

MIE1628: Cloud-Based Data Analytics

Alexander Cheng – 1001634298 – Assignment 1

MIE1628 Assignment 1: MapReduce

Question 1: Implement a Map Reduce Function program for counting the number of lines in a document. Use 'shakespeare.txt' file, download it from Quercus. Please submit output files with code.

Output from Command Prompt:

```
Administrator: Command Prompt

C:\Users\alex_\Downloads\Assignment 1\WordCount>type shakespeare-1.txt |python mapper.py |python reducer.py
Number of Lines:
58483

C:\Users\alex_\Downloads\Assignment 1\WordCount>
```

Output from Hadoop: Issue with Hadoop, do not have enough time to troubleshoot it.

```
Administrator: Command Prompt

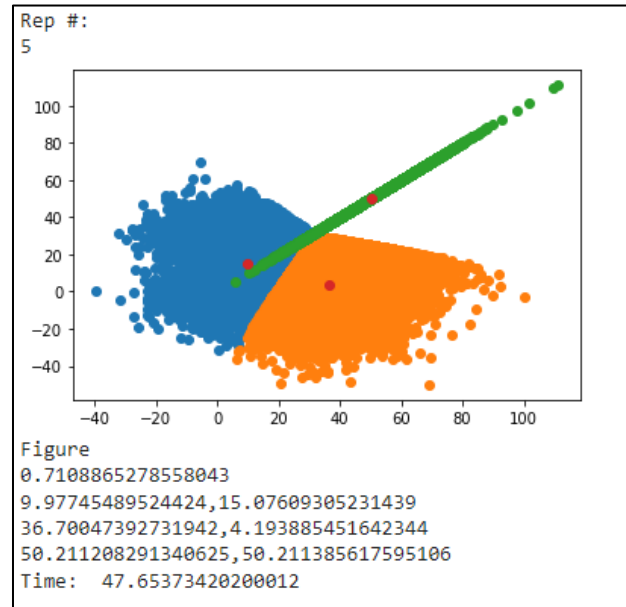
2022-09-26 00:48:51,681 INFO mapreduce.Job: Task Id : attempt_1664165500665_0001_m_000001_2, Status : FAILED
Error: java.lang.RuntimeException: Error in configuring object
    at org.apache.hadoop.util.ReflectionUtils.setJobConf(ReflectionUtils.java:115)
    at org.apache.hadoop.util.ReflectionUtils.setConf(ReflectionUtils.java:81)
    at org.apache.hadoop.mapred.MapTask.runOldMapper(MapTask.java:462)
    at org.apache.hadoop.mapred.MapTask.run(MapTask.java:349)
    at org.apache.hadoop.mapred.YarnChild$2.run(YarnChild.java:178)
    at java.security.AccessController.doPrivileged(Native Method)
    at javax.security.auth.Subject.doAs(Subject.java:422)
    at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1845)
    at org.apache.hadoop.mapred.YarnChild.main(YarnChild.java:172)
Caused by: java.lang.reflect.InvocationTargetException
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.ReflectionUtils.setJobConf(ReflectionUtils.java:112)
    ... 9 more
Caused by: java.lang.RuntimeException: Error in configuring object
    at org.apache.hadoop.util.ReflectionUtils.setJobConf(ReflectionUtils.java:115)
    at org.apache.hadoop.util.ReflectionUtils.setConf(ReflectionUtils.java:81)
    at org.apache.hadoop.mapred.MapRunner.configure(MapRunner.java:38)
    ... 14 more
Caused by: java.lang.reflect.InvocationTargetException
    at sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
    at sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
    at sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.ReflectionUtils.setJobConf(ReflectionUtils.java:112)
    ... 17 more
Caused by: java.lang.RuntimeException: java.io.FileNotFoundException: HADOOP_HOME and hadoop.home.dir are unset. -see https://wiki.apache.org/hadoop/WindowsProblems
    at org.apache.hadoop.util.Shell.getEnvPath(Shell.java:726)
    at org.apache.hadoop.util.Shell.setSetPermissionCommand(Shell.java:271)
    at org.apache.hadoop.fs.FileUtil.chmod(FileUtil.java:1111)
    at org.apache.hadoop.fs.FileUtil.chmod(FileUtil.java:1092)
    at org.apache.hadoop.streaming.PipeMapper.configure(PipeMapRed.java:180)
    at org.apache.hadoop.streaming.PipeMapper.configure(PipeMapper.java:66)
    ... 22 more
Caused by: java.io.FileNotFoundException: HADOOP_HOME and hadoop.home.dir are unset. -see https://wiki.apache.org/hadoop/WindowsProblems
    at org.apache.hadoop.util.Shell.fileNotFoundException(Shell.java:948)
    at org.apache.hadoop.util.Shell.getHadoopHomeDir(Shell.java:569)
    at org.apache.hadoop.util.Shell.getQualifiedBin(Shell.java:592)
    at org.apache.hadoop.util.Shell.<clinit>(Shell.java:489)
    at org.apache.hadoop.util.StringUtils.<clinit>(StringUtils.java:78)
    at org.apache.hadoop.conf.Configuration.getBoolean(Configuration.java:1689)
    at org.apache.hadoop.security.SecurityUtil.setConfigurationInternal(SecurityUtil.java:104)
    at org.apache.hadoop.security.SecurityUtil.<clinit>(SecurityUtil.java:98)
    at org.apache.hadoop.security.UserGroupInformation.initialize(UserGroupInformation.java:312)
    at org.apache.hadoop.security.UserGroupInformation.setConfiguration(UserGroupInformation.java:366)
    at org.apache.hadoop.mapred.YarnChild.main(YarnChild.java:164)
Caused by: java.io.FileNotFoundException: HADOOP_HOME and hadoop.home.dir are unset.
    at org.apache.hadoop.util.Shell.checkHadoopHomeInner(Shell.java:468)
    at org.apache.hadoop.util.Shell.checkHadoopHome(Shell.java:439)
    at org.apache.hadoop.util.Shell.<clinit>(Shell.java:516)
    ... 7 more

2022-09-26 00:48:51,681 INFO mapreduce.Job: map 100% reduce 100%
2022-09-26 00:48:55,297 INFO mapreduce.Job: Job Job 1664165500665_0001 failed with state FAILED due to: Task failed task_1664165500665_0001_m_000000
Job failed as tasks failed. failedMaps:1 failedReduces:0 killedMaps:0 killedReduces:0

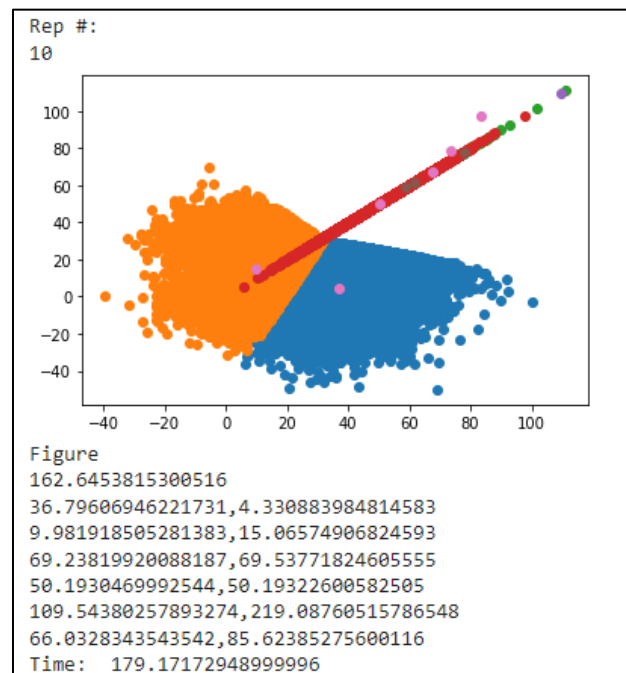
2022-09-26 00:48:55,488 INFO mapreduce.Job: Counters: 14
Job Counters
    Failed map tasks=7
    Killed map tasks=1
    Killed reduce tasks=1
    Launched map tasks=8
    Other local map tasks=0
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=32366
    Total time spent by all reduces in occupied slots (ms)=0
    Total time spent by all map tasks (ms)=32366
    Total vcore-millisecseconds taken by all map tasks=32366
    Total megabyte-millisecseconds taken by all map tasks=33142784
Map-Reduce Framework
    CPU time spent (ms)=0
    Physical memory (bytes) snapshot=0
    Virtual memory (bytes) snapshot=0
2022-09-26 00:48:55,489 ERROR streaming.StreamJob: Job not successful!
Streaming Command Failed!
```

Question 2: Apply K-Means Clustering on MapReduce using k=3 and k=6 clusters on the given dataset, list the cluster labels or centroids, the number of iterations for convergence or use maximum iterations = 10 and time/duration.

I tried implementing it on Java, but since I am not familiar with Java I switched to Python. I had trouble implementing MapReduce on Python using the command prompt. I used Google Colab and wrote what can be implemented into either a Mapper or a Reducer function for the command prompt and Hadoop.



Output for K= 3 Clusters



Output for K=6 Clusters

First step is to import the initial centroid positions text file and convert it into an array. We then make clusters based on which data points are closest to which clusters and assign each data point the cluster label. The new centroids are calculated by taking the average of all the data points that belong in the cluster label. We then edit the text file to those new centroid positions and loop the map and reducer function. There are two stopping criteria used, either the clusters converge in which the clusters do not change much between each iteration, or a maximum number of iterations (10) has been reached. The time for each calculation has also been reached, we can observe that having 6 clusters takes significantly more time to complete compared to the program with only 3 clusters. We can see that as the number of clusters goes up, it will take exponentially more time to run in which we need help of Hadoop to run in parallel.

Question 3: Explain advantages and disadvantages of using K-Means Clustering with MapReduce:

Advantages:

- MapReduce separates the task into two tasks – Map and Reduce that can be used iteratively and in parallel when multiple clusters are required.
- The reducer function does the shuffling of input pairs with the key values
- Computationally much faster

Disadvantages:

- Hard to optimize hyperparameter k in MapReduce
- Iterative process of MapReduce is hard to implement
- MapReduce requires the input data format in the same format

Question 4: Can we reduce the number of distance comparison by applying Canopy Selection? Which distance metric should we use for the canopy selection and why?

We can reduce the number of distance comparisons by applying Canopy Selection. Canopy Selection allows us to significantly reduce the number of distance comparisons by filtering out the number of calculations, we would only have to calculate the distances between the points that are within the same canopy in which a canopy is when we divide the data into overlapping subsets.

This is significantly less calculations than calculating the distance between all the data points of the dataset. The Canopy Selection can be ran using a cheap and approximate similarity measure, such as the inverted index, which is commonly used as a cheap distance metric and provides high accuracy with lower computational cost when a small number of features are sufficient to build canopies.

Question 5: Is it possible to implement Canopy Selection on MapReduce? If yes, then explain in words, how would you implement it?

It is possible to implement Canopy Selection on MapReduce. We can implement it by sending the list of data points as an input for the Mapper function. We can also set the threshold values T1 and T2 inside the

mapper with $T1 > T2$. We can use a while loop and randomly select a data point in which we find the distance between that data point and the rest of the data points. The points within the distance of $T1$ is set as a Canopy, and the points with the distance less than $T2$ will be grouped together with the data point to form a key-value pair in which those data points are removed from the data set list. The while loop is run again to form a new cluster until there are no more data points inside the list. The output of the mapper will be the cluster location and the key value pairs (data points that belong to which cluster).

Question 6: Is it possible to combine the Canopy Selection with K-Means on MapReduce? If yes, then explain in words, how would you do that?

It is possible to combine the Canopy Selection with K-Means on MapReduce. The Mapper function will take the list of data points and the number of clusters (k) with initial centroid positions as an input value, the initial centroid positions can be randomly selected as one of the data points. We can also set the threshold values $T1$ and $T2$ inside the mapper with $T1 > T2$. We can use a while loop and randomly set an initial centroid which we find the distance between that centroid and the rest of the data points. The points within the distance of $T2$ is set as a Canopy, and the points with the distance less than $T1$ will be grouped together with the data point to form a cluster in which those data points are removed from the data set list and the number of clusters k goes up. The while loop is run again to form a new cluster until there are no more data points inside the list or if the number of clusters k has been met. If the list isn't empty after k clusters, the distance between the remaining data points and the selected centroids will be calculated and they will be assigned to the closest centroid. The output of the mapper will be the centroid locations and the data points that belong to each cluster. The reducer will take the intermediate key/value pairs as an input and the new centroid for each cluster will be calculated. The list of centroids will also be updated and sent as a loop for MapReduce and will repeat until the stopping criteria is met (such as no data points changes, number of iterations passed, etc.).