing service role | ○ Create a service role |
|---|---|
| Select a default service role or a custom role with IAM policies attached so that your cluster can interact with other AWS services. | Let Amazon EMR create a new service role so that you can grant and restrict access to resources in other AWS services. |

**Service role**

| EMR_DefaultRole ▼ | ⟳ |
|---|---|

### EC2 instance profile for Amazon EMR
The instance profile assigns a role to every EC2 instance in a cluster. The instance profile must specify a role that can access the resources for your steps and bootstrap actions.

| ● Choose an existing instance profile | ○ Create an instance profile |
|---|---|
| Select a default role or a custom instance profile with IAM policies attached so that your cluster can interact with your resources in Amazon S3. | Let Amazon EMR create a new instance profile so that you can specify a custom set of resources for it to access in Amazon S3. |

**Instance profile**

| EMR_EC2_DefaultRole ▼ | ⟳ |
|---|---|

2. Creación del bucket

## Create bucket Info

Buckets are containers for data stored in S3. Learn more ⧉

### General configuration

Bucket name

st0263eatorresm

Bucket name must be globally unique and must not contain spaces or uppercase letters. See rules for bucket naming ⧉

AWS Region

US East (N. Virginia) us-east-1 ▼

Copy settings from existing bucket - *optional*
Only the bucket settings in the following configuration are copied.

**Choose bucket**

**Buckets (1)** Info
Buckets are containers for data stored in S3. Learn more ⧉

🔍 Find buckets by name

| Name | | AWS Region | | Access | | Creation date | |
|---|---|---|---|---|---|---|---|
| ○ st0263eatorresm | | US East (N. Virginia) us-east-1 | | Bucket and objects not public | | May 7, 2023, 11:55:24 (UTC-05:00) | |

3. Configuración de seguridad.

| Custom TCP ▼ | TCP | 8888 | Anywh... ▼ | 🔍 | | Delete |
|---|---|---|---|---|---|---|
| | | | | 0.0.0.0/0 ✕ | | |
| Custom TCP ▼ | TCP | 9443 | Anywh... ▼ | 🔍 | | Delete |
| | | | | 0.0.0.0/0 ✕ | | |
| Custom TCP ▼ | TCP | 8890 | Anywh... ▼ | 🔍 | | Delete |
| | | | | 0.0.0.0/0 ✕ | | |
| Custom TCP ▼ | TCP | 22 | Anywh... ▼ | 🔍 | | Delete |
| | | | | 0.0.0.0/0 ✕ | | |

4. Conexión con el clúster, Hue y Zeppelin.

Search data and saved documents...

**📄 File Browser**

Search for file name    ⚙ Actions ▾    ⚡ Delete forever    ⊕ Upload    ⊕ New ▾

🏠 Home  /  user  /  **hadoop**

| | Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|
| ☐ | 📁 ↰ | | hdfs | hdfsadmingroup | drwxr-xr-x | May 07, 2023 11:08 AM |
| ☐ | 📁 . | | hadoop | hdfsadmingroup | drwxrwxrwx | May 07, 2023 12:09 PM |
| ☐ | 📁 datasets | | hadoop | hdfsadmingroup | drwxr-xr-x | May 07, 2023 12:09 PM |

Show 45 ▾ of 1 items    Page 1 of 1  ⏮ ⏪ ⏩ ⏭

---

MySQL
Databases (0) ↻
*Error loading databases.*

Search data and saved documents...

'server_name'

**📄 File Browser**

Search for file name    ⚙ Actions ▾    ⚡ Delete forever    ⊕ Upload    ⊕ New ▾

☁ us-east-1    s3a://

| | Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|
| ☐ | ☁ . | | | | drwxrwxrwx | |
| ☐ | ☁ datasetseatorresm | | | | drwxrwxrwx | |
| ☐ | ☁ st0263eatorresm | | | | drwxrwxrwx | |

Show 45 ▾ of 2 items    Page 1 of 1  ⏮ ⏪ ⏩ ⏭

**prueba**  ▷ ⤢ 🗐 ✎ 🗗 ⬆ ⬆    📄 ⚙ ⇄ Head    🔍    🗑    ⌨ ⚙ 🔒 default ▼

```
%spark.pyspark
spark
```
FINISHED ▷ ⤢ 🗐 ⚙

```
<pyspark.sql.session.SparkSession object at 0x7fe1ffb81290>
```

Took 2 sec. Last updated by anonymous at May 07 2023, 2:14:24 PM.

```
%spark.pyspark
sc
```
FINISHED ▷ ⤢ 🗐 ⚙

```
<SparkContext master=yarn appName=Zeppelin>
```

Took 0 sec. Last updated by anonymous at May 07 2023, 2:14:42 PM.

```
%sql
show databases
```
FINISHED ▷ ⤢ 🗐 ⚙

📊 📈 🌐 📊 📈 📉    ⬇ ▼    settings ▼

| namespace | ⌄ | ≡ |
|-----------|---|---|
| default | | |

Took 6 sec. Last updated by anonymous at May 07 2023, 2:15:08 PM. (outdated)