# Scalable ML and AWS Practical 2024

**By Xianyuan Liu [Adapted from Michael Smith 2021-2023]**

You will have been sent an email with a login link to the Amazon console.
Please **skip** the steps marked with **(SKIP)**, as they do not apply to our lab.


## Logging in and the console

Once logged in, you will be placed in the AWS Console Home.


### 1. Locations

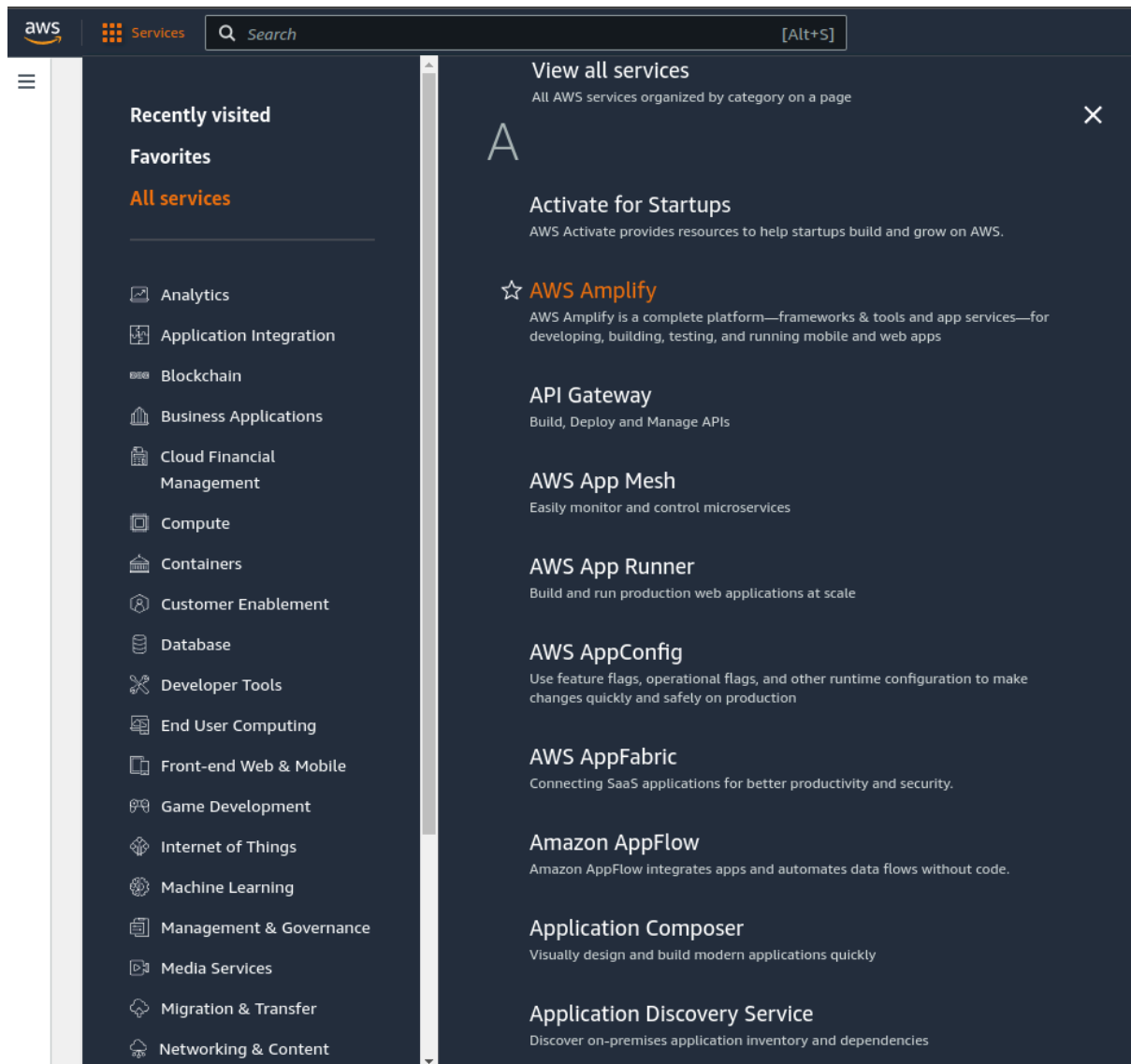Please note that AWS is hosted in <u>multiple locations</u> worldwide.



From Amazon:
> *"These locations are composed of Regions and Availability Zones. Each Region is a separate geographic area. Each Region has multiple, isolated locations known as Availability Zones. Amazon EC2 provides you with the ability to place resources, such as instances, and data in multiple locations. Although rare, failures can occur that affect the availability of instances that are in the same location. If you host all your instances in a single location that is affected by such a failure, none of your instances would be available."*

We will use one region today, **Europe (Ireland) eu-west-1**. You can see (and select) the region from the drop-down list in the top right of the console:

## 2. Services

You can click "Services" in the top left of the console to see the list of services available…
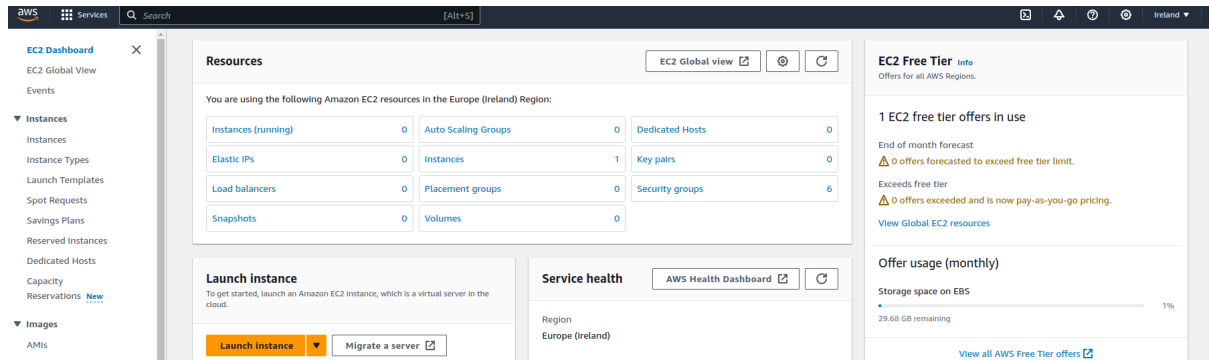


Under 'all services' one can select which tool of AWS one wishes to use. We will restrict ourselves for today to Elastic Compute Cloud (**EC2**), Simple Storage Service (**S3**), Elastic Map Reduce (**EMR**) and **Lambda**.

# Elastic Compute Cloud (**EC2**)

First, we'll explore EC2, set up an 'instance' (virtual machine) and connect to it.
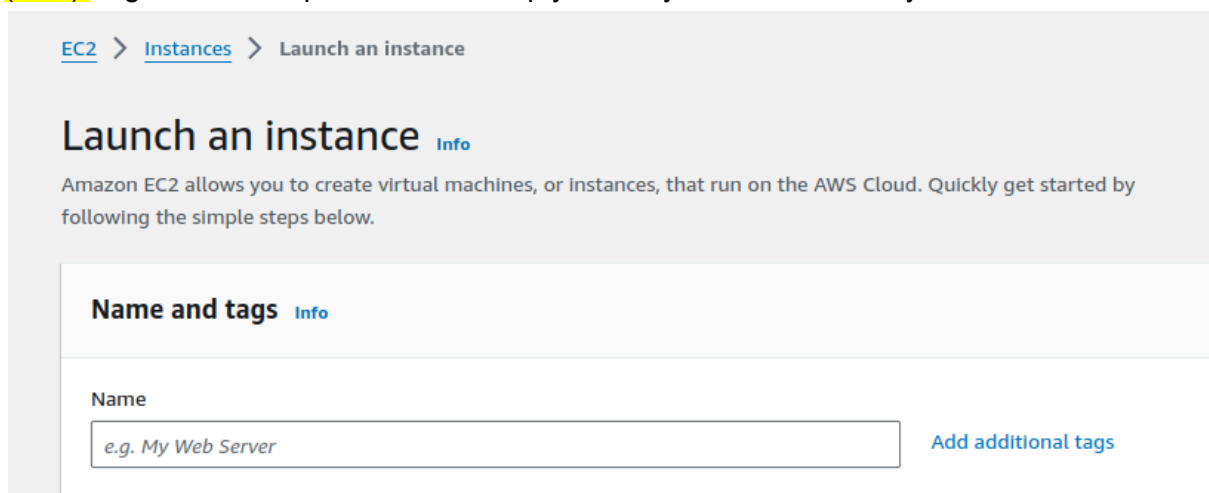- Click on **EC2** in services (or search for it in the top search bar), opening the EC2 console:



- Click on 'Instances' from the left pane. You will see an empty table and a big '**Launch Instance**' button - click it!



## 1. Launching instance

### Step 1: Name instance

As everyone is using the same account, it is necessary to **label** your **instance**. It might be useful to enter a **unique** name by **adding your ID**, e.g. the first part of your email address ('xliu'). A good and unique name can help you find your instance easily.

## Step 2: Application and OS images

The second step is deciding which machine image to use for your new instance. For an easy start, consider using one of the images created by Amazon. <mark>Select the Ubuntu Server 24.04 LTS (HVM) machine image.</mark>



## Step 3: Choose the Instance type

The third step is selecting the (virtual) hardware the instance will run on. For this practical, we ask that you <mark>use t2.micro (preferred type), t2.nano, or t2.small instance type</mark> (as we have increased the limit on this account to allow you all to start one). (The current limit is set to 512 vCPUs can be instantiated by the class).

Step 4: Key pair

A quick tour, to explain key pairs and ssh… (quoting from ssh.com):

> *"Each SSH key pair includes two keys:*
>
> *A **public key** that is copied to the SSH server(s). Anyone with a copy of the public key can encrypt data which can then only be read by the person who holds the corresponding private key. Once an SSH server receives a public key from a user and considers the key trustworthy, the server marks the key as authorized in its authorized_keys file. Such keys are called authorized keys.*
>
> *A **private key** that remains (only) with the user. The possession of this key is proof of the user's identity. Only a user in possession of a private key that corresponds to the public key on the server will be able to authenticate successfully. The private keys need to be stored and handled carefully, and no copies of the private key should be distributed. The private keys used for user authentication are called identity keys."*

AWS uses a key pair because it is typically more secure, allows for automation, and is generally the standard method for secure communication.

- Click '**create a new key pair**' and enter a **\*\*unique\*\* key pair name**, e.g. your instance name + key ('xliukey'). [Use the defaults: **RSA**]

  Please note: in **Private key file format**, **Linux** users should choose .**pem** file, while **Windows** users should choose a .**ppk** file for use with PuTTY in EMR (The following example uses .pem. Windows users should replace .pem with .ppk).

- Click '**Create key pair**'. You'll receive a file called 'xliukey.pem' (or 'xliukey.ppk' or whatever). You'll need this to SSH into this new instance. Keep it safe and secret.

*Step 5: Other options (**SKIP**)*

*There is no need to modify any of this. But a couple of items of interest:*

- *Allow SSH traffic from anywhere.*
- *Security groups are how EC2 organises access to the instances you create. I've already created one called '**justssh**' which gives access to the SSH port from anywhere. Typically one would restrict this to be from just your IP address, for example. Feel free to either use a security group that already exists or create a new one. You'll need to be able to SSH into the server later.*

***Spot pricing** Typically AWS will not be using all its computational resources. To make use of this 'spare' hardware, AWS offer a service called 'spot pricing' which is typically considerably cheaper than the on-demand price but comes at the cost of an instance that might be terminated with two minute warning.*

***Storage** Simply leave it an 8GB general-purpose SSD.*
*Please note:*
*There are many types of storage provided by AWS:*
- *Low cost, slow access: Amazon Glacier*
- *Elastic Block Store: This is the type of storage you need in your EC2 instance. ([More info](#)) This comes in four flavours,*
    - *slowest/cheapest: sc1 (cold HDD, solid state)*
    - *still cheap: st1 (throughput-optimised HDD)*
    - *solid-state: gp2,gp3 (general purpose SSD)*
    - *fast/expensive: io1,io2.*

Finally, click '**Launch instance**' on the right-hand side:



## 2. SSHing into your new instance

You'll be shown a summary stating your instances are launching. Either **click** the link to the instance or return to the list of instances and filter by the tag your username or email address you entered.
You might need to wait a few seconds while the instance starts... Also, you might need to click on the refresh button:



Click on the instance and then right-click (or use the button at the top) select 'Connect', then click 'SSH client'. This will give instructions on how to SSH in. The example is as below. Please locate your .pem/.ppk file.

## For Linux users...

In Linux and iOS, the key file's permissions should be set to read-only using the command line:

**chmod 400 filename.pem**

or the file browser (right-click the key file, and choose properties->permissions):



then ssh:

**ssh -i "filename.pem"
ubuntu@ec2-54-78-205-229.eu-west-1.compute.amazonaws.com**

It will ask if you're sure (as it can't guarantee that there's not a man-in-the-middle-attack):
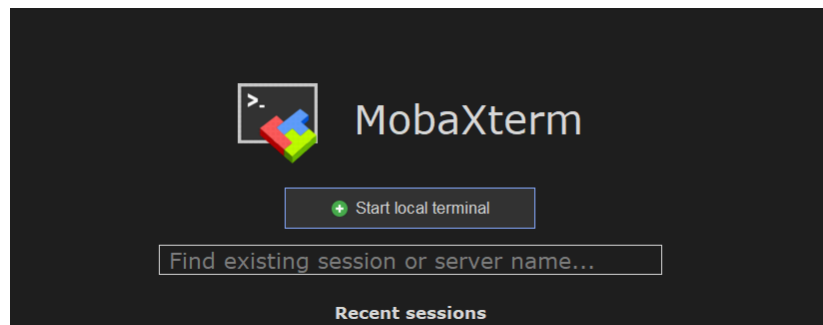
Tip: One usually keeps the .pem files in the .ssh folder in one's home directory.

## For Windows users...

In Windows, things are more complicated. We recommend using MobaXterm (AWS has instructions on how to get this working with putty).

## Option 1

- Click '**Start local terminal**' after you open MobaXterm.



- Follow the guidance of 'For Linux users' step by step. The commands are
    1) changing permissions by
  **chmod 400 filename.ppk**
    2) then sshing:
  **ssh -i "filename.ppk"**
  **ubuntu@ec2-54-78-205-229.eu-west-1.compute.amazonaws.com**



## Option 2

Follow these instructions to set up a new session for reuse.

If you can't get that working, please try alternative SSH clients. *(We have not enabled IAM access to "EC2 instance connect" otherwise anyone could access any other EC2 instance).*

*Also, this slightly awkward process is a step that hopefully shows how the key approach to security works etc.*

## 3. Disconnect and stop
- Enter '**exit**' in the terminal to log out of the instance
- Click '**Instance state**' and '**Stop instance**' to power off the instance <mark>(Please double-check the instance name to make sure you are stopping your instance)</mark>

# Elastic Map Reduce (EMR)

We'll set up a cluster to run a simple Recommender system (a little like the one you set up earlier in the course). Here we're not interested in the algorithm, more in some of the aspects of the cloud.

We're using Elastic Map Reduce.

> "The node types in Amazon EMR are as follows:
> - Primary node: A node that manages the cluster by running software components to coordinate the distribution of data and tasks among other nodes for processing. The primary node tracks the status of tasks and monitors the health of the cluster. Every cluster has a primary node, and it's possible to create a single-node cluster with only the primary node.
> - Core node: A node with software components that run tasks and store data in the Hadoop Distributed File System (HDFS) on your cluster. Multi-node clusters have at least one core node.
> - Task node: A node with software components that only runs tasks and does not store data in HDFS. **Task nodes are optional.**"
> – [AWS documentation](#).

## 1. Creating cluster

We will make a cluster with one primary node and about three core nodes.

- Click on **Services**, and find **[EMR](#)** (under analytics), then click **Create cluster**.

### Step 1: Name and applications

Enter a unique name (that won't be confused with someone else's cluster), e.g. the first part of your email address, e.g. 'xliu.cluster'. **Select the 'Spark' application bundle.**

## Step 2: Cluster configuration

- Edit **BOTH** the **Primary** and **Core** to use the "**r5a.xlarge**" instance type (we've increased the AWS limits on this type. You'll find that using other instance types will probably fail).
- **Remove** the "**Task 1 of 1**" instance group (we don't need Task nodes).

- Set the Core group cluster size to 2 (or 3).



You can save us money if you select 'Use spot purchasing option', but there is a really small risk your instances will get terminated by AWS during a spike in demand.

- Leave the rest of the page as its default: Don't enable cluster scaling. Auto-termination can be left on. EBS Root volume can be left at 15Gb.

- Tags: add the same name you gave before, i.e. xliu.cluster.

## Step 2: Bootstrap actions

Bootstrap actions can help us to install software or customise your instance configuration.
We have created and added the install-numpy.sh before the lab, so you can set the right
directory to load it:
- Choose **Bootstrap actions**
- Choose **Add**
- Enter
  **InstallNumpy** for Name and
  **s3://install-numpy/install-numpy.sh** for Script location
- Choose **Add bootstrap action**

## Step 4: Security configuration and EC2 key pair

**IMPORTANT: In this section you need to select the EC2 keypair you created**
**previously.** (If you don't do this you won't be able to connect to the cluster!). Please make
sure you choose the correct one!



## Step 5: Identity and Access Management (IAM) roles

The cluster and its constituent instances need to be able to do things in your AWS
environment. We use roles to control what they can or can't do. For our use case,
- select the **EMR-Sevice-Roles** for the Amazon EMR service role, and
- **EMR-EC2-InstanceProfile-Roles** for the EC2 instance profile for Amazon EMR:

## Step 6: Launch

- Finally, click "**Create cluster**" and wait.

# 2. Connecting

You will be able to select the cluster from the EMR Clusters page:



- **Click** your cluster's ID (or select it and click "view details") to check the details. Please be patient, it takes a while to get the cluster spun up. Initially, SSH won't be available, and even when it is PySpark etc might not be installed when you first arrive - wait a minute and click the refresh button again.

- Click '**Connect to the Primary node using SSH**'. You can find guidance on a connection using SSH on different systems, which you can follow accordingly.

# 3. Authorise inbound traffic

Authorising for primary node and core and task nodes. The guidance refers to
https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-connect-ssh-prereqs.html
- The **Properties** tab on this page will be pre-selected.
- Under '**Network and security**' in the Properties tab, select the **arrow next to EC2 security groups (firewall)** to expand this section



If you would like to update the **Primary** node, please select the **security group link** below primary node. You can update **core and task** nodes with the same way as follows. Please update both in this lab.
- Choose the **Inbound rules tab** and then choose **Edit inbound rules**.
- Check for an inbound rule that allows public access with the following settings. If it exists, choose **Delete** to remove it, which is (Type: SSH, Port: 22, Source: Custom 0.0.0.0/0)
- **Scroll** to the **bottom** of the list of rules and choose **Add Rule**.
- For Type, select **SSH**. This selection automatically enters TCP for Protocol and 22 for Port Range.
- For source, select **My IP** to automatically add your IP address as the source address.
- Choose **Save**.

# For Linux users…



## When you connect you will see…





# For Windows users…

Please follow the guidance provided by AWS EMR

**Connect to the primary node using SSH**                                    ✕

You can connect to the Amazon EMR primary node using SSH to perform actions like running interactive queries, examining log files, submit Linux commands and view web interfaces hosted on Amazon EMR clusters. Find out more 🔗

**Windows**    Mac/Linux

1. Download PuTTY.exe to your computer from:
   https://www.chiark.greenend.org.uk/~sgtatham/putty/latest.html 🔗
2. Start PuTTY.
3. In the Category list, choose Session.
4. In the Host Name field, enter `hadoop@ec2-18-200-244-5.eu-west-1.compute.amazonaws.com`
5. In the Category list, expand Connection > SSH and then choose Auth.
6. For Private key file for authentication, choose Browse and select the private key file ( `xliukey.ppk` ) that you used to launch the cluster.
7. Select Open.
8. Select Yes to dismiss the security alert.

View web interfaces hosted on Amazon EMR clusters 🔗

**Close**

- Complete the **Host Name** field, like



- Select the **private key file**, like



- Select **Accept** to dismiss the security alert

The result is



# The AWS command line interface

We now have a cluster running and ready for our commands.

Let's try out the **AWS command line interface**. It is often the case that one will want to perform a console operation repeatedly or describe a complex action programmatically to make it easy to see what has been done. To this end, Amazon provides both a command line tool and an API. All the actions you can do with the console can be done via these alternative interfaces. As we're now logged into a machine that's got the AWS CLI installed, we could quickly try it out.

## 1. s3 access

First, we're going to look at Amazon's simple storage service (s3). This has a [web interface](), but we'll be accessing it using the command line. On your new instance type:

**aws s3 ls**

This will list the bucket(s):

To list the contents of the bucket:

**aws s3 ls ml10m100k**

The result is



Then to copy the data to our instance,

**aws s3 cp s3://ml10m100k . --recursive**



# 2. ec2 via the command line

Also, try `aws ec2 describe-instances` to list all the instances currently in this account. Enter **q** to exit.

The command line interface allows you to create and destroy instances, configure users, launch clusters, and do any other operation that you can also do with the web interface. You may find you don't have access to some of these commands, however.

==Be careful, you can terminate your fellow students' instances as you are all in the same account and have sufficient permissions==. For more help on a particular aspect, you can type something like,

**aws ec2 help**

## 3. Control via the API

We can also access AWS using the API.

On the command line, we need to install boto3, a library for accessing the AWS API from Python, by entering

      **sudo pip3 install boto3**

While still logged into the cluster, start a Python terminal via

      **python**

The result will be

```
[hadoop@ip-172-31-20-251 ~]$ python
Python 3.9.16 (main, Sep  8 2023, 00:00:00)
[GCC 11.4.1 20230605 (Red Hat 11.4.1-2)] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>>
>>> ▯
```

We can then initiate a session to control EC2 by connecting to the AWS API. This is done by entering the following commands:

```
import boto3
sess = boto3.client('ec2')
reservations = sess.describe_instances()['Reservations']

for res in reservations:
   for j in res['Instances']:
      print(j['InstanceType'])
      if 'PublicIpAddress' in j: print(j['PublicIpAddress'])
      if 'Tags' in j: print(j['Tags'])
      print('')
```

```
                          hadoop@ip-172-31-20-251:~              Q  =  _  □  ✕

>>> import boto3
>>> sess = boto3.client('ec2')
>>> reservations = sess.describe_instances()['Reservations']
>>>
>>> for res in reservations:
...     for j in res['Instances']:
...         print(j['InstanceType'])
...         if 'PublicIpAddress' in j: print(j['PublicIpAddress'])
...         if 'Tags' in j: print(j['Tags'])
...         print('')
...
t2.micro
[{'Key': 'Name', 'Value': 'xliu'}]

r5a.xlarge
[{'Key': 'aws:elasticmapreduce:job-flow-id', 'Value': 'j-I6UAFNX5QYZM'}, {'Key':
 'aws:elasticmapreduce:instance-group-role', 'Value': 'MASTER'}, {'Key': 'for-us
e-with-amazon-emr-managed-policies', 'Value': 'true'}]

r5a.xlarge
[{'Key': 'for-use-with-amazon-emr-managed-policies', 'Value': 'true'}, {'Key': '
aws:elasticmapreduce:instance-group-role', 'Value': 'CORE'}, {'Key': 'aws:elasti
cmapreduce:job-flow-id', 'Value': 'j-I6UAFNX5QYZM'}]

r5a.xlarge
[{'Key': 'aws:elasticmapreduce:instance-group-role', 'Value': 'CORE'}, {'Key': '
for-use-with-amazon-emr-managed-policies', 'Value': 'true'}, {'Key': 'aws:elasti
cmapreduce:job-flow-id', 'Value': 'j-I6UAFNX5QYZM'}]

r5a.xlarge
[{'Key': 'aws:elasticmapreduce:instance-group-role', 'Value': 'CORE'}, {'Key': '
aws:elasticmapreduce:job-flow-id', 'Value': 'j-I6UAFNX5QYZM'}, {'Key': 'for-use-
with-amazon-emr-managed-policies', 'Value': 'true'}]

r5a.xlarge
3.250.110.188
[{'Key': 'aws:elasticmapreduce:job-flow-id', 'Value': 'j-2LHM74TBPTTUM'}, {'Key'
: 'aws:elasticmapreduce:instance-group-role', 'Value': 'MASTER'}]

r5a.xlarge
34.244.33.98
[{'Key': 'aws:elasticmapreduce:instance-group-role', 'Value': 'CORE'}, {'Key': '
aws:elasticmapreduce:job-flow-id', 'Value': 'j-2LHM74TBPTTUM'}]

r5a.xlarge
176.34.205.73
[{'Key': 'aws:elasticmapreduce:instance-group-role', 'Value': 'CORE'}, {'Key': '
aws:elasticmapreduce:job-flow-id', 'Value': 'j-2LHM74TBPTTUM'}]

r5a.xlarge
52.208.86.111
[{'Key': 'aws:elasticmapreduce:instance-group-role', 'Value': 'CORE'}, {'Key': '
aws:elasticmapreduce:job-flow-id', 'Value': 'j-2LHM74TBPTTUM'}]

>>> ▯
```

This will print out a list of all the instance types, and their tags.
Leave the Python interface [ **ctrl-D** or **exit()** ] so we can use PySpark instead...

## 4. Word counting using PySpark

We need to put the datafiles we're going to use within reach of Hadoop….We simply run (on the command line). Referring to https://spark.apache.org/examples.html.
First, copy the data to our instance

**aws s3 cp s3://some-texts . --recursive**

Then, put the data within reach of Hadoop (Please try again, if error '8020 failed on connection exception: java.net.ConnectException: Connection refused' occurs.)

**hadoop fs -put *.txt /**



This copies the data to the [Hadoop Distributed File System](#).

## Spark RDD Example

You would like to compute the count of each word in the text file. Spark allows for efficient execution of the query because it parallelizes this computation. Many other query engines aren't capable of parallelizing computations.

First, **Enter 'pyspark'** to open the interactive python/spark command line (not just python!).

Then, performing the computation with Spark RDDs:

```
text_file = spark.sparkContext.textFile("/texts.txt")

counts = (
    text_file.flatMap(lambda line: line.split(" "))
    .map(lambda word: (word, 1))
    .reduceByKey(lambda a, b: a + b)
)

counts.collect()
```

```
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 3.5.0-amzn-1
      /_/

Using Python version 3.9.16 (main, Sep  8 2023 00:00:00)
Spark context Web UI available at http://ip-172-31-17-106.eu-west-1.compute.inte
rnal:4040
Spark context available as 'sc' (master = yarn, app id = application_17145644309
42_0003).
SparkSession available as 'spark'.
>>> text_file = spark.sparkContext.textFile("/texts.txt")
>>>
>>> counts = (
...     text_file.flatMap(lambda line: line.split(" "))
...     .map(lambda word: (word, 1))
...     .reduceByKey(lambda a, b: a + b)
... )
>>>
>>> counts.collect()
[('The', 5), ('Wind', 6), ('', 5), ('One', 1), ('conversation', 1), ('of', 3), (
'them', 1), ('was', 2), ('settle', 1), ('argument,', 1), ('spotted', 1), ('trave
ler', 4), ('down', 1), ('wearing', 1), ('boasted,', 1), ('"I', 1), ('am', 1), ('
sure', 1), ('take', 2), ('his', 3), ('than', 2), ('in', 1), ('gentle', 1), ('str
ength,', 1), ('let', 1), ('go', 1), ('blew', 1), ('as', 2), ('but', 1), ('blew,'
, 1), ('tighter', 1), ('wrapped', 1), ('around', 1), ('several', 1), ('attempts,
', 1), ('gave', 1), ('turn.', 1), ('forceful,', 1), ('gently', 1), ('upon', 1),
('warmer,', 1), ('felt', 1), ('sun', 1), ('decided', 1), ('stronger', 1), ('blus
ter.', 1), ('Sun', 4), ('and', 6), ('the', 16), ('day,', 1), ('were', 1), ('havi
ng', 1), ('a', 3), ('about', 1), ('which', 1), ('one', 1), ('stronger.', 1), ('T
o', 1), ('they', 1), ('walking', 1), ('road,', 1), ('warm', 1), ('coat.', 2), ('
I', 1), ('can', 2), ('get', 1), ('that', 2), ('to', 3), ('off', 2), ('coat', 2),
 ('faster', 1), ('you', 1), ('can!"', 1), ('Sun,', 1), ('confident', 1), ('its',
 1), ('agreed', 1), ('first.', 1), ('hard', 1), ('it', 4), ('could,', 1), ('hard
er', 1), ('him.', 1), ('After', 1), ('grew', 2), ('frustrated', 1), ('up.', 1),
('Then,', 1), ("Sun's", 1), ('Instead', 1), ('being', 1), ('shone', 1), ('earth.
', 1), ('As', 1), ('pleasant', 1), ('warmth', 2), ('had', 1), ('proven', 1), ('g
entleness', 1), ('be', 1), ('force', 1)]
>>>
```

- Exit by using **Ctrl-D**

## 5. Make a recommendation using PySpark

We need to put the datafiles we're going to use within reach of Hadoop….We simply run (on the command line).

**hadoop fs -put *.dat /**

This copies the data to the [Hadoop Distributed File System](#).
You can check if the files there, via

**hadoop fs -ls /**

The result is

```
[hadoop@ip-172-31-16-149 ~]$ hadoop fs -put *.dat /
[hadoop@ip-172-31-16-149 ~]$ hadoop fs -ls /
Found 8 items
drwxr-xr-x   - hdfs    hdfsadmingroup          0 2024-05-02 19:54 /apps
-rw-r--r--   1 hadoop hdfsadmingroup     522197 2024-05-02 20:10 /movies.dat
-rw-r--r--   1 hadoop hdfsadmingroup  265105635 2024-05-02 20:10 /ratings.dat
-rw-r--r--   1 hadoop hdfsadmingroup    3584119 2024-05-02 20:10 /tags.dat
-rw-r--r--   1 hadoop hdfsadmingroup        841 2024-05-02 20:09 /texts.txt
drwxrwxrwt   - hdfs    hdfsadmingroup          0 2024-05-02 19:55 /tmp
drwxr-xr-x   - hdfs    hdfsadmingroup          0 2024-05-02 19:54 /user
drwxr-xr-x   - hdfs    hdfsadmingroup          0 2024-05-02 19:54 /var
[hadoop@ip-172-31-16-149 ~]$
```

## The PySpark command line

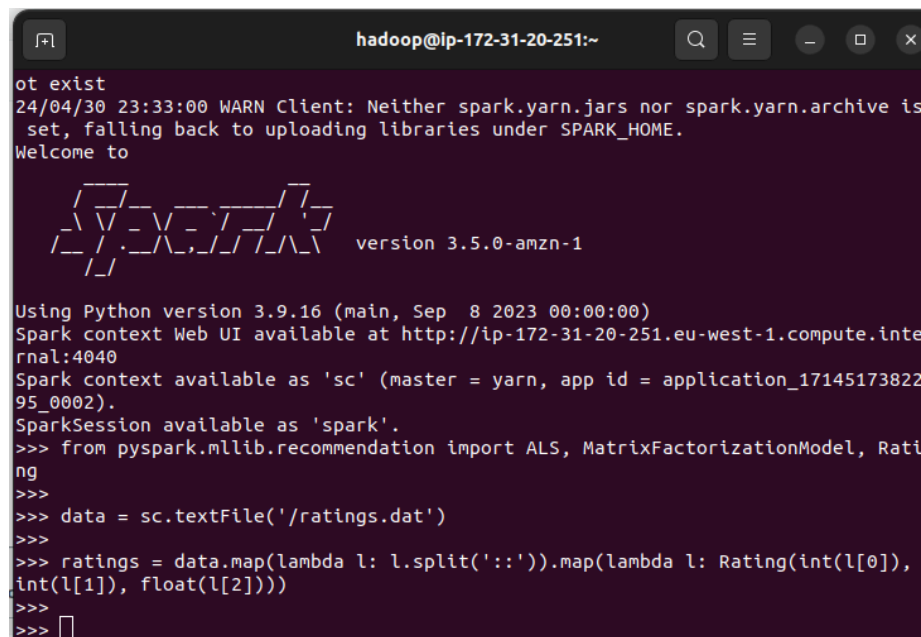**Enter '`pyspark`'** to open the interactive python/spark command line <mark>(not just python!).</mark>

We load the rating data from the Hadoop file system. This is a '**::**' separated table of ratings. The lambda/map functions are to split by these '**::**' and assign the types to the three columns.

> **from pyspark.mllib.recommendation import ALS, MatrixFactorizationModel, Rating**
>
> **data = sc.textFile('/ratings.dat')**
>
> **ratings = data.map(lambda l: l.split('::')).map(lambda l: Rating(int(l[0]), int(l[1]), float(l[2])))**

This will be very quick as <mark>Spark only runs things when it needs to</mark> (hence none of the above is run until the following lines are run...). The result is



Let's ask for the first item, to check it's working:

> **ratings.first()**

## A quick recap of recommender-systems/collaborative filtering

Let's set up the model...we're trying to guess the ratings people give to particular movies. We assume for this model that the full matrix of all reviews by everyone of every movie could be approximated by multiplying two low-rank matrices… for example, if there are N people and M movies, we believe there is an NxM matrix we want to approximate potentially. We could approximate this with the product of an Nx2 and a 2xM matrix (to help intuition imagine if the first column of the Nx2 is how much each person likes funny movies, and the second column is how much they like scary movies, then the first row of the 2xM matrix is how funny

each movie is and the second row is how scary, realistically we don't anticipate looking for such patterns here).

## Back to the application

We split our dataset randomly into a training and test set, then build a model using the Alternating Least Squares approach[1] to factorising a matrix...

```
test, train = ratings.randomSplit(weights=[0.3, 0.7], seed=1)
rank = 2
numIterations = 2
model = ALS.train(train, rank, numIterations)

testdata = test.map(lambda p: (p[0],p[1]))
predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
ratesAndPreds = test.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
print("Mean Squared Error = " + str(MSE))
```

The result shows the value of the Mean Squared Error.

---

[1] *Alternating Least Squares (ALS) matrix factorization.*

*ALS attempts to estimate the ratings matrix R as the product of two lower-rank matrices, X and Y, i.e. X \* Yt = R. Typically these approximations are called 'factor' matrices. The general approach is iterative. During each iteration, one of the factor matrices is held constant, while the other is solved for using least squares. The newly-solved factor matrix is then held constant while solving for the other factor matrix.*

*From the spark API reference.*

```
Welcome to


      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 3.5.0-amzn-1
      /_/

Using Python version 3.9.16 (main, Sep  8 2023 00:00:00)
Spark context Web UI available at http://ip-172-31-28-255.eu-west-1.compute.inte
rnal:4040
Spark context available as 'sc' (master = yarn, app id = application_17145920872
09_0002).
SparkSession available as 'spark'.
>>> from pyspark.mllib.recommendation import ALS, MatrixFactorizationModel, Rati
ng
>>>
>>> data = sc.textFile('/ratings.dat')
>>>
>>> ratings = data.map(lambda l: l.split('::')).map(lambda l: Rating(int(l[0]),
int(l[1]), float(l[2])))
>>>
>>> ratings.first()
Rating(user=1, product=122, rating=5.0)
>>> test, train = ratings.randomSplit(weights=[0.3, 0.7], seed=1)
>>> rank = 2
>>> numIterations = 3
>>> model = ALS.train(train, rank, numIterations)
>>> testdata = test.map(lambda p: (p[0],p[1]))
>>> predictions = model.predictAll(testdata).map(lambda r: ((r[0], r[1]), r[2]))
>>> ratesAndPreds = test.map(lambda r: ((r[0], r[1]), r[2])).join(predictions)
>>> MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
>>> print("Mean Squared Error = " + str(MSE))
Mean_Squared Error = 0.7625960347158568
>>>
```

Finally, a quick look at the MSE on the training data;

**traindata = train.map(lambda p: (p[0],p[1]))**
**predictionstrain = model.predictAll(traindata).map(lambda r: ((r[0], r[1]), r[2]))**
**ratesAndPreds = train.map(lambda r: ((r[0], r[1]), r[2])).join(predictionstrain)**
**MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()**
**print("Mean Squared Error = " + str(MSE))**

The result is

```
>>>
>>> traindata = train.map(lambda p: (p[0],p[1]))
>>> predictionstrain = model.predictAll(traindata).map(lambda r: ((r[0], r[1]),
r[2]))
>>> ratesAndPreds = train.map(lambda r: ((r[0], r[1]), r[2])).join(predictionstr
ain)
>>> MSE = ratesAndPreds.map(lambda r: (r[1][0] - r[1][1])**2).mean()
>>> print("Mean Squared Error = " + str(MSE))
Mean_Squared Error = 0.7388751588161084
>>>
```

## 6. Terminate

- Enter '**exit**' in the terminal or [Ctrl + d] to exit pyspark
- Enter **'exit'** to log out of the instance
- Click **'Terminate'** to power off the cluster

Updated 7 minutes ago     ↻     **Terminate**     **Clone in AWS CLI**     **Clone**     ✕

est-1/elasticmapreduce

## Status and time

Status
⊘ **Waiting**

Creation time
2 May 2024 20:50 (UTC+01:00)

Elapsed time
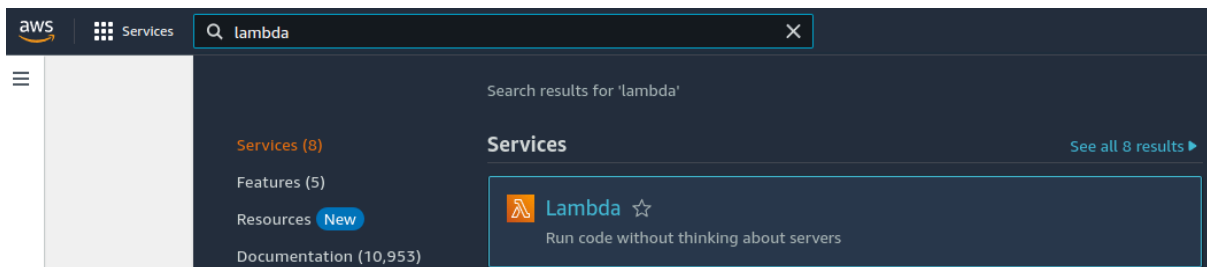28 minutes

-1.compute.amazonaws.com

ing SSH

ing SSM ↗

# Using Lambda

Another exercise: You can create a lambda function!

## 1. Create function

Find the lambda service…



- Click '**Create function**'
- Author from scratch
- Function name: pick something **unique**, such as adding your ID
- Runtime: Python 3.12



- Click '**Change default execution role**', choose '**use an existing role**' and then choose 'service-role/testfunction-role-X'.



Note: you need to change the execution role, as it, by default, tries to create a new role with lambda permissions (**the IAM profile you're using doesn't have permissions for that**). Example error message:

So please use the 'testfunction' role.
- Click



## 2. Add Trigger
- Click **Add Trigger** -> select **'API Gateway'** from the drop-down list
- Select **'Create a new API'**
- For Security select **'open'**
- Click **'Add'**. Here's the API Gateway configuration



- Click 'Add'

## 3. Update code
- Click **'code'**

Here's some example code for the lambda that adds together two numbers.

```
import json

def lambda_handler(event, context):
    if ('inputA' in event['queryStringParameters']) and ('inputB' in
event['queryStringParameters']):
```
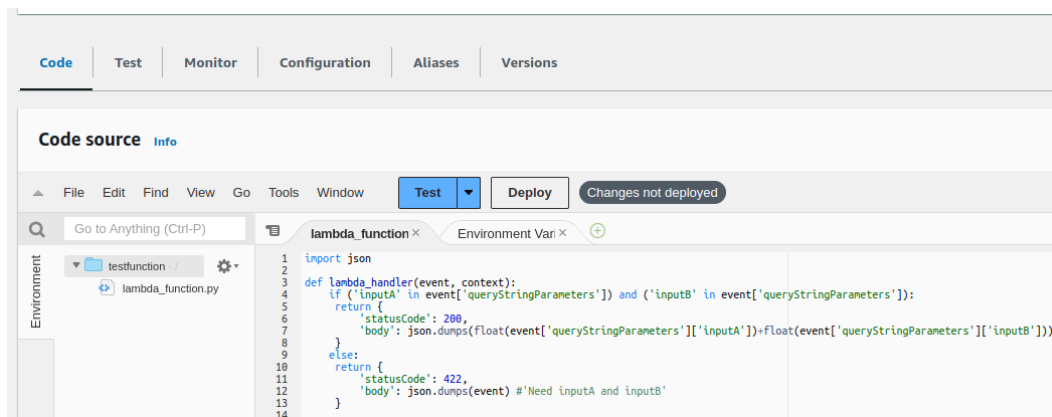
```
        return {
                'statusCode': 200,
                'body':
json.dumps(float(event['queryStringParameters']['inputA'])+float(event['queryStringPa
rameters']['inputB']))
        }
    else:
        return {
                'statusCode': 422,
                'body': json.dumps(event) #'Need inputA and inputB'
        }
```
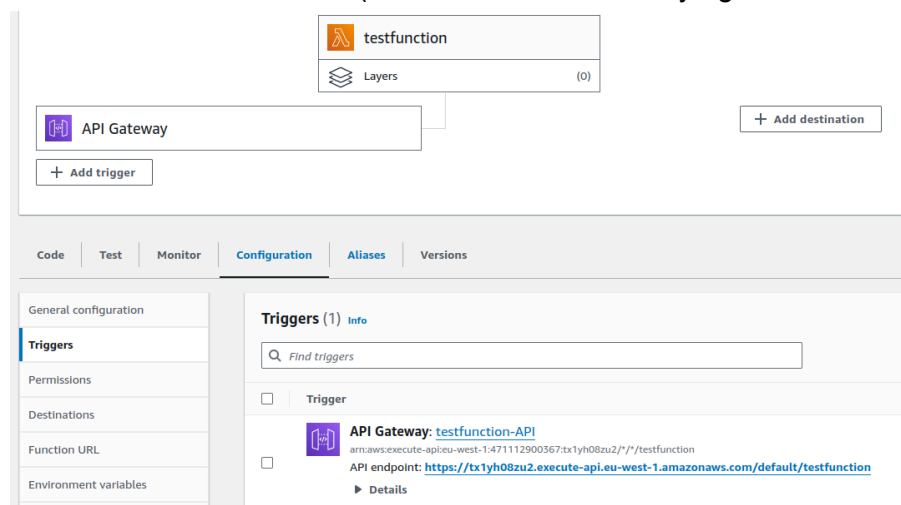
- Click **'Deploy'**



# 4. Test

## Option 1

- Visit the API interface address (click on the API Gateway again to find the URL):



- Click the link to the '**API endpoint**'.

- You'll probably get an "{"message":"Internal Server Error"}", as the code we wrote expects two arguments.
- Edit the URL by adding arguments after the URL, e.g.
  If the endpoint was at
  https://tx1yh08zu2.execute-api.eu-west-1.amazonaws.com/default/xliu,
  I will add
  **?inputA=42&inputB=58**
  to the end of the URL:
  https://tx1yh08zu2.execute-api.eu-west-1.amazonaws.com/default/xliu?inputA=3&inputB=4
  That seems to work! The result is 100.

Tip: Copying and pasting from the Google Doc sometimes jumbles the indentation in the code - so that might need sorting.

## Option 2

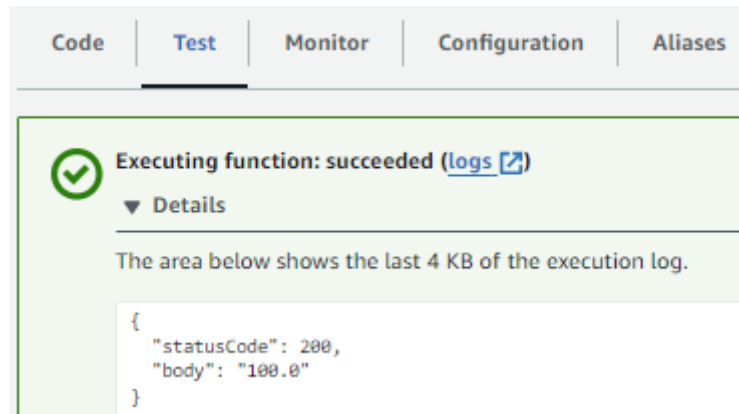You can also configure a test event (see Test tab).



Test event:
**{**
**"queryStringParameters": {**
  **"inputA": "42",**
  **"inputB": "58"}**
**}**

- Scroll down Details, and check the output: 100

```
{
    "statusCode": 200,
    "body": "100.0"
}
```

# Additional exercises (SKIP)

Trigger on a file being added to an s3 bucket…
https://docs.aws.amazon.com/lambda/latest/dg/with-s3-example.html

Back to the command line…

A global network of air pollution sensors is aggregated by openAQ. They archive this data on s3 (see here: https://openaq-data.s3.amazonaws.com/index.html)

(tip: you can leave PySpark using Ctrl-D).

We can copy one of these csv files to our instance from the command line using:

```
aws s3 cp s3://openaq-data/2019-01-01.csv .
```

Tip: This might be faster if the instance were in the same region as the data!

Or, from inside PySpark, can read the data directly:

```
df = spark.read.csv("s3a://openaq-data/2019-01-01.csv")
```

This could be improved by treating the headers properly for a start. Add options to do this (hint set `header=True`)

Feel free to try analysing this data. For example: count the number of measurements in Uganda...

Hint the 'country' column would equal 'UG'...

```
df[df['country']=='UG'].count()
```

Note that without the `.count()` the result is immediate - spark doesn't do any computation until it has to.

You'll find the operation is far quicker the second time thanks to caching. Also look at `.persist()` [here](#).


## Starting and stopping EC2 instances from the command line

Example code, modify to start an instance yourself.

```
aws ec2 run-instances --image-id ami-xxxxxxxx --count 1
--instance-type t2.micro --key-name MyKeyPair --security-group-ids
sg-903004f8 --subnet-id subnet-6e7f829e
```

Think about where you can get the **security-group-id**. (hint: either from the web console or for bonus credit from the command line interface!
hint: `aws ec2 describe-security-groups`).

# Preparation (SKIP, set by Xianyuan before the lab)

## Creating policies

**(Please Skip this process, as this is for the lecturer)**
*First, we will create suitable IAM policies and roles for our lab.*
  - *Click on **Services**, and find **IAM***
  - *Choose **Policies***

### List Buckets

  - *Choose **Create policy***
  - *Choose **s3** in the service*
  - *In Filter Actions, choose **List***
  - *Tick **ListBucket**, **ListAllMyBuckets***
  - *Tick **any** in **bucket***
  - *Choose **Next***
  - *Enter Policy name: **s3ListBuckets***
  - *Choose **Create policy***

### Describe Instances

  - *Choose **Create policy***
  - *Choose **ec2** in the service*
  - *In Filter Actions, enter **DescribeInstances***
  - *Tick **DescribeInstances***
  - *Choose **Next***
  - *Enter Policy name: **ec2DescribeInstances***
  - *Choose **Create policy***

### Create Role for EMR Sevices

  - *Choose **Create role***
  - *Choose **AWS service**, and in the Use case, choose **EMR**, and click **Next***
  - *Keep the default settings, and click **Next***
  - *In the Role name, enter **EMR-Sevice-Roles**, and click **Create role***

### Create Role for EMR EC2 Instance Profile

  - *Choose **Create role***
  - *Choose **AWS service**, and in the Use case, choose **EC2**, and click **Next***
  - *Add **ec2DescribeInstances** and **s3ListBuckets**, and click **Next***
  - *In the Role name, enter **EMR-EC2-InstanceProfile-Roles**, and click **Create role***

## Install_Numpy bootstrap actions

  - *Go to s3, Create bucket, Name: install-numpy, Create bucket*
  - *Click install-numpy, upload install-numpy.sh*

- *In install-numpy.sh, the code is*
  **#!/bin/bash**
  **sudo pip3 install numpy**

## *Upload text data for word counting*

- *Go to s3, Create bucket, Name: some-texts, Create bucket*
- *Click some-texts, upload texts.txt*
- *The texts.txt includes some paragraph found or generated by ChatGPT.*

## *Data downloading*

1. *Download MovieLens 10M Dataset (ml10m100k) from*
   [https://grouplens.org/datasets/movielens/10m/](https://grouplens.org/datasets/movielens/10m/).
2. *Upload to s3*