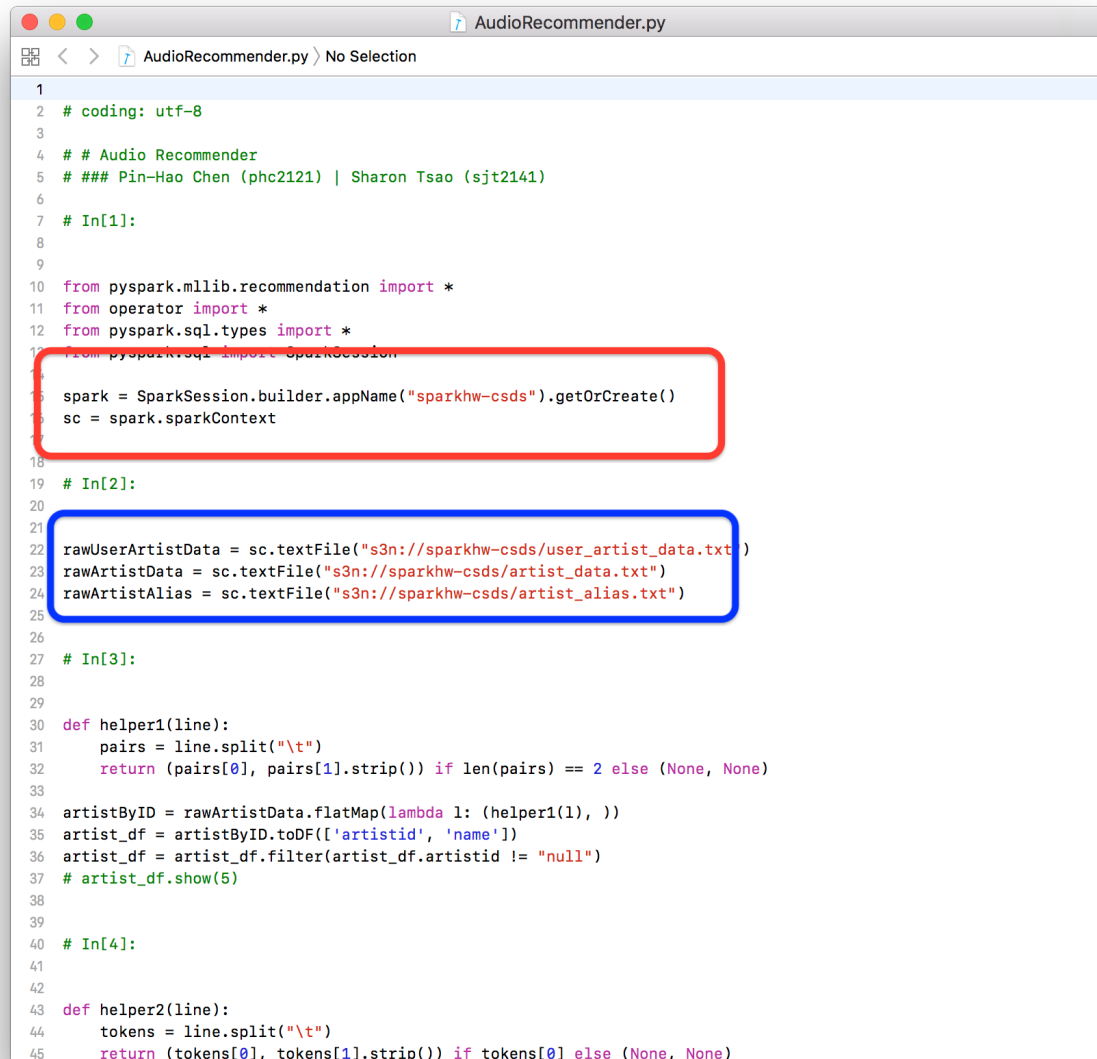Pin-Hao Chen (phc2121)
Sharon Tsao (sjt2141)

# Audio Recommender on AWS

1) We're uploading the files to our bucket, sparkhw-csds

2) In order to run our script on AWS, we've added a few lines to the code to initialize SparkContext (sc). This includes importing and creating Sparksession, as shown in the image below in red. Additionally, we changed the path to the input file location of the bucket on AWS, as shown below in blue. The file names are in the format of "s3n://bucketname/filename".

```python
# coding: utf-8

# # Audio Recommender
# ### Pin-Hao Chen (phc2121) | Sharon Tsao (sjt2141)

# In[1]:


from pyspark.mllib.recommendation import *
from operator import *
from pyspark.sql.types import *
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("sparkhw-csds").getOrCreate()
sc = spark.sparkContext


# In[2]:


rawUserArtistData = sc.textFile("s3n://sparkhw-csds/user_artist_data.txt")
rawArtistData = sc.textFile("s3n://sparkhw-csds/artist_data.txt")
rawArtistAlias = sc.textFile("s3n://sparkhw-csds/artist_alias.txt")


# In[3]:


def helper1(line):
    pairs = line.split("\t")
    return (pairs[0], pairs[1].strip()) if len(pairs) == 2 else (None, None)

artistByID = rawArtistData.flatMap(lambda l: (helper1(l), ))
artist_df = artistByID.toDF(['artistid', 'name'])
artist_df = artist_df.filter(artist_df.artistid != "null")
# artist_df.show(5)


# In[4]:


def helper2(line):
    tokens = line.split("\t")
    return (tokens[0], tokens[1].strip()) if tokens[0] else (None, None)
```
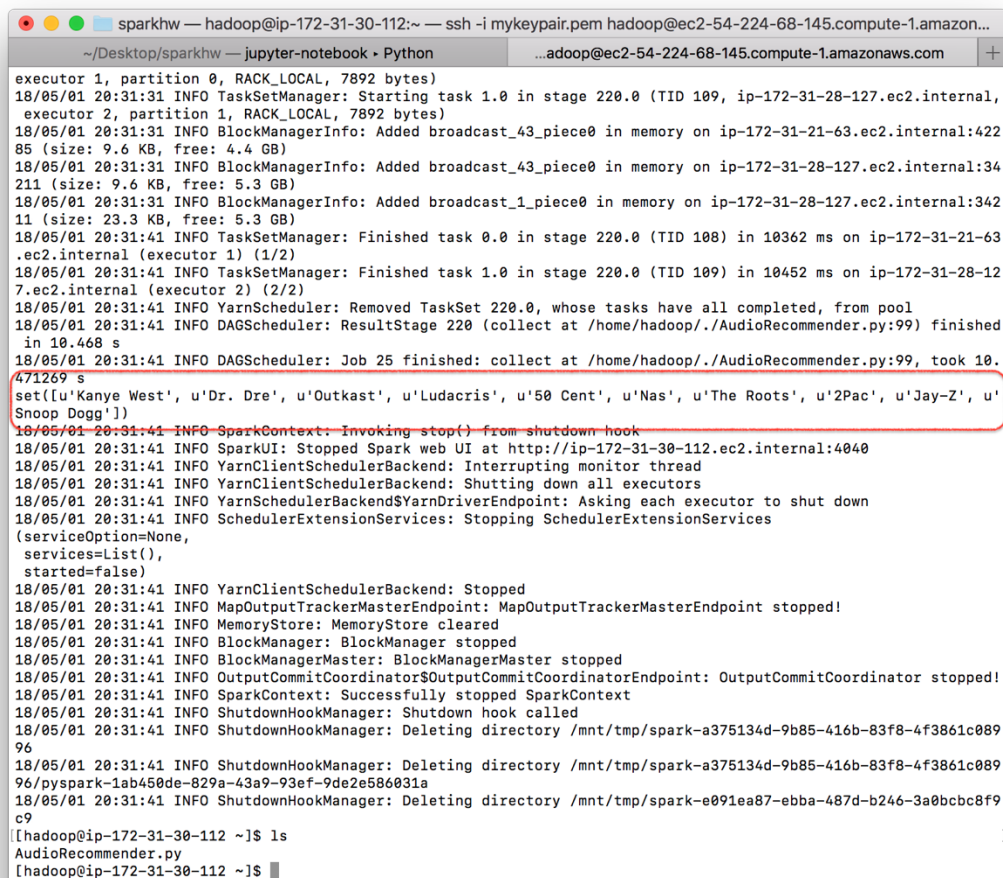
Pin-Hao Chen (phc2121)
Sharon Tsao (sjt2141)

# Audio Recommender on AWS

3) These are the files we've used in our bucket.



4) Here we are connecting to the instance through SSH by the permission key. We also changed the "mykeypair.pem" file to the 400 mode. Next we ran the following code in terminal.



```
1  ssh -i "mykeypair.pem" hadoop@ec2-54-224-68-145.compute-1.amazonaws.com
2
```

5) Screenshot of connecting to AWS EMR.

Pin-Hao Chen (phc2121)
Sharon Tsao (sjt2141)

# Audio Recommender on AWS

6) We copied over our python script and executed spark-submit.

```
[hadoop@ip-172-31-30-112 ~]$ aws s3 cp s3://sparkhw-csds/AudioRecommender.py .
[download: s3://sparkhw-csds/AudioRecommender.py to ./AudioRecommender.py
[[hadoop@ip-172-31-30-112 ~]$ ls
[AudioRecommender.py
[hadoop@ip-172-31-30-112 ~]$ spark-submit ./AudioRecommender.py
[18/05/01 20:26:33 INFO SparkContext: Running Spark version 2.3.0
18/05/01 20:26:33 INFO SparkContext: Submitted application: sparkhw-csds
18/05/01 20:26:33 INFO SecurityManager: Changing view acls to: hadoop
18/05/01 20:26:33 INFO SecurityManager: Changing modify acls to: hadoop
```

7) The top 10 Audio Recommender artists are printed, as circled in red in the image below.