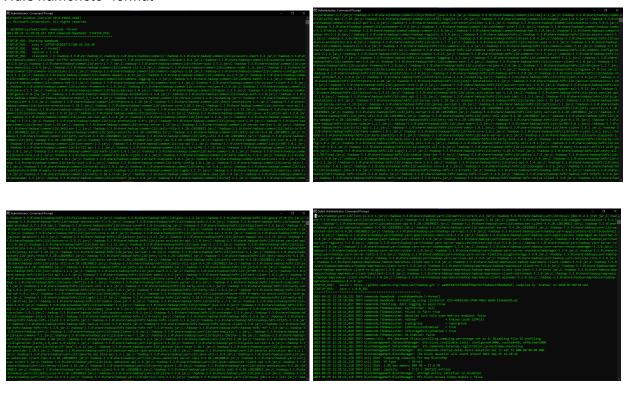
Assignment 1- Khaled Gaber - #1004144302

Questions:

1) [Marks: 15] Implement a Map Reduce program for counting the number of lines in a document. Use the 'shakespeare.txt' file and download it from Quercus. Please submit input/output files with code

The number of lines (total including blank lines) is 58483

Hdfs namenote -format





```
This control was a second form of the property of the second form of the property of the prope
```

start-dfs.cmd start-yarn.cmd

Hdfs dfs -mkdir -p inputtext

```
See Administration of the Control of
```

2) [Marks: 45] Apply K-means clustering on Map Reduce 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 = 20 and time/duration.

Here it is working locally: type data_points.txt | python3 mapper.py | sort | python3 reducer.py

k=3

```
Microsoft Windows [Version 10.0.22621.2283]
(c) Microsoft Corporation. All rights reserved.

D:\Documents\MIE1628\A1Q2\local testing>python3 driver_local.py
For K = 3
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
Iteration 6
Iteration 7
Iteration 8
K-means clustering completed!
Converged after 8 iterations.
['0,9.963012568253967,15.103543765690867\r', '1,49.99793261993235,30.103869640898814\r', '2,35.01195298036856,1.77652196
2410613']
Time taken: 53.50 seconds

D:\Documents\MIE1628\A1Q2\local testing>
```

k=6

```
D:\Documents\MIE1628\A1Q2\local testing>python3 driver_local.py
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 6
Iteration 7
Iteration 8
Iteration 9
Iteration 10
Iteration 11
Iteration 12
Iteration 13
Iteration 14
Iteration 15
Iteration
Iteration 17
Iteration 18
Iteration 19
Iteration 20
 K-means clustering completed!
Reached maximum iterations (20) without convergence.
['0,9.855654929627354,10.298995233546416\r', '1,39.46324964371489,13.166620703765759\r', '2,35.1832318085002,-6.17827973
44764315\r', '3,50.19781046472902,30.506998555793306\r', '4,34.64635359882776,2.5917956197877707\r', '5,10.0087964462727
31,18.19712902815192']
Time taken: 151.59 seconds
D:\Documents\MIE1628\A1Q2\local testing>
```

Here it is on MapReduce

k=1 (Command Testing)

hadoop jar C:\hadoop-3.3.0\share\hadoop\tools\lib\hadoop-streaming-3.3.0.jar -mapper "C:\Users\khale\AppData\Local\Microsoft\WindowsApps\python3.exe

C:\MIE1628\A1Q2\mapper.py" -reducer

"C:\Users\khale\AppData\Local\Microsoft\WindowsApps\python3.exe

C:\MIE1628\A1Q2\reducer.py" -input /input/kmeans/* -output /output/kmeans/results

```
| April | Deciding | Per | Valuation | April | Deciding | Per | Valuation | April | Deciding | Per | Valuation | Per | Per | Deciding | Per | Pe
```

```
| 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.100.00 | 20.1
```

```
2023-09-28 22:05:06,865 INFO streaming.StreamJob: Output directory: /output/kmeans/results
 C:\MIE1628\A1Q2>hdfs dfs -ls \
 ound 4 items

      drwxr-xr-x
      - khale supergroup
      0 2023-09-28 21:59 /input

      drwxr-xr-x
      - khale supergroup
      0 2023-09-28 17:55 /output

      drwx-----
      - khale supergroup
      0 2023-09-28 17:42 /tmp

      drwxr-xr-x
      - khale supergroup
      0 2023-09-28 17:26 /user

C:\MIE1628\A1Q2>hdfs dfs -get \output\* C:\MIE1628\A1Q2\Results
C:\MIE1628\A1Q2>hdfs dfs -get \input\* C:\MIE1628\A1Q2\Results
C:\MIE1628\A1Q2>hdfs dfs -cat \output\kmeans\part-00000
cat: `outputkmeanspart-00000': No such file or directory
C:\MIE1628\A1Q2>hdfs dfs -ls \output\kmeans
Found 1 items
                 - khale supergroup
                                                        0 2023-09-28 22:05 /output/kmeans/results
drwxr-xr-x
C:\MIE1628\A1Q2>hdfs dfs -ls \output\kmeans\results
 ound 2 items
 -rw-r--r-- 1 khale supergroup
-rw-r--r-- 1 khale supergroup
                                                          0 2023-09-28 22:05 /output/kmeans/results/_SUCCESS
                                                 119 2023-09-28 22:05 /output/kmeans/results/part-00000
 C:\MIE1628\A1Q2>hdfs dfs -cat \output\kmeans\results\part-00000
0,9.921463822213996,14.994793104788538
1,43.62540411580502,19.74851471321919
2,37.01713640846856,-2.726666152642699
```

k=3

This took 9 iterations in 120.76 seconds, and the finals centroids were:

The Final Centroids are: [

('0', (9.962931682425607, 15.10369315620118)),

('1', (49.99793261993191, 30.103869640898743)),

('2', (35.01195298036899, 1.7765219624106197))]

Note due to issues with running this, I followed the solution on Piazza, question @49, changing the property value from yarn to local.

The initial centroids were taken as the first 3 points from "data points.txt"

The output for this is very long, and so I only show screenshots of the input, one iteration, and the final below:

```
09-28 22:53:48,259 INFO mapreduce.lob: maj 0% reduce 0%
09-28 22:53:48,259 INFO mapreduce.lob: maj 0% reduce 0%
09-28 22:53:48,469 INFO streaming.PipeMapRed: R/W/S=300000/297646/0 in:30
        MilliantOpythoma driver,y
thing Main()
addig the data_points.txt file
addig the data_points.txt file
inlining intid centroid value
ting old results file
'/output/weems/results' in bosh file or directory
undigs_txt.depart/results' in Such file or directory
undigs_txt.depart/results' in Such file or directory
undigs_txt.depart/results' in Such file exists
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  -28 22:53:48,791 INFO streaming.PipeMapRed: R/W/S=400000/396925/0 in:40
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  -28 22:53:49,116 INFO streaming.PipeMapRed: R/W/S=500000/496418/0 in:500000=500000/1 [rec/s] out:496418=496418
        //octput/messent/results* I Ms such file or directory
//octput/messent/results* I Ms such file or directory
//optput/messent/results* I file edits
ting iteration i for k = 3
//optput/messent/results* I file edits
ting iteration i for k = 3
//optput/messent/results* I file edits
ting iteration i for k = 3
//optput/messent/results* I file edits
//optput/messent/results
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  -28 22:53:50,098 INFO streaming.PipeMapRed: R/W/S-800000/797936/0 in:400000-8000000/2 [rec/s] out:398968-79793
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              //s]
0-32 131510,702 DNO Streaming_Pleasaged(HittrorThread done
0-32 231510,702 DNO Streaming_Pleasaged(HittrorThread done
0-32 231510,702 DNO Streaming_College(Hittle)
0-32 231510,702 DNO STREAMING_COLLEGE(HITLE)
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             // Journal of the control of the con
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       ;
28 22:53:51,007 INFO mapred.MapTask: Finished spill 0
28 22:53:51,014 INFO mapred.Task: Task:attempt_local405568936_0001_m_0000000_0 is done. And is in
        -89-28 (25:25:47,380 HP mprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomprocomproco
                    09-28 22:53:47,356 INFO Configuration.deprecation: mapred.task.id is deprecated. Instead, use mapredu
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      : Faise
sedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux
Using ResourceCalculatorProcessTree : org.apache.hadoop.yarn.util.WindowsB
                                             .78 22:53:47,358 TNFO Configuration.deprecation: mapred.job.id is deprecated. Instead, use mapreduce.job.id
28 22:53:47,358 TNFO Configuration.deprecation: user.name is deprecated. Instead, use mapreduce.job.user.name
28 22:53:47,358 TNFO Configuration.deprecation: mapred.task.partition is deprecated. Instead, use mapreduce.tu
                                             22 33347,483 IMO STRENBERG PHENDERS (MAYS-100) EINE [FE/5] OLT ME [FE/5]
            0-72 2/3514/09 EDG treaming Piphengheid Piphengheid (Piphengheid Ster [1,1107-1443] Oct.NA [recf] Oc
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      urs
3-09-28 22:53:51,286 INFO mapred.Merger: Merging 2 sorted segments
3-09-28 22:53:51.286 INFO mapred.Merger: Down to the last merge-pass, with 2 segments left of to
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              -09-28 22:53:51,294 INFO mapreduce.Job: map 100% reduce 0%
-09-28 22:53:51,451 INFO reduce.MergeManagerImpl: Merged 2 segments, 40938093 bytes to disk
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              tt

-09-28 22:53:51,452 DNFO reduce.MergeManagerEmpl: Merging 1 files, 48938095 bytes from disk
-09-28 22:53:51,452 DNFO reduce.MergeManagerEmpl: Merging 0 segments, 0 bytes from memory into
-09-28 22:53:51,452 DNFO mapred.Merger: Merging 1 sorted segments
-09-28 22:53:51,452 DNFO mapred.Merger: Down to the last merge-pass, with 1 segments left of t
13-09-28 2233313,178 1MF0 magned.MepTask. Spilling map output do. 72; bufreld a justicide | 188257680 | 189-28 2233313,178 1MF0 magned.MepTask Spilling map output do. 72; bufreld a justicide | 189-28 2233313,173 1MF0 magned.MepTask Spilling map output do. 72; bufreld a justicide | 189-28 223331,183 1MF0 magned.MepTask Finithed spill a | 189-28 223331,183 1MF0 magned.MepTask Finithed spill a | 189-28 223331,183 1MF0 magned.MepTask Finithed spill a | 189-28 223331,183 1MF0 magned.MepTask Finithed | 189-28 233331,183 1MF0 magned.MedTask Finithed | 189-28 233331,183 1MF0 ma
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          22 22 23 23 34 49 UNIC Configuration degreeation superal job resider is deprecated. Interest 22 22 23 23 34 49 UNIC Configuration degreeation superal job resider in general process. The process of the 
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          8
214 INFO mapred.ReduceTask: Using ShuffleConsumerPlugin: org.apache.hadoop.m
                                                                                                  .59555,215 WARM impl.HetricsSystemImpl: ObbTracker metrics system already initialized in 51551,225 WARD reduce.MeroglemagerEmpl: ReggerHamager: memoryLimit=34938464, massIngl: Modice28665932, doortexforts], demindemeroglemagerinemologi-36 modice28665932, doortexforts], demindemeroglemagerinemologi-36 modice28665932, doortexforts; massIngl: Local86569592, doortexforts
                                                                                                      .53:51,277 INFO reduce.Mergedwangerisp1: closeInMemoryFile -> map-output of Size: 40938013, in 1, committemory -> 8, usselmeory ->40938013 in 1, committemory -> 8, usselmeory ->80938013 in 1, committemory -> 8, usselmeory 
                    .9
99-28 22:53:51,280 INFO reduce.WergeManagerImpl: closeInMemoryFile -> map-output of size: 80, in
>> 2, commitMemory -> 40938013, usedMemory -> 40938093
90-28 22:53:51,281 INFO reduce.EventFetcher: EventFetcher is interrupted.. Returning
```

```
DAY, DESCRIPTION OF THE PROPERTY OF THE PROPER
```

```
2023-09-28 22:55:40,754 INFO streaming.StreamJob: Output directory: /output/kmeans/results
New_centroids_output: 0,9.962931682425607,15.10369315620118
1,49.99793261993191,30.103869640898743
2,35.01195298036899,1.7765219624106197
New_centroids: [('0', (9.962931682425607, 15.10369315620118)), ('1', (49.99793261993191, 30.103869640898743)), ('2', (3 5.01195298036899, 1.7765219624106197))]
Convergence: True
Iteration 9 completed in 9.318159341812134 seconds
9 Iterations in: 120.76522064208984 seconds
The Final Centroids are: [('0', (9.962931682425607, 15.10369315620118)), ('1', (49.99793261993191, 30.103869640898743)), ('2', (35.01195298036899, 1.7765219624106197))]
```

Output files were saved using:

hdfs dfs -get /output/kmeans/results/* C:\MIE1628\A1Q2\Final_Results\k3

k=6

```
This did not converge, maxing out at 20 iterations. 20 Iterations in: 291.4041225910187 seconds
The Final Centroids are: [
('0', (9.855663875034324, 10.30023604782646)),
('1', (39.46324964371434, 13.166620703765755)),
('2', (35.18323180850043, -6.17827973447644)),
('3', (50.19781046472971, 30.506998555793103)),
('4', (34.64635359882755, 2.5917956197877845)),
('5', (10.008599091020962, 18.198393029289274))]
```

The initial centroids were taken as the first 6 points from "data points.txt"

```
C:\MIE1628\A1Q2>python3 driver.py
Starting Main()
uploading the data_points.txt file
put: `/input/kmeans/data_points.txt': File exists
initializing initial centroid values
Deleting old results file
Deleted /output/kmeans/results
uploading the centroids.txt file
put: `/input/kmeans/centroids.txt': File exists
Starting iteration 1 for k = 6
```

```
New_centroids: [('0', (9.855663875034324, 10.30023604782646)), ('1', (39.46324964371434, 13.166620703765755)), ('2', (35.18323180850043, -6.17827973447644)), ('3', (50.19781046472971, 30.506998555793103)), ('4', (34.64635359882755, 2.5917956197877845)), ('5', (10.008599091020962, 18.198393029289274))]
Convergence: False
Iteration 20 completed in 10.352383136749268 seconds
Saving Output file
20 Iterations in: 295.6279318332672 seconds
The Final Centroids are: [('0', (9.855663875034324, 10.30023604782646)), ('1', (39.46324964371434, 13.166620703765755)), ('2', (35.18323180850043, -6.17827973447644)), ('3', (50.19781046472971, 30.506998555793103)), ('4', (34.64635359882755, 2.5917956197877845)), ('5', (10.008599091020962, 18.198393029289274))]
```

3) [Marks: 10] Explain the advantages and disadvantages of using K-Means Clustering with MapReduce.

Advantages:

- Able to handle expansive datasets that are beyond the storage capacity of a single machine.
- Work parallelization that results in significantly enhanced clustering speed.
- Robust fault tolerance mechanism: in case a node experiences a failure during computation, the framework autonomously redistributes the workload to other operational nodes, ensuring the processing advances without disruption.

Disadvantages:

- Increased complexity: Implementing K-means using MapReduce is slightly more complex than running it locally
- Overhead costs: The process entails multiple phases including data partitioning, mapping, shuffling, and reducing, which incur additional computational and communication overheads. In situations where data processing tasks are small or simplistic, the overhead associated with configuring and operating the MapReduce framework could just not be worth it
- Since we need to divide the tasks into 2 parts, the Map and Reduce functions, this might not be suitable for all types of tasks, possibly limiting its suitability for certain clustering or analytical challenges.

Please read the paper which is provided with the assignment in the Quercus and answer the following questions.

- 4) [Marks: 10] Can we reduce the number of distance comparisons by applying the Canopy Selection? Which distance metric should we use for the canopy clustering and why?
- Yes, Canopy Selection can significantly diminish the number of distance comparisons required. The primary concept revolves around economically partitioning the data into subsets initially, followed by computing the distances solely between points within the same subset, thereby lowering the overall number of distance computations needed.
- The distance metric we decide to use is dependant on the nature of the data, as is task-specific

5) [Marks: 10] Is it possible to apply Canopy Selection on MapReduce? If yes, then explain in words, how would you implement it.

Yes it is, we can structure it in two primary phases: Canopy creation and allocation of data points.

Canopy Creation

- Mappers process segments of data points to generate initial canopies using a specified algorithm.
 - A random point is selected to initiate a canopy.
 - The algorithm iterates through remaining points:
 - Points within distance < T1 from the initial point are added to the canopy.
 - If the distance is also < T2, the point is removed from the set.
 - This loop continues until all points are either clustered or removed, generating a set of canopies.
 - Mappers return the centroids of these initial canopies.
- Reducers receive the initial centroids:
 - Apply the canopy measure and thresholds once more to refine the canopy centroids.
 - Output a final set of canopy centroids for the next phase.

Clustering Phase:

- In this phase, Mappers are fed the finalized canopy centroids, utilizing them to classify the data points, subsequently returning pairs composed of (canopy centroid identifier, data point).
- Reducers then consolidate data based on the canopy centroids' identifiers, culminating in the output of the final canopies.
- 6) [Marks: 10] Is it possible to combine the Canopy Selection with K-Means on MapReduce? If yes, then explain in words, how would you do that.

Yes it is possible, and would likely result in a speedup. The main potential steps are outlined below:

- 1. Canopy Selection Implementation:
 - Execute the canopy selection as described in Question 5 to partition the data into rough clusters or canopies initially.
- 2. K-means Initialization:
 - Within each canopy, initialize k prototypes.
- 3. K-means Iterations:
- Mappers:

- Each data point is associated with a list of canopies it belongs to, and consequently, a list of candidate prototypes.
- Mappers calculate the closest prototype for each point, outputting pairs of the closest prototype ID and the data point.

Combiners:

- This stage can aggregate data from mappers to reduce the amount of data transferred to reducers. They can compile lists of points associated with each prototype ID.

- Reducers:

- Post-combination, reducers receive pairs consisting of a prototype ID and a list of points associated with that prototype.
- The new centroids are computed by calculating the mean of the points within each cluster.
- Output the new centroids (prototype ID, mean of the points) for the iterations.