

# MSDS 7346

## Cloud Computing

### Mini Project 6 – Log Analysis using AWS EMR

Name: Mooyoung Lee

#### Log Analysis Using EMR

In the live session we will review the building blocks of the MapReduce application. The application should take syslogformatted log records as input and determine the frequency of log events using Amazon EMR to count the number of records per second. I have provided a sample syslog-formatted file as an input.

Although we are keeping this application primarily focused on log analysis, but counting and frequency analysis has many known uses in other data analysis situations. The MapReduce application is performing what is considered a summarization design pattern by simply summing up the values of a common key. Other real-world applications of this technique are:

#### Load or usage analysis

Many times it is useful to know how many users access a server or a website throughout a time period. Web access logs or application logs that include the time stamps of user events could be imported and processed with a similar MapReduce application to determine usage frequency. At this time we will not be doing load or usage analysis in this lab.

Question 1 : The objective of this lab is to get hands-on experience with developing and deploying MapReduce applications on AWS Elastic Map Reduce. In this lab you are asked to develop two simple MapReduce applications using Pig and Python

1) Develop MapReduce Application to count number of events per second in the the sample log file (provided) using Python.

- Create an EMR and add Steps.
- Select 'Streaming program' as a Step Type
- Select appropriate files from S3 bucket

The 'Add step' dialog box is shown with the following fields and values:

- Step type:** Streaming program
- Name:** Streaming program
- Mapper:** s3://logemr2018/mapper2.py
- Reducer:** s3://logemr2018/reducer2.py
- Input S3 location:** s3://logemr2018/Syslog-Sample
- Output S3 location:** s3://logemr2018/output\_py.txt
- Arguments:** (Empty text area)
- Action on failure:** Continue

Buttons: Cancel, Add

2) Develop MapReduce Application to count number of events per second in the the sample log file (provided) using Pig.

- Add Steps.
- Select 'Pig program' as a Step Type
- Select appropriate files from S3 bucket

The 'Add step' dialog box is shown with the following fields and values:

- Step type:** Pig program
- Name:** Pig program
- Script S3 location:** s3://logemr2018/WordCount.pig
- Input S3 location:** s3://logemr2018/Syslog-Sample
- Output S3 location:** s3://logemr2018/output\_pig.txt
- Arguments:** (Empty text area)
- Action on failure:** Continue

Buttons: Cancel, Add

3) Deploy both applications in AWS EMR cluster

Cluster: firstEMR Running Running step

Summary Application history Monitoring Hardware Events Steps Configurations Bootstrap actions

Add step Clone step Cancel step No job available

Steps

Filter: All steps  7 steps (all loaded)

	ID	Name	Status	Start time (UTC-5)	Elapsed time
	s-I5IK4R35EYU	Streaming program	Pending		
	s-1LAU383WDVKB7	Pig program	Running	2018-07-03 13:15 (UTC-5)	24 seconds
	s-3HJ6KETM2MC5Z	pythonMapReduce	Failed	2018-07-03 13:14 (UTC-5)	6 seconds
	s-35U9884AJ2PMQ	Pig program	Failed	2018-07-03 13:12 (UTC-5)	16 seconds
	s-20CN8D VWASKQN	Streaming program	Failed	2018-07-03 13:06 (UTC-5)	8 seconds
	s-1K7FGVOWO5NLW	Custom JAR	Failed	2018-07-03 12:56 (UTC-5)	0 seconds
	s-C71S69Z2HNZA	Setup hadoop debugging	Completed	2018-07-03 12:15 (UTC-5)	2 seconds

#### 4) Submit screenshots

Submission: Submit different screen shots to show completion of each steps

Collaborators: None.

Resources: Class lecture