

AI 620 Emerging Topics in Artificial Intelligence

HOS06A Inference Pipeline in SageMaker

03/14/2023 Developed by Yared Shewarade

09/26/2024 Updated by Anh Nguyen

10/7/2024 Reviewed by Jonathan Koerber

School of Technology and Computing (STC) @City University of Seattle (CityU)

Before You Start

- The directory path shown in screenshots may be different from yours.
- Some steps are not explained in the tutorial. If you are not sure what to do:
 1. Consult the resources listed below.
 2. If you cannot solve the problem after a few tries, ask a student worker for help.

Learning Outcomes

Students will be able to learn:

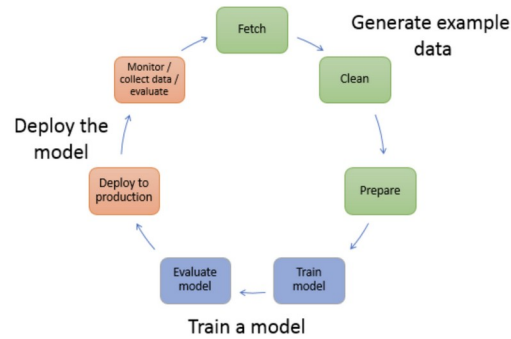
- Introduction to Amazon SageMaker
- Preprocessing big data through Spark EMR

Resources

- Tripuraneni, S., & Song, C. (2019). *Hands-on artificial intelligence on amazon web services: Decrease the time to market for AI and ML applications with the power of AWS* (1st ed.). Packt.

Introduction to Amazon SageMaker

Amazon SageMaker is a fully managed service that enables quick and easy integration of machine learning-based models into applications to generate inferences in real time and at scale. This section provides an overview of machine learning and explains how SageMaker works. If you are a first-time user of SageMaker, we recommend that you read the following sections in order:



- Fetch the data**— Pull the dataset or datasets into a single repository.
- Clean the data**— Inspect the data and clean it as needed.
- Prepare or transform the data**— Perform additional data transformations.
- Training the model**— Train a model, you need an algorithm or a pre-trained base model. The algorithm you choose depends on several factors. For a quick, out-of-the-box solution, you might be able to use one of the algorithms that SageMaker provides.
- Evaluating the model**—After you have trained your model, you evaluate it to determine whether the accuracy of the inferences is acceptable.

Preprocessing big data through Spark EMR

Wrangling a big dataset in Jupyter notebooks results in out-of-memory errors. Our solution is to employ AWS EMR (Elastic MapReduce) clusters to conduct distributed data processing.

Amazon Elastic MapReduce (Amazon EMR) is a web service that makes it easy to process big data. Hadoop will be used as the underlying distributed filesystem while Spark will be used as the distributed computing framework. It provides a managed notebook environment, based on Jupyter Notebook which can be used to interactively wrangle large data, visualize the same, and prepare analytics-ready datasets. These EMR notebooks can also be saved periodically to a persistent data store, S3, so the saved work can be retrieved later.

Note: For submission, take the screenshot for all steps and save it in your local repository along with your code.

A. Prepare input data on Amazon S3

1. Sign in to your AWS Management Console. Go to S3 and create a bucket named **“hos06emrnotebook-<yourname>”**.

Amazon S3 > Buckets > Create bucket

Create bucket [Info](#)

Buckets are containers for data stored in S3.

General configuration

AWS Region
US East (N. Virginia) us-east-1

Bucket type [Info](#)

☒ **General purpose**
 Recommended for most use cases and access patterns. General purpose buckets are the original S3 bucket type. They allow a mix of storage classes that redundantly store objects across multiple Availability Zones.

☐ **Directory**
 Recommended for low-latency use cases. These buckets use only the S3 Express One Zone storage class, which provides faster processing of data within a single Availability Zone.

Bucket name [Info](#) Add your name to make the bucket name unique

hos06emrnotebook-honganh

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

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

[Choose bucket](#)

Format: s3://bucket/prefix

- Download [books.csv](#) and [ratings.csv](#) from [here](#) and upload it to the above created S3 bucket.

Amazon S3 > Buckets > hos06emrnotebook-honganh

hos06emrnotebook-honganh [Info](#)

[Objects](#) | [Properties](#) | [Permissions](#) | [Metrics](#) | [Management](#) | [Access Points](#)

Objects (2) [Info](#)

[Copy S3 URI](#)
[Copy URL](#)
[Download](#)
[Open](#)
[Delete](#)
[Actions](#)
[Create folder](#)
[Upload](#)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

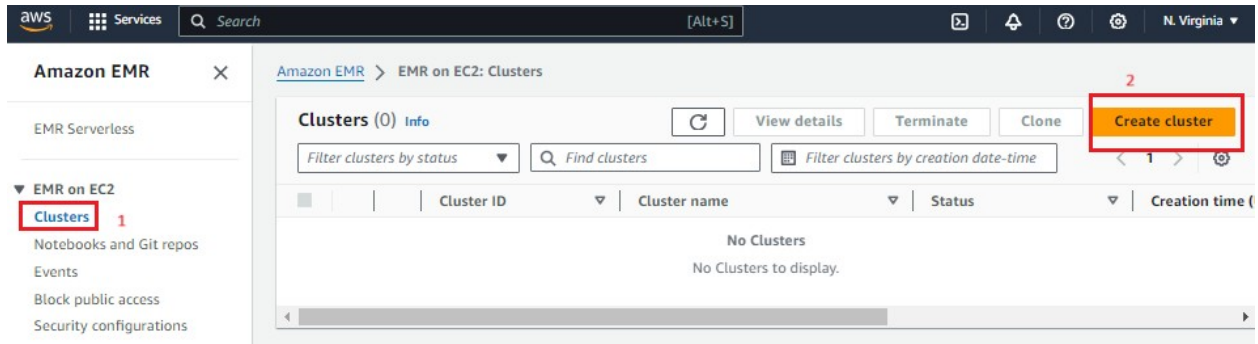
Find objects by prefix

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	books.csv	csv	September 22, 2024, 12:13:37 (UTC-07:00)	3.1 MB	Standard
<input type="checkbox"/>	ratings.csv	csv	September 22, 2024, 12:13:39 (UTC-07:00)	11.9 MB	Standard

****Screenshot Of Bucket. *****

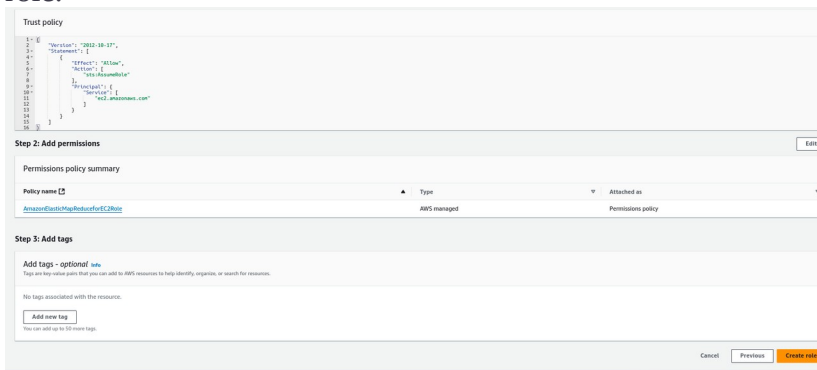
B. Launch an Amazon EMR Cluster

- Under EMR on EC2 in the left navigation pane, choose **Clusters**, and then choose **Create cluster**.



2. Under **Cluster termination and node replacement**, make sure Automatically terminate cluster after idle time is chosen. The idle time should be 1 hour.
3. Under Identity and Access management (IAM) roles:
 - Choose Create a service role
 - In EC2 instance profile for Amazon EMR:
 - Choose an existing instance profile
 - Instance profile: open the dropdown and choose the created profile “ecsInstanceRole”.

****If you have not created an IAM role with EMR default access you can create one.** Navigate to the AWS Identity and Access Management (IAM) by search IAM in the search bar. In the left gutter select Roles. On the left side of the page click create role. This will be for an AWS service so select that on the next screen enter EMR into service or use case you. Select EMR Role for EC2 click next and next again adding AmazonElasticMapReducerforEC2Role permission and name your role. Then create a role.



▼ Identity and Access Management (IAM) roles - required [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. **1**

☐ Choose an existing 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.

☒ Create a service role
Let Amazon EMR create a new service role so that you can grant and restrict access to resources in other AWS services.

Networking resources
We've already added the resources that you configured in the Networking section. Choose the VPC, subnet, and security groups that the service role can access.

Virtual Private Cloud (VPC)
Choose one or more VPCs
- vpc-0d0f8d11d20dce092

Subnet
Choose one or more subnets
- subnet-01a0eb9544622ae92

Security group
Choose one or more security groups

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
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.

☐ Create an instance profile
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
ecsInstanceRole **2**

Custom automatic scaling role - optional
When a custom automatic scaling role triggers, Amazon EMR assumes this role to add and terminate EC2 instances. [Learn more](#)

Custom automatic scaling role
Choose IAM role

Create IAM role [↗](#)

- Click **the Create Cluster** button. Wait for the status to change from Starting to Running to Waiting. cluster status changes to Waiting when the cluster is up, running, and ready to accept work.

You may need to choose the refresh icon on the right or refresh your browser to see status updates.

****Screenshot of running Cluster****

C. Create Workspace (previously called Notebook)

- Go to your S3 bucket **"hos06emrnotebook-<yourname>"**
- Go back to Amazon EMR, under EMR on EC2, select Notebooks and Git repos > Go to Workspaces (Notebooks).

Click Create Studio. Leave the configurations as default and click Create studio and launch Workspace.

Amazon EMR > EMR Studio: Studios > Create Studio

Create a Studio Info

Setup options Info

☒ Interactive workloads
 ☐ Batch jobs
 ☐ Custom

Studio settings Info Edit

Studio name
Studio_1

S3 location for Workspace storage
We'll create a new bucket and use the location `s3://aws-emr-studio-730335403677-us-east-1/1727378704632`.

Service role to let Studio access your AWS resources
We'll create a new service role named `AmazonEMRStudio_ServiceRole_1727378704632`.
[View permission details](#)

Workspace settings Info Edit

Workspace name
Studio_1_Workspace_1

EMR Serverless application settings Info Edit

Application name
Serverless_Interactive_App_1727378704632

► Default settings for interactive workloads

Runtime role
We'll create a new runtime role named `AmazonEMRStudio_RuntimeRole_1727378704632`.
[View permission details](#)

Cancel
Create Studio
Create Studio and launch Workspace

- Go back to Notebooks and Git repos. You will see that there is a workspace created. Click that workspace > Attach Cluster.

Now, you can use Jupyter Notebook to write Spark code.

****Screenshot of Notebook****

D. Clean up

- Go to **Clusters**. Choose the clusters you created > Click Terminate to delete the cluster. The status will change from Terminating to Terminated.

Amazon EMR ×

Amazon EMR > EMR on EC2: Clusters

EMR Serverless

▼ EMR on EC2

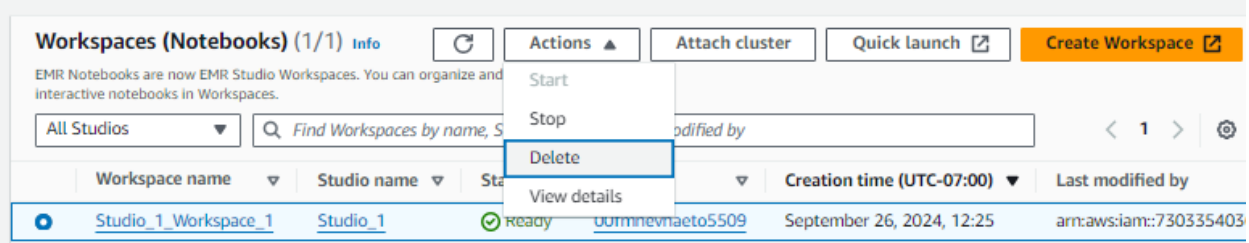
- Clusters**
- Notebooks and Git repos
- Events
- Block public access
- Security configurations

Clusters (1/1) Info
Refresh
View details
Terminate
Clone
Create cluster

Filter clusters by status Find clusters Filter clusters by creation date-time < 1 >

	Cluster ID	Cluster name	Status	Creation time (UTC-07:00)
<input checked="" type="checkbox"/>	j-2JVGAUFU3LHQD5	My cluster	Waiting Ready to run steps	September 26, 2024, 12:08

- Go to **Notebooks and Git repos**. Choose the workspace just created. Under the **Action** button, choose Delete. The workspace status will go from Deleting to Deleted, and eventually will be removed from the table.



- In the search bar, type S3 and go to Amazon S3. You will see that there are several buckets created during this exercise.
 - aws-emr-studio*
 - Aws-logs*
 - Hos06emrnotebook-<yourname>

Empty and delete all these buckets.

HOS submission instructions:

- Please install the GitHub Desktop: https://cityuseattle.github.io/docs/git/github_desktop/
- Clone, organize, and submit your work through GitHub Desktop: <https://cityuseattle.github.io/docs/hoporhos>