# CIS 4930 Introduction to Hadoop and Big Data

# Write and Run a Spark Application and Configure a Spark

# Application

**Julie Reyes | u76631122**

**Introduction:**

The purpose of this lab is to write a self-contained Python Spark application to submit to the cluster. This lab differs from previous labs, in that an RDD will be created via a python application, rather than utilizing the interactive spark application. The program for this lab is a simple python program that counts the number of JPG requests in a web log file on the Hadoop cluster.

**Lab:**

A python application was created with the text editor gedit**.** The following is the code used to create the application named CountJPGs.py:

```
import sys
from operator import add

from pyspark import SparkContext

if __name__ == "__main__":
    if len(sys.argv) != 2:
        print >> sys.stderr, "Usage: countJPGs <file>"
        exit(-1)
    sc = SparkContext(appName="CountJPGs")
    logRDD = sc.textFile(sys.argv[1], 1)
    output = logRDD.filter(lambda line: ".jpg" in line).count()
    print output
sc.stop()
```

The CountJPGs.py program performs in a similar manner to the previous assignment as it counts the number of JPGs in the weblog file. In the CountJPGs.py program, a logRDD is created from the file at the path entered in the command line. The output is calculated by filtering the logRDD for lines that contain jpg and counting the total amount of JPG lines, which was completed in the line that states:

**output = logRDD.filter(lambda line: ".jpg" in line).count()**

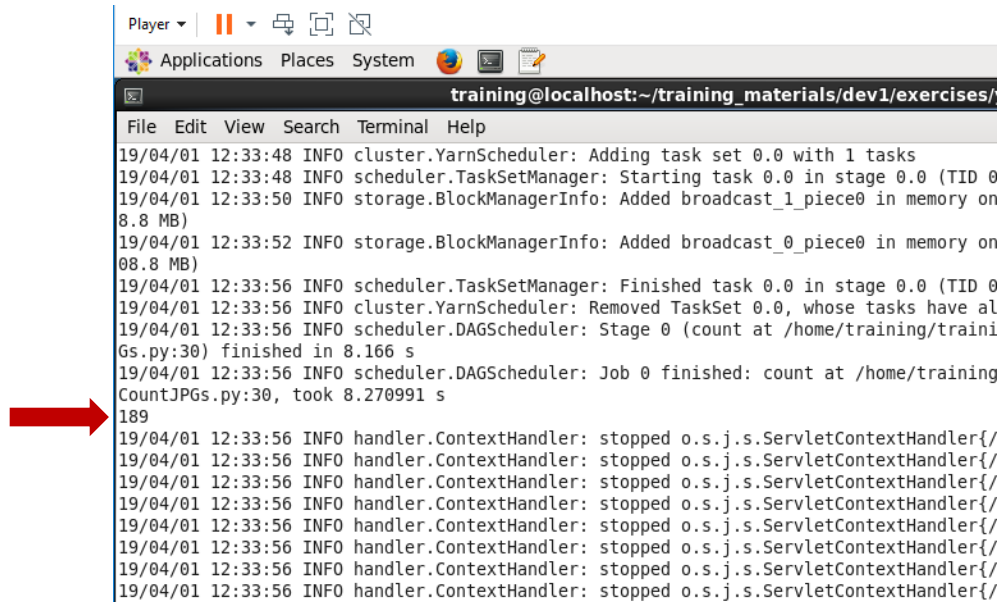Then, the final total count is printed to the terminal by the line:

**print output**

Lastly, the program stops the spark context to terminate the application with the sc.stop().To submit the application the following command line was entered into the terminal:
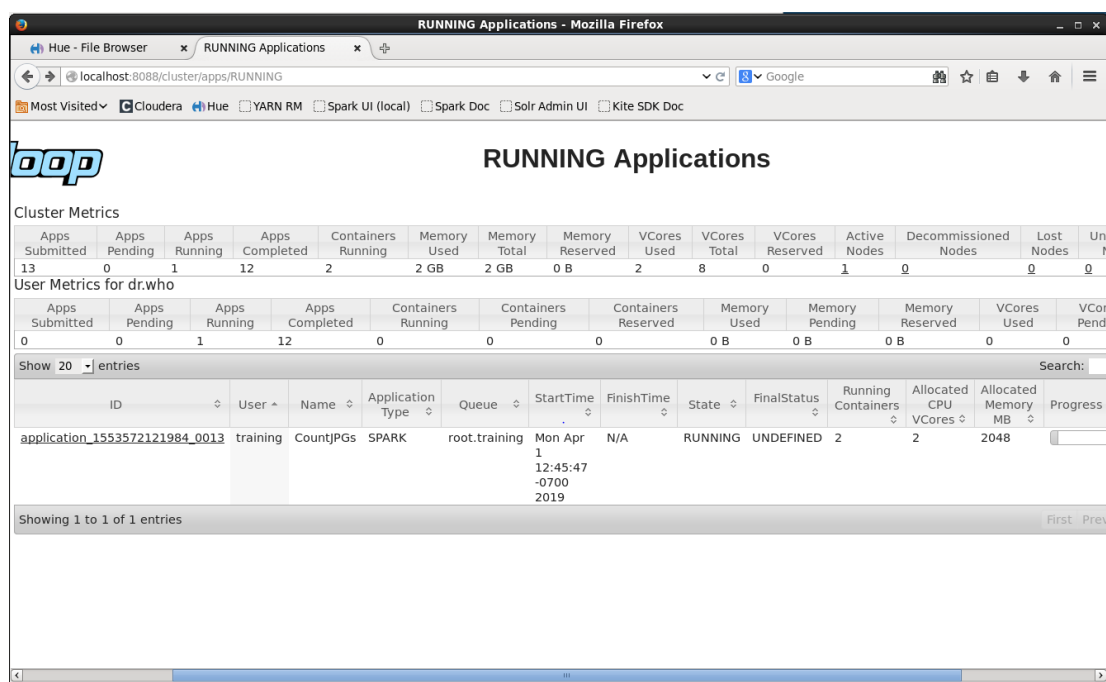
**spark-submit CountJPGs.py /loudacre/weblogs/***

This job ran locally produced the following output:



Therefore, there were a total of 189 webpages that were of type JPG in the weblogs directory.

Next, the job was submitted to the cluster and tracked via the yarn resource manager UI. The following is a snapshot of the job being ran on the cluster and tracked through the Yarn resource manager:

You can see that the application 1553572121984_0013 is running when the job is submitted to the cluster. To submit the application to the cluster the following command was used:

**spark-submit   --master yarn-client   CountJPGs.py /loudacre/weblogs/***