# Using AWS Educate For CSCE 676

This document is a quick guide on how to apply for an AWS Educate account, create a Hadoop/Spark cluster, create an S3 storage bucket and how to create a Jupyter notebook to run code against the created cluster.

Thanks to Amazon for providing free credits for this course!

## Table of Contents

# Applying to AWS Educate

To apply to join AWS Educate, follow the URL given in class and choose "Texas A&M University "and "Data Mining and Analysis - CSCE 676" course from the corresponding dropdown menus as show below. Click on an email confirmation link that will be sent to you. Your application should be approved within an hour.



After your application gets approved, follow the link sent to you to set a password and to access your AWS Educate account. Once you have created your account and are able to access the starting page as shown below click on "Use an AWS Educate Starter Account" and then click "Create Starter Account"

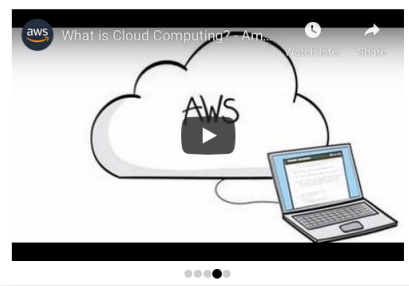**Majid Alfifi**   Consecutive Days: **1**   Pathways Completed: **0**   Badges Earned: **0**

Preferred Language:
English

Cloud technology is everywhere, creating over 18 million cloud jobs worldwide (source: Wanted Analytics). AWS Educate introduces you to lucrative 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. Coupled with courses at your school or through online providers, AWS Educate puts you on the pathway to your dream job in the clouds.
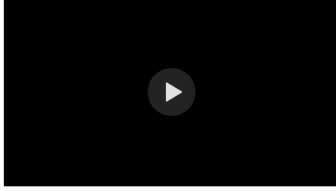
Begin your journey today!

What is Cloud Computing? - Am...

**Access AWS Services**

Use a personal AWS Account →

Use an AWS Educate Starter Account →

Which option is right for me?

Learn More →

Suggested Jobs

Portfolio   Career Pathways   Learn   Badges

---

# I'd like to use an AWS Educate Starter Account

Choose an AWS Educate Starter Account to receive access to an AWS account with a preset limit on your spend on AWS services. An AWS Educate Starter Account is run and managed by a third party (Vocareum, Inc.) and the Starter Account runs in the Vocareum's environment on AWS. Starter Accounts are subject to a separate agreement between you and Vocareum under separate terms and conditions.

The AWS Educate Starter Account provides access to most but not all AWS services. Students at AWS Educate member institution will receive up to $75 (US) of AWS credit per year in their AWS Educate Starter Account, and students at non-member institution will receive up to $30 (US) of AWS credit per year.

You don't need a credit card to use a Starter Account because AWS promotional credits are already available in the account. When your usage of AWS services exceed the balance on the account, the account is closed and any running services or other resources on the account are lost.

**Create Starter Account**

or choose another option

Clicking on "AWS Educate Starter Account" takes you to the following page which has the access to AWS Console! You can always access your account your email and the password you created above by visiting the following link:

https://www.awseducate.com/signin

# Creating A Hadoop/Spark Cluster

Click on"AWS Console" on the home page shown above. You may need to allow pop-ups for the AWS Console to open. On the AWS Management Console click on All services under "Find Services" and select EMR (Elastic MapReduce). You can also search for it in the search box.

This will open the page below. Under EMR, we will be mainly working on two tabs:
* Clusters where we will build the Hadoop/Spark cluster, and
* Notebooks where we will create a Jupyter notebook to access the cluster and operate on it.



The steps to create a Hadoop/Spark cluster are as follows:

1. Click on the Create cluster button above.
2. Click on "Go to advanced options"

3.  Chose the Hadoop, Hive, and Spark and also copy the following configuration line and paste it under "Edit software settings" a

```
classification=spark-defaults,properties=[spark.jars.packages=graphframes:graphframes:0.7.0-spark2.4-s_2.11]
```



4.  On "Step 2: Hardware" you can use the defaults (1 master node m5.xlarge, and 2 data nodes m5.xlarge). Use these for experimentation and learning but when you are ready to run your code on a bigger dataset, then you can create more core nodes (e.g. 10 nodes) which will cost more of your credits.

5.  On "Step 3: General Cluster Settings", click on "Bootstrap Actions". From the drop-down menu, choose "Custom action" and then click "Configure and add"



In the script location type `s3://caverlee/init.sh`



Click add. This script will install Jupyter on the master node and download "start-jupyter.sh" script to use later.

6. On "Step 4: Security". We need to do two things:
   1. Create an EC2 Key Pair, and
   2. Allow SSH access.
To create an EC2 Key Pair, click on the link shown and follow instructions.



After you have created the key select it in the EC2 key pair dropdown list:

To allow SSH access to your cluster click on "Create a security group".



Select the "default" security group and choose edit inbound rules as shown below



Then click "Edit Rules" -> "Add Rule" and type ~~22~~ 0-65535 in the Port Range and choose "Anywhere" for the Source. Click "Save rules". Note there is security risk but since this will live for only few hours after which you will terminate the cluster, it should be fine for this homework.

Finally in "Step 4: Security" click on "Additional security groups" links and add the "default" security group you just edited



7. Now you are ready to click "Create cluster"

The cluster will take few minutes to start.

# Accessing the Cluster with SSH

You can SSH to the master node (or any of the worker nodes if you allow SSH access) and run any unix commands you like, or investigate logs, and so on. To do so, click on the "SSH" link as shown below



This will open a window like the following that will show you the domain name of the master node so you can ssh to it

# Accessing Hadoop/Spark Web Interfaces

Hadoop/Spark provide web interfaces to investigate their status and logs. There are couple of ways to access those interfaces but the easiest maybe to use SSH tunnel and SOCKS proxy.

The steps are as follows:
1.  When connecting to the master node with SSH, us -D option as follows ssh -D 1050
    ```
    ssh -D 1050 -i csce676.pem hadoop@ec2-34-239-180-133.compute-1.amazonaws.com
    ```
2.  In your browser, forward all traffic through this port (1050) which makes your browser work as though it was running on the remote Amazon server. I use Firefox for this purpose because it's easy to setup SOCKs and I have this browser dedicated for my cluster work. In Firefox, set the connection settings as follows:



3.  Find out the IP of the master node and use it to access the web interfaces shown in the table below. (Hint: you can find the ip after you login in the prompt name. For example,
    `[hadoop@ip-172-31-62-185] $`

| Name of interface | URI |
|---|---|
| YARN ResourceManager | http://*master-ip*:8088/ |
| Hadoop HDFS NameNode | http://*master-ip*:50070/ |
| Spark HistoryServer | http://*master-ip*:18080/ |

Note: if you want to consider other ways to access the web interfaces and for a longer list of interfaces, consult the AWS documentation.

# Accessing the cluster with a Jupyter Notebook
Ignore this section (doesn't work any more). Check the last section in this guide for an alternative way.

Click on the Notebooks Tab and then click on the "Create a notebook" button:



Choose a name for your notebook and choose the cluster you created in the previous step as the cluster associated with this notebook. Finally click "create notebook".



Note: your notebook will be saved in S3 and you will always find it there. However, you should create a cluster only when you need to work on your notebook and then terminate the cluster

immediately after you are done to avoid wasting credits on an unused cluster. You can stop the notebook and leave it there for your next time to continue working where you left.



After opening the notebook, choose PySpark Kernel and you are now ready to run Spark code against your cluster!



You will find a sample notebook in Piazza Resources to get you started.

```
In [1]: tweets = spark.read.json("s3://us-congress-tweets/congress-20181003.json.gz")
```

▸ Spark Job Progress

Starting Spark application

| ID | YARN Application ID | Kind | State | Spark UI | Driver log | Current session? |
|----|---------------------|------|-------|----------|------------|------------------|
| 0 | application_1569008117578_0001 | pyspark | idle | Link | Link | ✔ |

SparkSession available as 'spark'.

```
In [3]: tweets.select("user.screen_name", "text").show()
```

▸ Spark Job Progress

```
+---------------+--------------------+
|    screen_name|                text|
+---------------+--------------------+
|    kikilezigoto|RT @namek237: - T...|
| chriswyoillini|@charliekirk11 Pl...|
|     michele5411|... she'll vote t...|
|        aneesajv|RT @SenatorDurbin...|
|  AcidRayneStorm|RT @Johnoco656060...|
|Jmooretrumpgirl|@lisamurkowski @S...|
|     AzLakeHouse|RT @ChuckGrassley...|
|        Srk1951mn|RT @RonWyden: Evi...|
|       jannsloan|RT @CheriJacobus:...|
|      burcham_don|RT @TODAYshow: "D...|
|        EhHannah|@Keith1156 @canda...|
|JESUSFALFONSO1|RT @TrulyTrumpett...|
|         Mo_An2016|RT @DananaMama: @...|
|       Stumpcuttr|#CoonsAndFlake Ne...|
|       shoop_judy|RT @ChuckGrassley...|
|        SheilaUtz1|@RepAdamSchiff @H...|
|guernsey_robert|RT @RepAdamSchiff...|
|   LymeLadytrump|RT @LawrenceBuck1...|
|       trumpATeam|@realDonaldTrump ...|
|   pamelasengle1|RT @SenSchumer: A...|
+---------------+--------------------+
only showing top 20 rows
```

# Creating An S3 Bucket

To avoid losing data when destroying a cluster, you can store your results in S3 but to do so you need to first create a bucket. To do that you go back to the AWS Management Console and this time choose S3 service rather than EMR. This is to store output of your operations. The Jupyter file itself will be stored by default in S3 without creating a specific bucket for it.



Click "Create bucket" and follow the steps. You can use defaults for all steps. You should then have a bucket like the following:



You can now store your data produced in the Jupyter notebook to this bucket which will survive cluster terminations. For example, in PySpark you could extract tweet ids and store them in S3 as follows:

```
tweets.select("id").write.csv("s3://mybucket/mytweets")
```

You could also write the output to the cluster HDFS as follows but keep in mind these will be deleted with you terminate the cluster:

```
tweets.select("id").write.csv("hdfs://mybucket/mytweets")
```

# Warning! Be a Terminator!

You should create a cluster before starting your work and terminate it immediately after you are done to avoid wasting credits. Your notebook will be save in S3 and next time you open the notebook there is an option to choose a new cluster to associate to the notebook.

# Running Homework Jupyter Notebook

ssh to the master node of your cluster and run the following commands. If you haven't started working on the homework, then you can download the homework from S3 to the master node as follows:

```
$ aws s3 cp s3://caverlee/homework2.ipynb .
```

After that, you can start Jupyter in the terminal as follows:

```
$ ./start-jupyter.sh
```

You can then access it on the web your master node public domain and the port is 8880. The password is csce676.







# Warning: Save homework2.ipynb before termination!

Remember to download your edited homework2.ipynb to your machine or create a bucket on s3 and upload there because when you terminal the cluster it will be deleted.