# CSP554—Big Data Technologies

## Assignment #10 (5 points extra credit)

If you don't do the extra credit exercise, you don't need to submit anything, as the assignment itself is just here to let you know what to read and is worth zero points.

If you want to try for up to 5 points of extra credit, then do exercise 1 and submit your results. This exercise is a good thing to try if you want to get more hands on about streaming technology, especially if you are considering using spark streaming for your project.

## Readings

Read Chapters 1-3 from: Pramod J. Sadalage and Martin Fowler. 2012. *NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot Persistence*. Addison-Wesley.(PS)

Note, it is time to purchase/obtain this book. The first three chapters are available in the "Free Books and Chapters" section of our blackboard site to get you started.

## 5 points extra credit

### Due by the start of the next class period

Assignments should be uploaded via the Blackboard portal

Exercise 1) 5 points extra credit

Follow the document "Spark Streaming Demo Instructions" included with this assignment and execute the demo code. Provide enough screen shots to indicate you have completed the document through section 4. Then remember to terminate your VM.

# Spark Streaming Demo

## Overview

This demo illustrates how to execute a pyspark (Python) spark streaming job. The job accepts a sequence of lines that the user types in onto one terminal window over a 10 second interval and then counts the number of distinct words in those lines and outputs the word count results to a second terminal window. This continues every 10 seconds. To do this we will set up a Spark EMR cluster and connect two terminal windows to it. In the first we will run the Linux 'nc' (Netcat) command. It will open a TCP socket on port 3333. After it does so, any line you then type will be sent out on that port. In another terminal window we will execute a pyspark word count program that will set up the spark streaming pipeline using DStreams. Our initial DStream will be connected to and read the lines from port 3333 and then go on to perform the word count process.

So on one terminal (connected to the EMR master node) you might see:

> [hadoop@ip-172-31-19-223 ~]$ nc -lk ec2-3-91-10-18.compute-1.amazonaws.com 3333
>
> this is a test of the the system         <- note

And output from the word count program running in the other terminal should look something like:

> -------------------------------------------
>
> Time: 2019-11-05 21:25:00
>
> -------------------------------------------
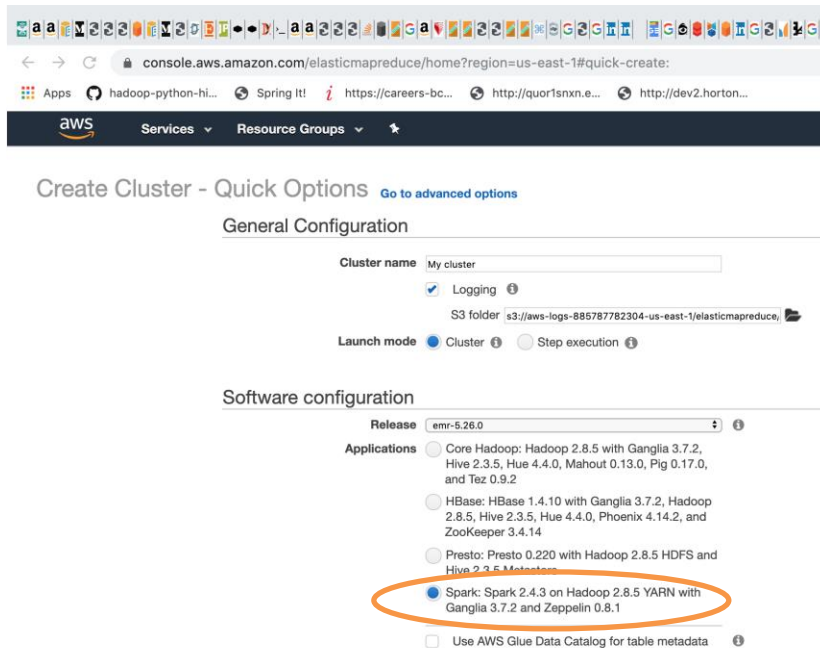>
> (u'a', 1)
>
> (u'this', 1)
>
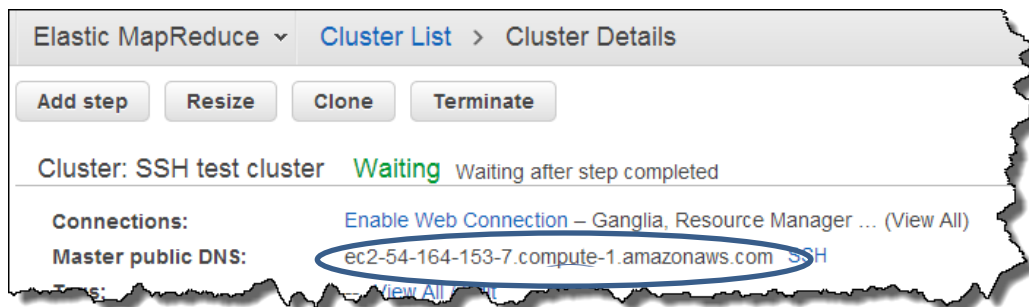> (u'is', 1)
>
> (u'test', 1)
>
> (u'the', 2)
>
> (u'of', 1)
>
> (u'system', 1)

## Running the Demo

1) Start up a Hadoop cluster as previously, but instead of choosing the "Core Hadoop" configuration chose the "Spark" configuration (see below), otherwise proceed as before.



2) At a later point in these instructions you will need to use the public DNS name of the master node of your EMR cluster. To retrieve it using the Amazon EMR console

   a) Find the EMR service page.

   b) On the **Cluster List** page, select the link for your cluster.

   c) Note the **Master public DNS** value that appears at the top of the **Cluster Details** page.



3) Download consume.py and log4j.properties files from the assignment to your local PC or MAC

4) There is one item you must change in consume.py. In the following line you must replace <Master public DNS> with your own public DNS name (found as described above)

   lines = ssc.socketTextStream("<Master public DNS>", 3333)

   For example:

   lines = ssc.socketTextStream("ec2-54-164-153-7.compute-1.amazonaws.com", 3333)

```python
1    from pyspark import SparkContext
2    from pyspark.streaming import StreamingContext
3
4    # Create a local StreamingContext with a batch interval of 10 seconds
5    sc = SparkContext("yarn", "NetworkWordCount")
6    ssc = StreamingContext(sc, 10)
7
8    # Create a DStream
9    lines = ssc.socketTextStream("ec2-3-15-178-106.us-east-2.compute.amazonaws.com", 3333)
10
11   # Split each line into words
12   words = lines.flatMap(lambda line: line.split(" "))
13
14   # Count each word in each batch
15   pairs = words.map(lambda word: (word, 1))
16   wordCounts = pairs.reduceByKey(lambda x, y: x + y)
17
18   # Print each batch
19   wordCounts.pprint()
20
21   ssc.start()                  # Start the computation
22   ssc.awaitTermination()       # Wait for the computation to terminate
```

5) scp this modified consume.py file to your EMR cluster master node. You may need to answer a security question with "Y/y" or "Yes".

```
yashmp.21197@Yashu MINGW64 /d/Academic/Sem-4/CSP-554-BDT/Assignments/Assignment-10
$ scp -i emr-key-pair.txt consume.py hadoop@ec2-3-15-178-106.us-east-2.compute.amazonaws.com:/home/hadoop
The authenticity of host 'ec2-3-15-178-106.us-east-2.compute.amazonaws.com (3.15.178.106)' can't be established.
ECDSA key fingerprint is SHA256:OZTduss399swOwznuStj/5rmrrfnCPD9lMbI4FHuPEO.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'ec2-3-15-178-106.us-east-2.compute.amazonaws.com,3.15.178.106' (ECDSA) to the list of known hosts.
consume.py                                                                    100%  697     14.1KB/s   00:00
```

6) Then scp the file log4j.properties to your EMR cluster master node.

```
yashmp.21197@Yashu MINGW64 /d/Academic/Sem-4/CSP-554-BDT/Assignments/Assignment-10
$ scp -i emr-key-pair.txt log4j.properties hadoop@ec2-3-15-178-106.us-east-2.compute.amazonaws.com:/home/hadoop
log4j.properties                                                              100% 3199     62.7KB/s   00:00
```

7) Open two terminal sessions to the EMR master node. We will call one the EC2-1 window and the other the EC2-2 window.

8) In the EC2-1 window enter the following:

sudo cp ./log4j.properties /etc/spark/conf/log4j.properties

This changes the logging properties to turn off "INFO" messages to allow easier viewing of the results of the stream processing job. But it is not something you always want to disable.

```
[hadoop@ip-172-31-36-29 ~]$ sudo cp ./log4j.properties /etc/spark/conf/log4j.properties
[hadoop@ip-172-31-36-29 ~]$ |
```

9) In the EC2-1 window enter the following command to open a TCP (socket) connection on port 3333

nc -lk 3333

```
[hadoop@ip-172-31-36-29 ~]$ nc -lk 3333
```

10) In the EC2-2 window enter the following command:

spark-submit consume.py

This takes a while to start up. So, wait for some messages issued to the console before continuing. Note, when you do this you might see a message beginning with "WARN StreamingContext:..." which you can ignore.

```
[hadoop@ip-172-31-36-29 ~]$ spark-submit consume.py
21/03/30 02:06:15 WARN StreamingContext: Dynamic Allocation is enabled for this application. Enabling Dynamic allocation for Spark Streaming applications can cause data loss if Write Ahead Log is not enabled fo
r non-replayable sources like Flume. See the programming guide for details on how to enable the Write Ahead Log.
```

11) Now in the EC2-1 window enter one or more lines of text and press Enter/Return after each one including the last. You should see the word count results scroll by in the EC2-2 window

```
my name is yash patel
i study in iit chicago
i am doing my big data assignment 10

-------------------------------------------
Time: 2021-03-30 02:08:00
-------------------------------------------
('name', 1)
('is', 1)
('patel', 1)
('my', 1)
('yash', 1)

-------------------------------------------
Time: 2021-03-30 02:08:10
-------------------------------------------
('i', 1)
('in', 1)
('study', 1)
('iit', 1)
('chicago', 1)

-------------------------------------------
Time: 2021-03-30 02:08:20
-------------------------------------------

-------------------------------------------
Time: 2021-03-30 02:08:30
-------------------------------------------
('i', 1)
('am', 1)
('10', 1)
('doing', 1)
('my', 1)
('big', 1)
('data', 1)
('assignment', 1)
```

12) Remember to terminate your EMR instance after you are done!