
Amazon AWS introduction

CS-543

Prerequisites

Amazon services we need:

Amazon provides a quite a few services. We only need:

- EMR
- S3

EMR is the service that allows you to instantiate a cluster, rented by amazon with all the tools you will need (hadoop - spark-shell).


S3 is persistent storage for your files accessible from your cluster machines. When you stop a VM you lose any data loaded to it, so if you want to store it someplace you should use S3.

Setting up the services

Setting up the services

With your newly created amazon educate account head over to <https://www.awseducate.com/signin/SiteLogin> and login. Once you have logged on press the “AWS account” button on the top right corner of the webpage and the orange “AWS educate starter account” orange button afterwards.

Setting up the services



My Classrooms

Portfolio

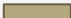

Career Pathways

Badges

Jobs

AWS Account

Logout




Consecutive Days: 1

Pathways Completed: 0

Badges Earned: 0

Preferred Language:
English

Cloud technology is everywhere, creating over 18 million cloud jobs worldwide (source: Worknet Analytics). AWS Educate introduces you to knowledge cloud enabled careers through more than 25 learning pathways, each with content from industry professionals, learning activities and labs, opportunities to earn AWS Educate Badges and Certificates of Completion, and access to the AWS Educate Job Board. Co-opted with courses at your school or through online providers, AWS Educate puts you on the pathway to your dream job in the clouds. Begin your journey today!



If you missed out the "Optimizing your AWS Educate Profile to Help You Find a Cloud Career" webinar and Q&A session, watch it here!

● ○ ○ ○ ○

Suggested Jobs

Solution Analyst - Deloitte Consulting, US Delivery Center
Deloitte Consulting
[more about this opportunity](#)

Solution Analyst - Deloitte Consulting, US Delivery Center
Deloitte Consulting
[more about this opportunity](#)

Solution Analyst - Deloitte Consulting, US Delivery Center
Deloitte Consulting
[more about this opportunity](#)

Deloitte Consulting Solutions Engineering Analyst
Deloitte Consulting
[more about this opportunity](#)

Intern - Software Engineer
Fannie Mae
[more about this opportunity](#)

See More

Setting up the services

My Classrooms

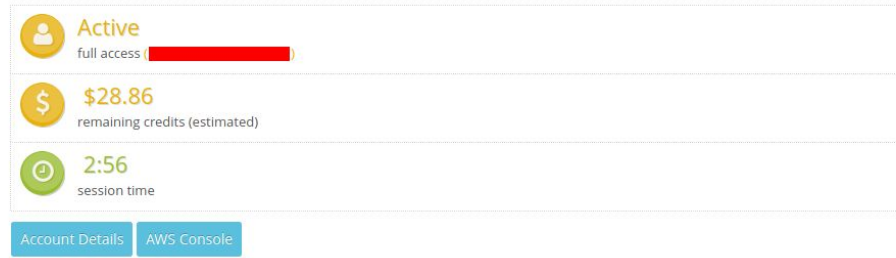
View your list of Classroom invitations and accept or decline the invitation. Access a Classroom by clicking Go to my classroom.

Course Name ↕	Description	Educator ↕	Course End Date ↕	Credit Allocated Per Student ↕	Status
Software Systems and Technologies for Big Data Applications	This is an introductory course to Big Data technologies. The course studies a series of problems, such as, distributed ETL and Machine Learning at scale, distributed data representation, columnar storage, NoSQL databases, eventual consistency, real time analytics (streaming), and cluster management	Christos Kozanitis	06/30/2021	\$40	Accepted Go to classroom ➔

Setting up the services

You should see a screen containing the following:

Your AWS Account Status



The screenshot displays the 'Your AWS Account Status' page. It features three rows of status information, each with a circular icon on the left and text on the right. The first row shows a person icon, the word 'Active', and a red progress bar with the text 'full access' below it. The second row shows a dollar sign icon, the amount '\$28.86', and the text 'remaining credits (estimated)' below it. The third row shows a clock icon, the time '2:56', and the text 'session time' below it. At the bottom of the screenshot, there are two blue buttons: 'Account Details' and 'AWS Console'.

Active	full access
\$28.86	remaining credits (estimated)
2:56	session time

Account Details AWS Console

Pressing the “Account Details” > “show” next to Amazon CLI will show you your aws credentials. We will need them later. Keep in mind these change upon **LOGIN**.

Pressing “AWS Console” will redirect you to the Amazon dashboard where you access all the services amazon provides.

Setting up the services

The screenshot displays the AWS Management Console interface. At the top, the navigation bar includes the AWS logo, 'Services', 'Resource Groups', and a star icon. On the right, it shows a user profile 'vocstartsoft/user(242527-...)', the region 'N. Virginia', and a 'Support' link.

The main content area is titled 'AWS Management Console'. It is divided into several sections:

- AWS services**:
 - Find Services**: A search bar with the placeholder text 'You can enter names, keywords or acronyms.' and an example 'Example: Relational Database Service, database, RDS'.
 - Recently visited services**: A list of icons for EC2, EMR, and S3.
 - All services**: A link to view all available services.
- Build a solution**: A section with the subtitle 'Get started with simple wizards and automated workflows.' containing eight cards:
 - Launch a virtual machine**: With EC2, 2-3 minutes. Icon: server rack.
 - Build a web app**: With Elastic Beanstalk, 6 minutes. Icon: cloud with people.
 - Build using virtual servers**: With Lightsail, 1-2 minutes. Icon: cloud with a plus sign.
 - Register a domain**: With Route 53, 3 minutes. Icon: shield with '53'.
 - Connect an IoT device**: With AWS IoT, 5 minutes. Icon: person with a cloud.
 - Start migrating to AWS**: With CloudEndure Migration, 1-2 minutes. Icon: cloud with an arrow.
 - Start a development project**: With CodeStar, 5 minutes. Icon: cloud with a plus sign.
 - Deploy a serverless microservice**: With Lambda, API Gateway, 2 minutes. Icon: cloud with a plus sign.
- Access resources on the go**: A section with a mobile app icon and text: 'Access the Management Console using the AWS Console Mobile App. [Learn more](#)'.
- Explore AWS**: A section with a link to 'Free Digital Training' and text: 'Get access to 350+ self-paced online courses covering AWS products and services. [Learn more](#)'.
- AWS IQ**: A section with text: 'Connect with AWS Certified third-party experts for on-demand consultations and project help. [Get started](#)'.
- Amazon GuardDuty**: A section with text: 'Protect your AWS accounts and workloads with intelligent threat detection. [Learn more](#)'.
- Amazon DynamoDB**: A section with text: 'Want more scale? Try a serverless NoSQL database service for your modern application. [Get started](#)'.
- Have feedback?**: A section with an email icon and text: 'Submit feedback to tell us about your experience with the AWS Management Console.'

Setting up a key pair

1. Go to AWS management console, and select EC2
2. Click on key-pairs
3. On top right corner press "Create key pair"
4. Enter a name for it. Feel free to use either ppk or pem, but i would suggest using pem if you have linux.
5. When you click "Create key pair" a new key will be downloaded automatically.
 - In order to use that key it needs to have 700 permissions. (chmod 700 key.pem)

Elastic Map Reduce - EMR

- Go to AWS Management Console and select EMR.
- Press Create Cluster > advanced options.

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name

☒ Logging ⓘ

S3 folder

Launch mode ☒ Cluster ⓘ ☐ Step execution ⓘ

Software configuration

Release ⓘ

Applications ☒ Core Hadoop: Hadoop 2.10.1, Hive 2.3.7, Hue 4.8.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.2

☐ HBase: HBase 1.4.13, Hadoop 2.10.1, Hive 2.3.7, Hue 4.8.0, Phoenix 4.14.3, and ZooKeeper 3.4.14

☐ Presto: Presto 0.240.1 with Hadoop 2.10.1 HDFS and Hive 2.3.7 Metastore

☐ Spark: Spark 2.4.7 on Hadoop 2.10.1 YARN and Zeppelin 0.8.2

☐ Use AWS Glue Data Catalog for table metadata ⓘ

Hardware configuration

Instance type ⓘ The selected instance type adds 64 GiB of GP2 EBS storage per instance by default. [Learn more](#)

Number of instances (1 master and 2 core nodes)

Cluster scaling ☐ scale cluster nodes based on workload

Security and access

EC2 key pair ⓘ [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 ⓘ

EC2 instance profile ⓘ

Cancel

Create cluster

Elastic Map Reduce - EMR

- In step 1: software and steps
 - Leave release on default.
 - Select Hadoop and spark. (ganglia and zeppelin might be useful to you but not required)
 - No need to change anything else
- In Step 2: Hardware

Elastic Map Reduce - EMR

- In Step 2: Hardware

Change the specs of each machine.

Cluster Nodes and Instances

Choose the instance type, number of instances, and a purchasing option. [Learn more about instance purchasing options](#)

Console options for automatic scaling have changed. [Learn more](#)

Node type	Instance type	Instance count	Purchasing option
Master Master - 1	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	1 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Core Core - 2	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	2 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Task Task - 3	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	0 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

Elastic Map Reduce - EMR

- In Step 2: Hardware

Change Space. SSDs are costly and not really needed use magnetics 30 / 40 GB should be ok

Cluster Nodes and Instances

Choose the instance type, number of instances, and a purchasing option. [Learn more about instance purchasing options](#)

Console options for automatic scaling have changed. [Learn more](#)

Node type	Instance type	Instance count	Purchasing option
Master Master - 1	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	1 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Core Core - 2	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	2 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Task Task - 3	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	0 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

Elastic Map Reduce - EMR

- In Step 2: Hardware

- Consider using spot instead of on-demand
- Lower cost, same machines.
- If the cost goes up you are kicked out effective immediately
- Better for small tests than paying full price

Cluster Nodes and Instances

Choose the instance type, number of instances, and a purchasing option. [Learn more about instance purchasing options](#)

Console options for automatic scaling have changed. [Learn more](#)

Node type	Instance type	Instance count	Purchasing option
Master Master - 1	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	1 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Core Core - 2	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	2 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price
Task Task - 3	m5.xlarge 4 vCore, 16 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	0 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

Elastic Map Reduce - EMR

- In Step 2: Hardware
 - Change EBS root Volume to 30 GB for starters, keep in mind you might have to increase that.

Elastic Map Reduce - EMR

- In Step 3: General Cluster Settings
 - Provide a name for your cluster.
 - Turn off logging, not needed.
 - In S3 folder, check your S3 folder.
 - No need for any tags.

Elastic Map Reduce - EMR

- In Step 4: Security
 - From the drop-down menu, select your previously created key.
 - Press create Cluster.

Elastic Map Reduce - EMR

- You should see this screen and your cluster booting up

Filter: All clusters <input type="text" value="Filter clusters ..."/> 2 clusters (all loaded) 						
	Name	ID	Status	Creation time (UTC+2)	Elapsed time	Normalized instance hours
<input type="checkbox"/>	hy543-cluster	j-1SQLBIW9Y1JCG	Terminated User request	2021-03-03 18:17 (UTC+2)	2 hours, 11 minutes	48

- Click on your cluster.
- Click on the link next to “security groups for master”, it should open a new tab with:

Security Groups (2) [Info](#)

Q

Filter security groups

search: sg-07f8cbf249f14b087

X

Clear filters

↶

1

↷

☐

Name

☐

-

☐

-

Security group ID

sg-07f8cbf249f14b087

sg-0afcc92a68877c149

Security group name

ElasticMapReduce-mas...

ElasticMapReduce-slave

VPC ID

vpc-800ea4fd

vpc-800ea4fd

Description

Master group for Elasti..

Slave group for Elastic ...

Owner

639883095772

639883095772

Inbound

21 Permis

6 Permissi

Click on the security group ID of the master

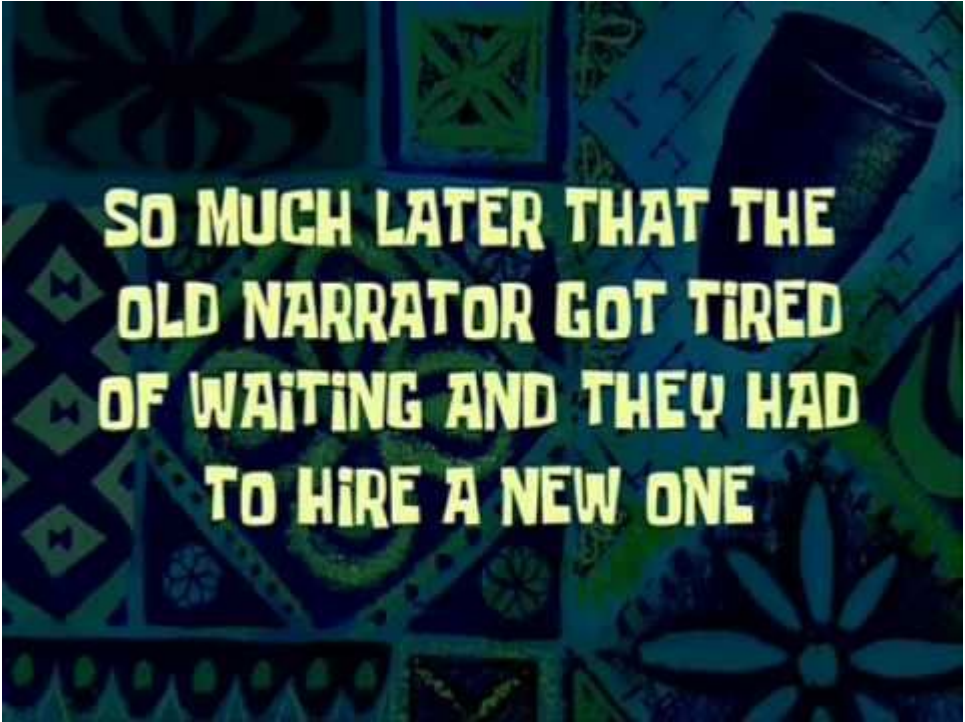
Elastic Map Reduce - EMR

- Scroll down a bit, you should see a bunch of Inbound rules. Press the “edit inbound rules”
- Press add rule
 - Select SSH
 - On the source drop-down menu select either your ip or anywhere. (to avoid any issues I would suggest anywhere.)
- Press save rules.

Elastic Map Reduce - EMR

- In order to connect to your machine cluster, first wait for it to boot up.

Elastic Map Reduce - EMR



**SO MUCH LATER THAT THE
OLD NARRATOR GOT TIRED
OF WAITING AND THEY HAD
TO HIRE A NEW ONE**

Elastic Map Reduce - EMR

- In your clusters info you should see the following:

Cluster: hy543-cluster **Waiting** Cluster ready after last step completed.

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-GT797F6MY5UT

Creation date: 2021-03-05 12:43 (UTC+2)

Elapsed time: 19 minutes

After last step completes: Cluster waits

Termination protection: Off [Change](#)

Tags: [View All / Edit](#)

Master public DNS: ec2-3-81-28-91.compute-1.amazonaws.com [🔗](#)

Connect to the Master Node Using SSH

Configuration details

Release label: emr-5.32.0

Hadoop distribution: Amazon

Applications: Spark 2.4.7, Zeppelin 0.8.2

Log URI: s3://aws-logs-639883095772-us-east-1/elasticmapreduce/ [📁](#)

EMRFS consistent view: Disabled

Custom AMI ID: --

Application user interfaces

Persistent user interfaces [🔗](#): [Spark history server](#), [YARN timeline server](#)

On-cluster user interfaces [🔗](#): Not Enabled [Enable an SSH Connection](#)

Network and hardware

Availability zone: us-east-1f

Subnet ID: [subnet-ad2c7fa3](#) [🔗](#)

Master: **Running** 1 m5.xlarge

Core: **Running** 1 m5.xlarge

Task: --

Cluster scaling: Not enabled

Security and access

Key name: hy543-emr

EC2 Instance profile: EMR_EC2_DefaultRole

EMR role: EMR_DefaultRole

Visible to all users: All [Change](#)

Security groups for Master: [sg-07f8cbb249f14b087](#) [🔗](#) (ElasticMapReduce-master)

Security groups for Core & Task: [sg-0afcc92a68877c149](#) [🔗](#) (ElasticMapReduce-slave)

Elastic Map Reduce - EMR

- Pressing that hyperlink you can see how you can ssh to your machine.

Connect to the Master Node Using SSH

You can connect to the Amazon EMR master node using SSH to run interactive queries, examine log files, submit Linux commands, and so on.

[Learn more](#) .

Windows

Mac / Linux

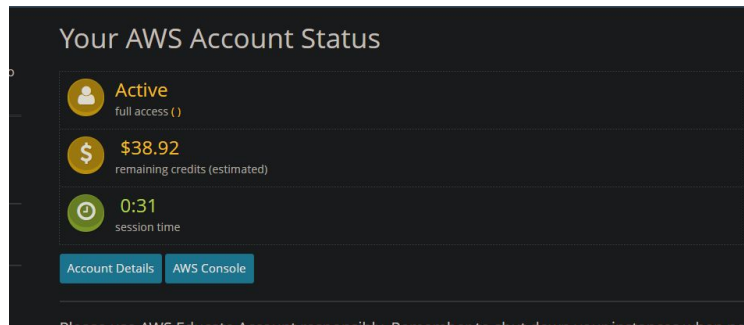
1. Open a terminal window. On Mac OS X, choose Applications > Utilities > Terminal. On other Linux distributions, terminal is typically found at Applications > Accessories > Terminal.
2. To establish a connection to the master node, type the following command. Replace ~/hy543-emr.pem with the location and filename of the private key file (.pem) used to launch the cluster.

```
ssh -i ~/hy543-emr.pem hadoop@ec2-3-81-28-91.compute-1.amazonaws.com
```

3. Type yes to dismiss the security warning.

Elastic Map Reduce - EMR

- Remember this ? ----->
- Click on account details>AWS CLI>show
- Copy that in your ~/.aws/credentials file when you ssh to the cluster master server.



EMR - S3 connectivity

Once you have set up your credentials in .aws folder you can then access your S3 bucket contents. For example you can run:

```
aws s3 ls
```

Which will show you a list from your buckets or

```
aws s3 cp s3://bucket-name/your-file ~/ \
```

In order to copy files from your S3 bucket to your EC2 machine

Amazon Services Pricing

Amazon services are not free of charge. With your educate account you get 40\$ to use, ration them well. Before starting VMs remember to choose a low cost location (e.d North Virginia). You can look up each service cost online:

EMR -> <https://aws.amazon.com/emr/pricing/>

S3 -> <https://aws.amazon.com/s3/pricing/>

Amazon Services Pricing

Explore various regions and find the minimum cost one!

Region: US East (Ohio) ▾					
	vCPU	ECU	Memory (GiB)	Instance Storage (GB)	Linux/UNIX Usage
General Purpose - Current Generation					
a1.medium	1	N/A	2 GiB	EBS Only	\$0.0255 per Hour
a1.large	2	N/A	4 GiB	EBS Only	\$0.051 per Hour
a1.xlarge	4	N/A	8 GiB	EBS Only	\$0.102 per Hour
a1.2xlarge	8	N/A	16 GiB	EBS Only	\$0.204 per Hour
a1.4xlarge	16	N/A	32 GiB	EBS Only	\$0.408 per Hour
a1.metal	16	N/A	32 GiB	EBS Only	\$0.408 per Hour
t3.nano	2	Variable	0.5 GiB	EBS Only	\$0.0052 per Hour
t3.micro	2	Variable	1 GiB	EBS Only	\$0.0104 per Hour

Amazon Services Pricing

S3 charges you not only based on the storage occupy but also on the data transfers.

Amazon S3 pricing

Pay only for what you use. There is no minimum fee. There are four cost components to consider when deciding on which S3 storage class best fits your data profile – storage pricing, request and data retrieval pricing, data transfer and transfer acceleration pricing, and data management features pricing.

Storage	Requests and data retrievals	Data transfer	Management and replication
<p>You pay for storing objects in your S3 buckets. The rate you're charged depends on your objects' size, how long you stored the objects during the month, and the storage class—S3 Standard, S3 Intelligent-Tiering, S3 Standard - Infrequent Access, S3 One Zone - Infrequent Access, S3 Glacier, and S3 Glacier Deep Archive, and Reduced Redundancy Storage (RRS). You pay a monthly monitoring and automation fee per object stored in the S3 Intelligent-Tiering storage class to monitor access patterns and move objects between access tiers in S3 Intelligent-Tiering.</p> <p>There are per-request ingest fees when using PUT, COPY, or lifecycle rules to move data into any S3 storage class. Consider the ingest or transition cost before moving objects into any storage class. Estimate your costs using the AWS Simple Monthly Calculator.</p> <p>Region: US East (Ohio) ▾</p>			
Storage pricing			
S3 Standard - General purpose storage for any type of data, typically used for frequently accessed data			
First 50 TB / Month			\$0.023 per GB
Next 450 TB / Month			\$0.022 per GB
Over 500 TB / Month			\$0.021 per GB

Amazon Services Pricing

Amazon S3 pricing

Pay only for what you use. There is no minimum fee. There are four cost components to consider when deciding on which S3 storage class best fits your data profile – storage pricing, request and data retrieval pricing, data transfer and transfer acceleration pricing, and data management features pricing.

Keep an eye out of
small details such as:

Storage	Requests and data retrievals	Data transfer	Management and replication
<p>You pay for all bandwidth into and out of Amazon S3, except for the following:</p> <ul style="list-style-type: none">• Data transferred in from the internet.• Data transferred out to an Amazon Elastic Compute Cloud (Amazon EC2) instance, when the instance is in the same AWS Region as the S3 bucket.• Data transferred out to Amazon CloudFront (CloudFront).			
<p>The pricing below is based on data transferred "in" and "out" of Amazon S3 (over the public Internet)†††. Transfers between S3 buckets or from Amazon S3 to any service(s) within the same AWS Region are free. You also pay a fee for any data transferred using Amazon S3 Transfer Acceleration. Learn more about AWS Direct Connect pricing.</p>			
Region: US East (Ohio) ▾			
			Price
Data Transfer IN To Amazon S3 From Internet			
All data transfer in			\$0.00 per GB
Data Transfer OUT From Amazon S3 To Internet			
Up to 1 GB / Month			\$0.00 per GB

Amazon Services Pricing

As far as the pricing is concerned, remember to look up the up to date pricing of each service and don't base your calculations only on what is shown in this tutorial !

Elastic Map Reduce - EMR

- Troubleshooting:
 - Your cluster may fail to start. This is most likely because the region you selected does not support the type of machine you selected.
 - For example, c1 medium is not available on us-east-1f
 - Follow the instructions given on the error and no further problems should occur.

Cluster: hy543-cluster2 **Terminated with errors** The requested instance type c1.medium is not supported in the requested availability zone. Learn more at https://docs.aws.amazon.com/console/elasticmapreduce/ERROR_noinstancetype

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-PV267U3PVTG7
Creation date: 2021-03-05 12:38 (UTC+2)
End date: 2021-03-05 12:38 (UTC+2)
Elapsed time: 38 seconds
After last step completes: Cluster waits
Termination protection: Off
Tags: --
Master public DNS: --

Configuration details

Release label: emr-5.32.0
Hadoop distribution: Amazon 2.10.1
Applications: Spark 2.4.7
Log URI: --
EMRFS consistent view: Disabled
Custom AMI ID: --

Application user interfaces

Persistent user interfaces: --
On-cluster user interfaces: --

Network and hardware

Availability zone: us-east-1f
Subnet ID: [subnet-ad2c7fa3](#)
Master: Terminated 1 c1.medium Spot (max \$0.060/hr)
Core: Terminated 1 c1.medium Spot (max \$0.060/hr)
Task: --
Cluster scaling: Not enabled

Security and access

Elastic Map Reduce - EMR

- Troubleshooting:
 - https://docs.aws.amazon.com/console/elasticmapreduce/ERROR_noinstancetype
- Amazon has very good troubleshooting guides and information. Don't forget to check carefully everything they throw at you.

Cluster: hy543-cluster2 **Terminated with errors** The requested instance type c1.medium is not supported in the requested availability zone. Learn more at https://docs.aws.amazon.com/console/elasticmapreduce/ERROR_noinstancetype

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-PV267U3PVTG7
Creation date: 2021-03-05 12:38 (UTC+2)
End date: 2021-03-05 12:38 (UTC+2)
Elapsed time: 38 seconds
After last step completes: Cluster waits
Termination protection: Off
Tags: --
Master public DNS: --

Configuration details

Release label: emr-5.32.0
Hadoop distribution: Amazon 2.10.1
Applications: Spark 2.4.7
Log URI: --
EMRFS consistent view: Disabled
Custom AMI ID: --

Application user interfaces

Persistent user interfaces: --
On-cluster user interfaces: --

Network and hardware

Availability zone: us-east-1f
Subnet ID: [subnet-ad2c7fa3](#)
Master: Terminated 1 c1.medium Spot (max \$0.060/hr)
Core: Terminated 1 c1.medium Spot (max \$0.060/hr)
Task: --
Cluster scaling: Not enabled

Security and access