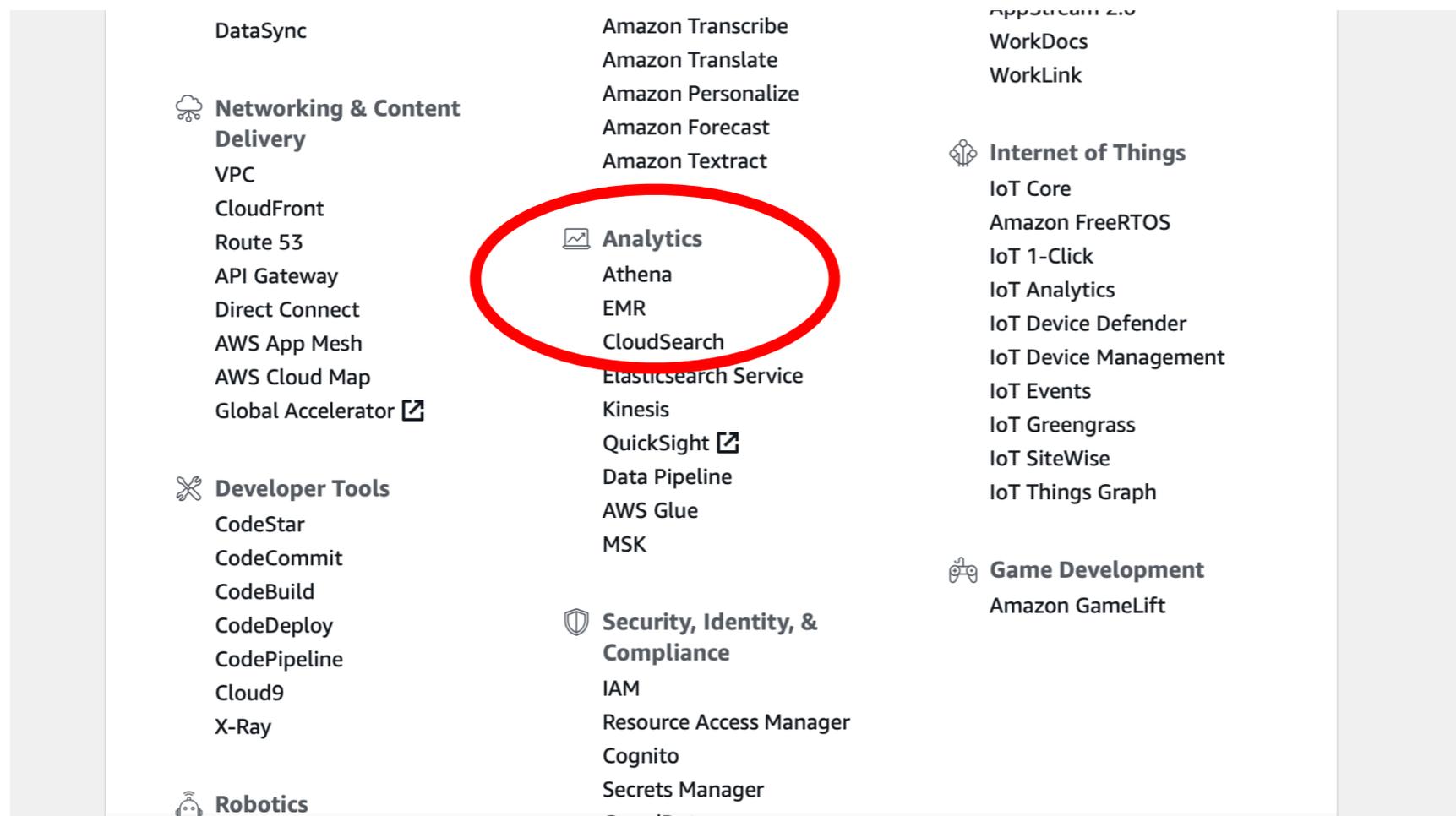# Project

Data at Scale

# Create a Cluster

- In AWS Management Console:

  - Select EMR under Analytics

# Creating a Cluster

- To avoid costs for moving data, use the Northern Virginia data center

# Creating a Cluster

- Go to AWS

# Creating a Cluster

- After five (!) minutes

# Creating a Cluster

- Click on the SSH link

**SSH**       ✕

## 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 ~/nvirginia.pem with the location and filename of the private key file (.pem) used to launch the cluster.

   ```
   ssh -i ~/nvirginia.pem hadoop@ec2-54-174-96-130.compute-1.amazonaws.com
   ```

3. Type yes to dismiss the security warning.

**Close**

# Connecting to the Cluster

- Enter the credential in your Command Window / Terminal

  - Be sure to specify the path to the private key that you used to obtain the cluster



```
● ● ●  🏠 thomasschwarz — hadoop@ip-172-31-38-58:~ — ssh -i ~/ssh_open/nvirginia....

Last login: Mon Apr 22 15:44:33 on ttys000
[MSCSs-MacBook-Pro-2:~ thomasschwarz$ ssh -i ~/ssh_open/nvirginia.pem hadoop@ec2-]
54-174-96-130.compute-1.amazonaws.com
The authenticity of host 'ec2-54-174-96-130.compute-1.amazonaws.com (54.174.96.1
30)' can't be established.
ECDSA key fingerprint is SHA256:Q5qDiJjnLHGdMfYE+t1jeLpxCJRWoOPRSpmuFQUebqE.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'ec2-54-174-96-130.compute-1.amazonaws.com,54.174.96.
130' (ECDSA) to the list of known hosts.
Last login: Mon Apr 22 21:49:47 2019
```

# Getting the data

- Go to Kaggle

  - Get the data sets Black Friday and Gun-violence

  - download them to a known location

  - sftp into your cluster

  - upload the files to a directory of your chosing

```
MSCSs-MacBook-Pro-2:BlackFriday thomasschwarz$ sftp -i ../../ssh_open/nvirginia.pem hadoop@ec2-54-162-49-218.compute-1.amazonaws.com
Connected to hadoop@ec2-54-162-49-218.compute-1.amazonaws.com.
sftp> ls
data
sftp> lls
BlackFriday.csv
sftp> put BlackFriday.csv
Uploading BlackFriday.csv to /home/hadoop/BlackFriday.csv
BlackFriday.csv                                                          100%  24MB  23.6MB/s   00:01
sftp> lls
BlackFriday.csv              gun-violence-data_01-2013_03-2018.csv
sftp> ls
BlackFriday.csv    data
```

# SFTP and FTP commands

# Getting the data

- In the master, organize your data

    - Manually delete the first row with the description of the fields

# Starting Pig

- Start pig

```
pig -x mapreduce
```

- Register piggybank

```
register file:/usr/lib/pig/lib/piggybank.jar
```

- Move data into Hadoop

```
grunt> sh ls data
BlackFriday.csv
gun-violence-data_01-2013_03-2018.csv
grunt> copyFromLocal ./data/BlackFriday.csv BlackFriday.csv
grunt> ls
hdfs://ip-172-31-40-153.ec2.internal:8020/user/hadoop/
BlackFriday.csv<r 1>  24956107
```

# Load data into Pig

```
blackfriday = LOAD 'blackfriday' USING PigStorage(',')
AS (UserID: int, ProductID: chararray, Gender:
chararray, Age: chararray, Occupation: int,
CityCategory: chararray, Stay_In_Current_City: int,
Marital_Status: int, Product_Category_1: chararray,
Product_Category_2: chararray, Product_Category_3:
chararray, Purchase: int);

illustrate blackfriday;

simplified = FOREACH blackfriday GENERATE UserID,
Gender, Age, Purchase;
```

# Getting the data

- Group by age and gender

```
gender_age = GROUP simplified by (Gender, Age);
```

- Calculate maximum and average purchase amounts per age-gender group

```
my_sum = FOREACH gender_age GENERATE group as
age_count, COUNT(simplified),
MAX(simplified.Purchase),
SUM(simplified.Purchase)/COUNT(simplified);
```

# Getting the data

- Display the results:

```
dump my_sum;
```

# Tasks

- Get the statistics from Black Friday for age, gender, age-gender, age-gender-CityCategory

- Use gun_violence_data to find

  - the city or county in Wisconsin with the highest number of gun incidents

  - the city in the country with the highest number of gun incidents