

Assignment 10 CSP 554

Performing SSH to the cluster

[illegible]

Making Changes in consume.py file.

```
Users > rishabhjain > Desktop > consume.py
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("ec2-44-203-246-127.compute-1.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
```

Rishabh Jain || A20495530 || rjain35@hawk.iit.edu

Performing SCP for sending the files in cluster

```
➔ Desktop scp -i /Users/rishabhjain/Desktop/new-key-pair-emr.pem /Users/rishabhjain/Desktop/consume.py hadoop@ec2-44-203-246-127.compute-1.amazonaws.com:/home/hadoop
consume.py
100% 691 1.8KB/s 00:00
➔ Desktop scp -i /Users/rishabhjain/Desktop/new-key-pair-emr.pem /Users/rishabhjain/Desktop/log4j.properties hadoop@ec2-44-203-246-127.compute-1.amazonaws.com:/home/hadoop
log4j.properties
100% 3199 79.2KB/s 00:00
```

Terminal 1 Input

```
E::::E      EEEEE M:::::M      M:::::M RR::::R      R::::R
E::::E      M:::::M::M      M::M:::::M      R:::R      R::::R
E:::::EEEEEEEEEE M:::::M M:::M M:::M M:::::M      R:::RRRRRR:~::~R
E:::::~::~E      M:::::M M:::M::M      M:::::M      R:::~::~RR
E:::::EEEEEEEEEE M:::::M      M:::::M      M:::::M      R:::RRRRRR:~::~R
E:::::E      M:::::M      M:::M      M:::::M      R:::R      R::::R
E:::::E      EEEEE M:::::M      MMM      M:::::M      R:::R      R::::R
EE:::::EEEEEEEE~::~E M:::::M      M:::::M      R:::R      R::::R
E:::::~::~E M:::::M      M:::::M RR::::R      R::::R
EEEEEEEEEEEEEEEEEE MMMMMM      MMMMMM RRRRRRR      RRRRRR

[hadoop@ip-172-31-0-139 ~]$ sudo cp ./log4j.properties /etc/spark/conf/log4j.prop
[hadoop@ip-172-31-0-139 ~]$ nc -lk 3333
My name is Rishabh
```

Terminal 2 Output

```
22/04/03 22:52:30 INFO BlockManagerInfo: Added broadcast_25_piece0 in memory on t
B)
22/04/03 22:52:30 INFO TaskSetManager: Finished task 0.0 in stage 46.0 (TID 93) t
22/04/03 22:52:30 INFO YarnScheduler: Removed TaskSet 46.0, whose tasks have all
22/04/03 22:52:30 INFO DAGScheduler: ResultStage 46 (runJob at PythonRDD.scala:15
22/04/03 22:52:30 INFO DAGScheduler: Job 23 finished: runJob at PythonRDD.scala:1
-----
Time: 2022-04-03 22:52:30
-----
('name', 1)
('is', 1)
('Rishabh', 1)
('My', 1)
```

Terminal 1 Input

```

E:::EEEEEEEEEE M:::M M::M::M M:::M R:::RRRRRR:::R
E:::EEEEEEEEEE M:::M M:::M M:::M R:::RRRRRR:::R
E:::E M:::M M::M M:::M R::R R:::R
E:::E EEEEE M:::M MMM M:::M R::R R:::R
EE:::EEEEEEEE:::E M:::M M:::M R::R R:::R
E:::EEEEEEEEEE M:::M M:::M RR:::R R:::R
EEEEEEEEEEEEEEEEEE MMMMMM MMMMMM RRRRRRR RRRRRR

[hadoop@ip-172-31-0-139 ~]$ sudo cp ./log4j.properties /etc/spark/conf/log4j.prop
[hadoop@ip-172-31-0-139 ~]$ nc -lk 3333
My name is Rishabh
I am graduate student at IIT Chicago

```

Terminal 2 Output

```

22/04/03 22:53:00 INFO TaskSetManager: Finished task 0.0 in stage 58.0 (TID 100)
22/04/03 22:53:00 INFO YarnScheduler: Removed TaskSet 58.0, whose tasks have all
22/04/03 22:53:00 INFO DAGScheduler: ResultStage 58 (runJob at PythonRDD.scala:1
22/04/03 22:53:00 INFO DAGScheduler: Job 29 finished: runJob at PythonRDD.scala:
-----
Time: 2022-04-03 22:53:00
-----
('am', 1)
('graduate', 1)
('student', 1)
('at', 1)
('IIT', 1)
('Chicago', 1)
('I', 1)

```