

CLOUD ANALYTICS

- Cloud Computing Overview
- AWS Services
- Amazon EMR
- **Amazon EMR Exercise**

AWS ACCOUNT

- URL:
https://ets-apps.ucsd.edu/dse230_sp21-custom-aws/
- Log in using UCSD Active Directory credentials

EMR NOTEBOOKS

- Jupyter notebooks on Amazon EMR
- Pre-configured with following kernels:
 - PySpark
 - Python3
 - Spark
 - SparkR
- Note
 - EMR notebook contents are saved in S3, separate from EMR cluster that executes code
 - Attach cluster to notebook to execute

OUTLINE OF STEPS FOR EMR

- Create EC2 key pair
 - Security credentials for connecting to EC2 instance
- Create S3 bucket
 - Store notebooks and instance logs
- Create EMR cluster
 - Specify software configuration
 - Add step to copy data from public S3 bucket for class
 - Specify number of master and worker nodes
 - Specify EC2 key pair
- Create EMR notebook
 - Connect to cluster
 - Specify notebook location
 - Create notebook
 - Read data from your S3 bucket or from public S3 bucket

OUTLINE OF STEPS FOR EMR

- Create EC2 key pair
 - Security credentials for connecting to EC2 instance
- Create S3 bucket
 - Store notebooks and instance logs
- Create EMR cluster
 - Specify software configuration
 - Add step to copy data from public S3 bucket for class
 - Specify number of master and worker nodes
 - Specify EC2 key pair
- Create EMR notebook
 - Connect to cluster
 - Specify notebook location
 - Create notebook
 - Read data from your S3 bucket or from public S3 bucket

CREATE KEY PAIR

- Services -> EC2 (under Compute)
- Under 'Network & Security', click on 'Key Pairs'
- Click on 'Create key pair'
- Enter name. Click on 'Create key pair'
- Save <key-pair-name>.pem on your computer

EC2 > Key pairs > Create key pair

Create key pair

Key pair
A key pair, consisting of a private key and a public key, is a set of security credentials that you use to prove your identity when connecting to an instance.

Name

The name can include up to 255 ASCII characters. It can't include leading or trailing spaces.

File format
☒ pem
For use with OpenSSH
☐ ppk
For use with PuTTY

Tags (Optional)
No tags associated with the resource.

You can add 50 more tags.

OUTLINE OF STEPS FOR EMR

- Create EC2 key pair
 - Security credentials for connecting to EC2 instance
- Create S3 bucket
 - To store notebooks and instance logs
- Create EMR cluster
 - Specify software configuration
 - Add step to copy data from public S3 bucket for class
 - Specify number of master and worker nodes
 - Specify EC2 key pair
- Create EMR notebook
 - Connect to cluster
 - Specify notebook location
 - Create notebook
 - Read data from your S3 bucket or from public S3 bucket

CREATE S3 BUCKET

- Services -> S3 (under Storage)
- Bucket name
 - Name must be unique across all S3 buckets
 - At least 3 characters and no more than 63 characters long
 - Cannot contain uppercase characters or underscores
 - Must start with lowercase letter or number
 - Cannot change name after bucket has been created
 - More on bucket naming
 - <https://docs.aws.amazon.com/AmazonS3/latest/dev/BucketRestrictions.html#bucketnamingrules>
- Click 'Create Bucket'

CREATE S3 BUCKET

- Enter bucket name
- Click 'Create bucket'

Amazon S3 > Create bucket

Create bucket

Buckets are containers for data stored in S3. [Learn more](#)

General configuration

Bucket name

Bucket name must be unique and must not contain spaces or uppercase letters. [See rules for bucket naming](#)

AWS Region

US West (Oregon) us-west-2

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

Choose bucket

► Advanced settings

i After creating the bucket you can upload files and folders to the bucket, and configure additional bucket settings.

Cancel **Create bucket**

CREATE FOLDER IN BUCKET

- Create folder in S3 bucket
 - Click on bucket name
 - Select 'Create folder'
 - Enter name for folder. Click 'Create folder'
 - Create folder for EMR notebooks called 'Notebooks'

The screenshot shows the Amazon S3 console interface for an empty bucket. At the top, it says 'Objects (0)' and provides a brief explanation of objects. Below this is a toolbar with buttons for 'Refresh', 'Delete', 'Actions', 'Create folder', and 'Upload'. A search bar is present with the placeholder text 'Find objects by prefix'. To the right of the search bar are navigation controls showing page 1 of 1 and a settings gear icon. Below the toolbar is a table header with columns: 'Name', 'Type', 'Last modified', 'Size', and 'Storage class'. The main content area displays 'No objects' and the message 'You don't have any objects in this bucket.' At the bottom center, there is an 'Upload' button.

Objects (0)

Objects are the fundamental entities stored in Amazon S3. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

1

	Name	Type	Last modified	Size	Storage class
--	------	------	---------------	------	---------------

No objects

You don't have any objects in this bucket.

UPLOAD DATA FILES TO S3

- Services -> S3
- Click into bucket and folder
- Click 'Upload'
- Select file to upload
 - 2 ways
 - Drag and Drop
 - Click on 'Add files' and select file
- Click 'Upload'
- Upload progress is shown at bottom of webpage.
 - Large files may take several minutes.

OUTLINE OF STEPS FOR EMR

- Create EC2 key pair
 - Security credentials for connecting to EC2 instance
- Create S3 bucket
 - Upload data to bucket
- Create EMR cluster
 - Specify software configuration
 - Specify number of master and worker nodes
 - Specify EC2 key pair
- Create EMR notebook
 - Connect to cluster
 - Specify notebook location
 - Create notebook
 - Read data from your S3 bucket or from public S3 bucket

CREATE CLUSTER


- Services -> EMR
- Click on 'Create cluster'
- Click on 'Go to advanced options'

[Go to advanced options](#)  **Go to advanced options**

General Configuration

Cluster name



☒ **Logging** ⓘ

S3 folder 

Launch mode ☒ **Cluster** ⓘ ☐ **Step execution** ⓘ

CREATE CLUSTER - SOFTWARE CONFIGURATION

Software Configuration

Release  

Select emr-6.2.0
Then select software options shown here

<input checked="" type="checkbox"/> Hadoop 3.2.1	<input type="checkbox"/> Zeppelin 0.9.0	<input type="checkbox"/> Livy 0.7.0
<input type="checkbox"/> JupyterHub 1.1.0	<input type="checkbox"/> Tez 0.9.2	<input type="checkbox"/> Flink 1.11.2
<input type="checkbox"/> Ganglia 3.7.2	<input type="checkbox"/> HBase 2.2.6-amzn-0	<input type="checkbox"/> Pig 0.17.0
<input type="checkbox"/> Hive 3.1.2	<input type="checkbox"/> Presto 0.238.3	<input type="checkbox"/> PrestoSQL 343
<input type="checkbox"/> ZooKeeper 3.4.14	<input checked="" type="checkbox"/> JupyterEnterpriseGateway 2.1.0	<input type="checkbox"/> MXNet 1.7.0
<input type="checkbox"/> Sqoop 1.4.7	<input type="checkbox"/> Hue 4.8.0	<input type="checkbox"/> Phoenix 5.0.0
<input type="checkbox"/> Oozie 5.2.0	<input checked="" type="checkbox"/> Spark 3.0.1	<input type="checkbox"/> HCatalog 3.1.2
<input type="checkbox"/> TensorFlow 2.3.1		

CREATE CLUSTER - CONFIGURE NODES

- Specify instance count for Master and Core (worker)
- <https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-master-core-task-nodes.html>

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: 64 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: 64 GiB Add configuration settings	<input type="text" value="1"/> 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: 64 GiB Add configuration settings	<input type="text" value="0"/> Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot Use on-demand as max price

CREATE CLUSTER - NAME AND EC2 KEY PAIR

- Enter name for 'Cluster name'
- Click Next

General Options

Cluster name

☒ **Logging** ⓘ

S3 folder 

☐ **Log encryption** ⓘ

☒ **Debugging** ⓘ

☒ **Termination protection** ⓘ

- Select EC2 key pair created earlier
- Click 'Create cluster'


Security Options

EC2 key pair ⓘ

☒ **Cluster visible to all IAM users in account** ⓘ

CREATE CLUSTER

- Creating cluster may take several minutes
 - To provision resources for master and worker nodes
- Click on Steps tab to monitor status of cluster creation and view logs
- Cluster is ready when Status shows 'Waiting. Cluster ready'

Create cluster		View details	Clone	Terminate
Filter:		All clusters ▼	Filter clusters ...	7 clusters (all loaded) ↻
		Name	ID	Status
<input type="checkbox"/>	▶	 EMR-cluster	j-QTCMAXIJGA2O	Waiting Cluster ready

OUTLINE OF STEPS FOR EMR

- Create EC2 key pair
 - Security credentials for connecting to EC2 instance
- Create S3 bucket
 - Upload data to bucket
- Create EMR cluster
 - Specify software configuration
 - Add step to copy data from public S3 bucket for class
 - Specify number of master and worker nodes
 - Specify EC2 key pair
- Create EMR notebook
 - Connect to cluster
 - Specify notebook location
 - Create notebook
 - Read data from your S3 bucket or from public S3 bucket

CREATE EMR NOTEBOOK

- Services -> EMR
- Click on Notebooks in menu on left
- Click on 'Create notebook'



The screenshot shows the Amazon EMR console interface. On the left sidebar, under 'Amazon EMR', the 'Notebooks' option is selected and highlighted with an orange bar. The main content area is titled 'Notebooks' and contains a description: 'Use EMR notebooks based on Jupyter to analyze data interactively with live code, narrative text clusters running Hadoop, Spark, and Livy. Notebooks run free of charge and are saved in Amazon S3. [Learn more](#)'. Below the description are four buttons: 'Create notebook' (blue), 'View details', 'Open in JupyterLab', and 'Open in Jupyter'. At the bottom, there is a filter bar with a dropdown menu set to 'All notebooks', a search input field with the placeholder 'Filter notebooks ...', and a status indicator showing '2 notebooks (all loaded)' with a refresh icon.

CREATE EMR NOTEBOOK

Create notebook

Name and configure your notebook

Name your notebook, choose a cluster or create one, and customize configuration options if desired. [Learn more](#)

Notebook name*

Names may only contain alphanumeric characters, hyphens (-), or underscores (_).

Description

256 characters max.

Cluster* ☒ Choose an existing cluster

☐ Create a cluster ?

Security groups ☒ Use default security groups ?

☐ Choose security groups

AWS service role* ?

Notebook location* Choose an S3 location where files for this notebook are saved.

☐ Use a location that EMR creates ?

☒ Choose an existing S3 location in us-west-2

▶ Git repository

▶ Tags ?

* Required

Add name for notebook

Choose cluster previously created

Choose Notebooks folder of S3 bucket previously created

OUTLINE OF STEPS FOR EMR

- Create EC2 key pair
 - Security credentials for connecting to EC2 instance
- Create S3 bucket
 - Upload data to bucket
- Create EMR cluster
 - Specify software configuration
 - Add step to copy data from public S3 bucket for class
 - Specify number of master and worker nodes
 - Specify EC2 key pair
- Create EMR notebook
 - Connect to cluster
 - Specify notebook location and create notebook
 - Read data from S3 bucket

WORKING WITH EMR NOTEBOOKS

- Wait for Notebook status to be 'Pending' or 'Ready'
- Click 'Open in JupyterLab'
- In JupyterLab,
 - Double click notebook to open it
 - Select 'PySpark' as kernel (not 'Python 3')

Select Kernel

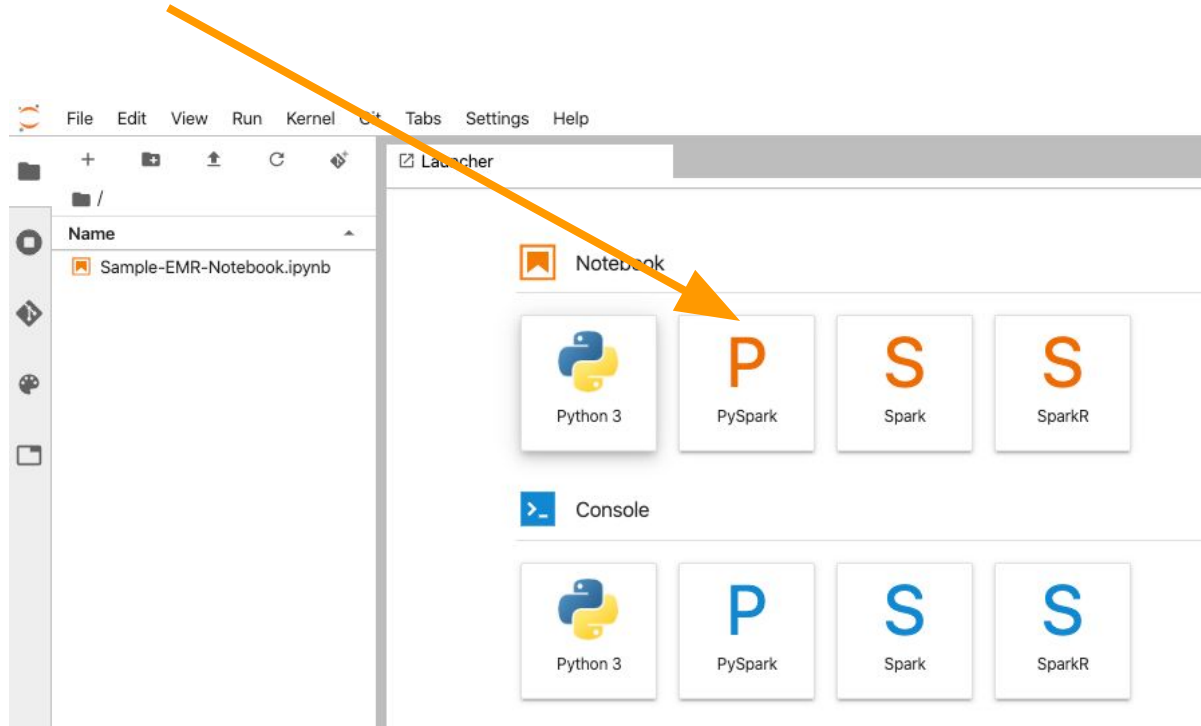
Select kernel for: "Sample-EMR-Notebook.ipynb"

PySpark ▼

Select

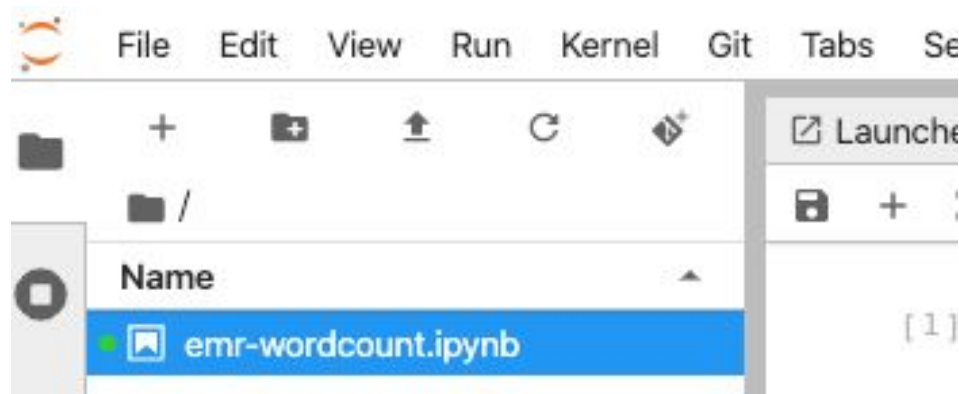
WORKING WITH EMR NOTEBOOKS

- To create new notebook from within JupyterLab
 - Select 'PySpark' as kernel (not 'Python 3')



EMR NOTEBOOK

- Rename notebook from 'Untitled.ipynb'
- Entering code
 - Can either copy and paste code from notebook shell or upload file directly.
 - Ensure kernel is PySpark



RUNNING EMR NOTEBOOK

- Enter code in notebook, then run as normal
- Print versions

```
import pyspark
print(spark.version, pyspark.version.__version__)
```

- Read data

```
dataFileName = "<file-location>"
textDF = spark.read.text(dataFileName).cache()
```

- Get number of rows

```
textDF.count()
```

- Stop Spark

```
spark.stop()
```

- File location -

<https://dse230-emr.s3-us-west-1.amazonaws.com/Shakespeare.txt>

DOWNLOAD NOTEBOOK AND RESULTS

- Download notebook from JupyterLab
 - File -> Download as Notebook (.ipynb)
 - File -> Download as HTML (.html)
- Download results from S3
 - Navigate to bucket
 - Navigate to folder
 - Right click on file to download


CLEANING UP

- Do this when done with current session
 - Important: Delete cluster and notebook to avoid accumulating fees!
 - **Note: There is a daily limit of \$5.00 and total course limit of \$50.00. You will not be able to run anything if you exceed these limits!**
- In JupyterLab
 - File -> Close and Shutdown Notebook
- In EMR
 - Click on Notebooks in the left menu
 - Select Notebook, then click Stop
 - Click on Clusters in the left menu
 - Select cluster and Click Terminate to delete cluster
 - Wait until cluster status is Terminated.
- S3 - Can leave S3 bucket up until end of course

STOP NOTEBOOK

- Cleaning up: Make sure notebook is stopped

Notebooks

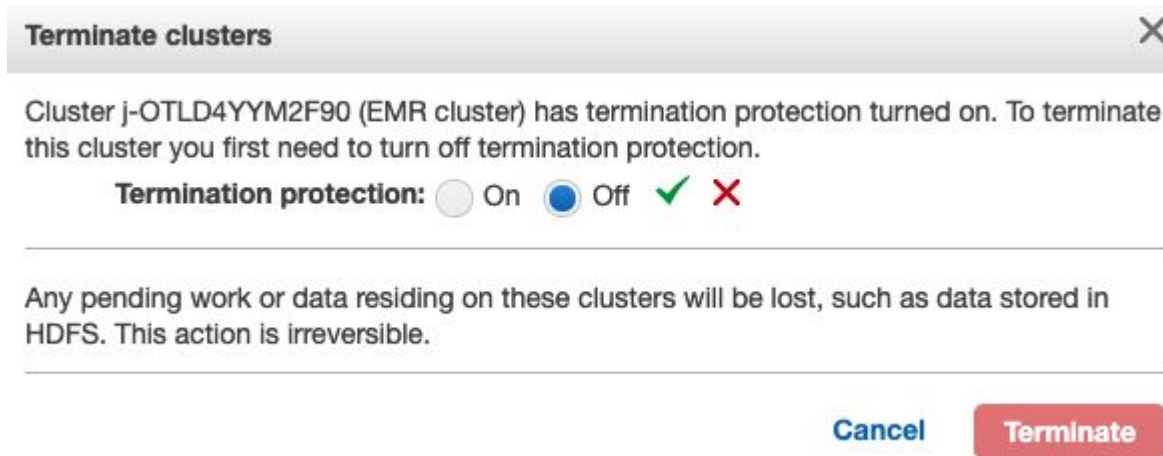
Use EMR notebooks based on Jupyter to analyze data interactively with live code, narrative text, visualizations, and more. Create and attach notebooks to Amazon EMR clusters running Hadoop, Spark, and Livy. Notebooks run free of charge and are saved in Amazon S3 independently of clusters. Standard billing for clusters and Amazon S3 apply. [Learn more](#) 

[Create notebook](#) [View details](#) [Open in JupyterLab](#) [Open in Jupyter](#) [Start](#) [Stop](#) [Delete](#)

Filter: All notebooks Filter notebooks ... 2 notebooks (all loaded) 					
	Name	Status	Cluster	Creation time (UTC-8)	Last modified 
	EMR-Notebook	Stopped	j-332GN8XPDB86V	2021-02-26 19:58 (UTC-8)	20 seconds ago

TERMINATE CLUSTER

- Select cluster. Click 'Terminate'
- In popup window
 - Click on 'Change' for Termination protection
 - Click on 'Off'. Then click on checkmark
 - Click on 'Terminate'
- Check that cluster's status is 'Terminated'



WORK ON EXISTING NOTEBOOK

- Start another cluster
 - Select 'Clusters' on the left
 - Select previous cluster and click Clone
 - Click 'Create Cluster'
- Open notebook in new cluster
 - Select 'Notebooks' on the left
 - Select 'Change cluster'
 - Select 'Choose an existing cluster' and select new cluster
 - Select 'Change cluster and start notebook'