

Lab Course: Distributed Data Analytics
Exercise Sheet 5
Group 2 - Monday

Submitted by: Sruthy Annie Santhosh, 312213

Topic:Distributed Machine Learning (Supervised)

- My System

Hardware and Software SetUp:

Processor	M1 chip
Operating System	MacOS Monterey
Number of Cores	8
RAM	8 GB
Python version	3.8.8

The code editor used in this sheet : VSCode and Jupyter Notebook (graph)

The output images attached are screenshots from the terminal.

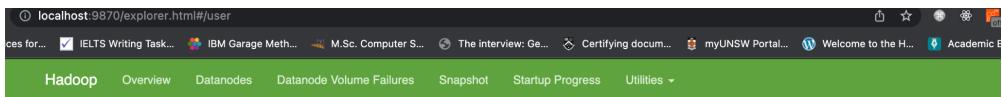
- Exercise 1: Preparing your Hadoop infrastructure

I have set up the Pseudo Distributed mode of Hadoop in my local system. Testing of the correct implementation was done by creating directories and putting files in the directories. The output of jps is as follows:

```
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % jps
1062 SecondaryNameNode
1257 ResourceManager
21722 Jps
1660 NameNode
1357 NodeManager
925 DataNode
(base) sruthysanthosh@Sruthys-MacBook-Air ~ %
```

```
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % hadoop version
Hadoop 3.3.3
Source code repository https://github.com/apache/hadoop.git -r d37586cbda38c338d9fe481addd5a05fb516f71
Compiled by stevel on 2022-05-09T16:36Z
Compiled with protoc 3.7.1
From source with checksum eb96dd4a797b6989ae0cdb9db6efc6
This command was run using /opt/homebrew/Cellar/hadoop/3.3.3/libexec/share/hadoop/common/hadoop-common-3.3.3.jar
(base) sruthysanthosh@Sruthys-MacBook-Air ~ %
```

Hadoop user-interface was also familiarised with:



Browse Directory

Browse Directory							
/user							
Go!							
Search: <input type="text"/>							
Show: <input type="button" value="25"/> entries	<input type="button" value=""/>						
<input type="checkbox"/>	<input type="button" value=""/>						
	Permission	Owner	Group	Size	Last Modified	Replication	Block Size
	drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 16:28	0	0 B
Showing 1 to 1 of 1 entries							
<input type="button" value="Previous"/> <input type="button" value="1"/> <input type="button" value="Next"/>							
Hadoop, 2022.							

Cluster Metrics																																																																															
<table border="1"> <tr> <td>Apps Submitted</td> <td>0</td> <td>Apps Pending</td> <td>0</td> <td>Apps Running</td> <td>1</td> <td>Apps Completed</td> <td>0</td> </tr> <tr> <td colspan="4">Cluster Nodes Metrics</td> <td colspan="4"></td> </tr> <tr> <td colspan="2">Active Nodes</td> <td colspan="2">Decommissioning Nodes</td> <td colspan="4"></td> </tr> <tr> <td colspan="2">1</td> <td colspan="2">0</td> <td colspan="4">0</td> </tr> <tr> <td colspan="8">Scheduler Metrics</td> </tr> <tr> <td colspan="4">Scheduler Type</td> <td colspan="4">Scheduling Resource Type</td> </tr> <tr> <td colspan="4">Capacity Scheduler</td> <td colspan="4"><memory:1024, vcores:1></td> </tr> <tr> <td colspan="8" style="text-align: center;">Show: 20 entries</td> </tr> <tr> <td>ID</td> <td>User</td> <td>Name</td> <td>Application Type</td> <td>Application Tags</td> <td>Queue</td> <td>Applicat</td> <td>Priorit</td> </tr> </table>								Apps Submitted	0	Apps Pending	0	Apps Running	1	Apps Completed	0	Cluster Nodes Metrics								Active Nodes		Decommissioning Nodes						1		0		0				Scheduler Metrics								Scheduler Type				Scheduling Resource Type				Capacity Scheduler				<memory:1024, vcores:1>				Show: 20 entries								ID	User	Name	Application Type	Application Tags	Queue	Applicat	Priorit
Apps Submitted	0	Apps Pending	0	Apps Running	1	Apps Completed	0																																																																								
Cluster Nodes Metrics																																																																															
Active Nodes		Decommissioning Nodes																																																																													
1		0		0																																																																											
Scheduler Metrics																																																																															
Scheduler Type				Scheduling Resource Type																																																																											
Capacity Scheduler				<memory:1024, vcores:1>																																																																											
Show: 20 entries																																																																															
ID	User	Name	Application Type	Application Tags	Queue	Applicat	Priorit																																																																								
<table border="1"> <tr> <td>application_1654625852682_0001</td> <td>sruthysanthosh</td> <td>streamjob4733257675330931132.jar</td> <td>MAPREDUCE</td> <td></td> <td>default</td> <td>0</td> <td></td> </tr> <tr> <td colspan="8" style="text-align: center;">Showing 1 to 1 of 1 entries</td> </tr> </table>								application_1654625852682_0001	sruthysanthosh	streamjob4733257675330931132.jar	MAPREDUCE		default	0		Showing 1 to 1 of 1 entries																																																															
application_1654625852682_0001	sruthysanthosh	streamjob4733257675330931132.jar	MAPREDUCE		default	0																																																																									
Showing 1 to 1 of 1 entries																																																																															

Word Count MapReduce Example

I have run the word count example using python mapper and reducer files. One mapper and one reducer were used. The given text file was first transferred to the hadoop directory using the -put command.

```
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as sruthysanthosh in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
localhost: datanode is running as process 925. Stop it first and ensure /tmp/hadoop-sruthysanthosh-datanode.pid file is empty before retry.
Starting secondary namenodes [Sruthys-MacBook-Air.local]
Sruthys-MacBook-Air.local: secondarynamenode is running as process 1062. Stop it first and ensure /tmp/hadoop-sruthysanthosh-secondarynamenode.pid file is empty before retry.
2022-06-07 20:18:27,754 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting resourcemanager
resourcemanager is running as process 1257. Stop it first and ensure /tmp/hadoop-sruthysanthosh-resourcemanager.pid file is empty before retry.
Starting nodemanagers
localhost: nodemanager is running as process 1357. Stop it first and ensure /tmp/hadoop-sruthysanthosh-nodemanager.pid file is empty before retry.
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % jps
1062 SecondaryNameNode
2217 Jps
1257 ResourceManager
1668 NameNode
1357 NodeManager
925 DataNode
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % cd Downloads
(base) sruthysanthosh@Sruthys-MacBook-Air Downloads % python version
python: can't open file 'version': [Errno 2] No such file or directory
(base) sruthysanthosh@Sruthys-MacBook-Air Downloads % python --version
Python 3.8.8
(base) sruthysanthosh@Sruthys-MacBook-Air Downloads % cd
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % hdfs dfs -mkdir /user
2022-06-07 20:28:39,077 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % hdfs dfs -mkdir /user/sruthy
2022-06-07 20:29:05,743 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % hdfs dfs -put /Users/sruthysanthosh/Downloads/text.txt /user/sruthy
```

In the mapper file

The lines from the input file are read and preprocessed to get each word. Then each word is sent as the output of the mapper along with the value 1. Mapper output- (word,1),(word,1).....

```
import sys

#Reading each line in the input file
for line in sys.stdin:
    line = line.strip() #Removing white space characters from beginning and end
    words = line.split() #Splitting based on space
    for word in words:
        #Sending each word and value 1 to reducer
        print('%s\t%s' % (word, 1))
```

Reducer file

In the reducer, mapper outputs after being sorted according to the key. Hence all same words would be together.

The counting of each word is done here.

```

# read the entire line from STDIN, which is the output of the mapper
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # splitting based on tab character as we wrote in mapper file
    word, count = line.split('\t'),1
    # convert count to int
    try:
        count = int(count)
    except ValueError:
        # count was not a number, discard
        continue

    # Hadoop sorts the map output by keys , ie, word before passing to reducer
    # Hence all same words will be together
    if current_word == word:
        #If the current word is same as new word, add count
        current_count += count
    else:
        if current_word:
            #end of one word, hence writing its count
            # write result to STDOUT
            #Current word and its count is overwritten after each word is over
            print('%s\t%s' % (current_word, current_count))
            #Take the new word as current words, and its count as current count
        current_count = count
        current_word = word

```

Output of Reducer is each word and its count

The command used for doing mapreduce is : The output is stored in a file in the hadoop directory names output. The paths of mapper and reducer are specified along with the file names.

```
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % hadoop jar /opt/homebrew/Cellar/hadoop/3.3.3/libexec/share/hadoop/tools/lib/hadoop-streaming-3.3.3.jar -input /user/sruthy/text.txt -output /user/sruthy/output -mapper mapper1.py -reducer reducer1.py -file /Users/sruthysanthosh/Downloads/mapper1.py -file /Users/sruthysanthosh/Downloads/reducer1.py
```

The console output is as follows:

```
(base) sruthysanthosh@Sruthys-MacBook-Air ~ % hadoop jar /opt/homebrew/Cellar/hadoop/3.3.3/libexec/share/hadoop/tools/lib/hadoop-streaming-3.3.3.jar -input /user/sruthy/text.txt -output /user/sruthy/output -mapper mapper1.py -reducer reducer1.py -file /Users/sruthysanthosh/Downloads/mapper1.py -file /Users/sruthysanthosh/Downloads/reducer1.py

2022-06-07 20:45:55,214 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
2022-06-07 20:45:55,502 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
packageJobJar: [/Users/sruthysanthosh/Downloads/mapper1.py, /Users/sruthysanthosh/Downloads/reducer1.py] [] /var/folders/9z/75yp7_69shg75lbf0tpmj9w000gn/T/streamjob1627353474031156597.jar tmpDir=null
2022-06-07 20:45:56,686 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2022-06-07 20:45:56,826 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2022-06-07 20:45:56,857 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2022-06-07 20:45:57,001 INFO mapred.FileInputFormat: Total input files to process : 1
2022-06-07 20:45:57,851 INFO mapreduce.JobSubmitter: number of splits:1
2022-06-07 20:45:58,128 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1343776106_0001
2022-06-07 20:45:58,129 INFO mapreduce.JobSubmitter: Executing with tokens: []
2022-06-07 20:45:58,518 INFO mapred.LocalDistributedCacheManager: Localized file:/Users/sruthysanthosh/Downloads/mapper1.py as file:/tmp/hadoop-sruthysanthosh/mapred/local/job_local1343776106_0001_96318fd-f1dd-464c-a99e-1cc55d6a65d6/mapper1.py
2022-06-07 20:45:58,557 INFO mapred.LocalDistributedCacheManager: Localized file:/Users/sruthysanthosh/Downloads/reducer1.py as file:/tmp/hadoop-sruthysanthosh/mapred/local/job_local1343776106_0001_be9a1578-5664-44c8-b061-989c96b2979a/reducer1.py
2022-06-07 20:45:58,776 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2022-06-07 20:45:58,781 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2022-06-07 20:45:58,783 INFO mapreduce.Job: Running job: job_local1343776106_0001
2022-06-07 20:45:58,786 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
2022-06-07 20:45:58,794 INFO mapred.FileOutputCommitter: File Output Committer Algorithm version is 2
2022-06-07 20:45:58,802 INFO mapred.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2022-06-07 20:45:58,802 INFO mapred.LocalJobRunner: Waiting for map tasks
2022-06-07 20:45:58,804 INFO mapred.LocalJobRunner: Starting task: attempt_local1343776106_0001_m_000000_0
2022-06-07 20:45:58,894 INFO mapred.FileOutputCommitter: File Output Committer Algorithm version is 2
2022-06-07 20:45:58,894 INFO mapred.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup failures: false
2022-06-07 20:45:58,904 INFO util.ProcsBasedProcessTree: ProcsBasedProcessTree currently is supported only on Linux.
2022-06-07 20:45:58,936 INFO mapred.MapTask: Using ResourceCalculatorProcessorTree : null
2022-06-07 20:45:58,936 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/user/sruthy/text.txt:9+665057
2022-06-07 20:45:58,994 INFO mapred.MapTask: numReduceTasks: 1
2022-06-07 20:45:59,022 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396/104857584)
2022-06-07 20:45:59,023 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2022-06-07 20:45:59,023 INFO mapred.MapTask: soft limit at 83886000
2022-06-07 20:45:59,023 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2022-06-07 20:45:59,023 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2022-06-07 20:45:59,027 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2022-06-07 20:45:59,044 INFO streaming.PipeMapper: PipeMapper exec (/Users/sruthysanthosh/..mapper1.py)
2022-06-07 20:45:59,108 INFO Configuration.deprecation: mapred.work.output.dir is deprecated. Instead, use mapreduce.task.output.dir
2022-06-07 20:45:59,108 INFO Configuration.deprecation: mapred.local.dir is deprecated. Instead, use mapreduce.cluster.local.dir
2022-06-07 20:45:59,166 INFO Configuration.deprecation: map.input.file is deprecated. Instead, use mapreduce.map.input.file
2022-06-07 20:45:59,166 INFO Configuration.deprecation: map.input.length is deprecated. Instead, use mapreduce.map.input.length
2022-06-07 20:45:59,166 INFO Configuration.deprecation: mapred.job.id is deprecated. Instead, use mapreduce.job.id
```

```
WRONG_MAP=0
WRONG_REDUCE=0
File Output Format Counters
Bytes Written=10024
022-06-10 11:16:26,093 INFO mapred.LocalJobRunner: Finishing task: attempt_local28801168
022-06-10 11:16:26,094 INFO mapred.LocalJobRunner: reduce task executor complete.
022-06-10 11:16:26,369 INFO mapreduce.Job: map 100% reduce 100%
022-06-10 11:16:26,370 INFO mapreduce.Job: Job job_local288011680_0001 completed successfully
022-06-10 11:16:26,392 INFO mapreduce.Job: Counters: 36
File System Counters
FILE: Number of bytes read=10659488
FILE: Number of bytes written=17300799
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=54129196
HDFS: Number of bytes written=10024
HDFS: Number of read operations=15
HDFS: Number of large read operations=0
HDFS: Number of write operations=4
HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
Map input records=450018
Map output records=450018
Map output bytes=4427247
Map output materialized bytes=5327289
Input split bytes=108
Combine input records=0
Combine output records=0
Reduce input groups=299
Reduce shuffle bytes=5327289
Reduce input records=450018
Reduce output records=298
Spilled Records=900036
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=107
Total committed heap usage (bytes)=1068498944
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=27864598
File Output Format Counters
Bytes Written=10024
022-06-10 11:16:26,392 INFO streaming.StreamJob: Output directory: /user/sruthy/output1
```

The output file in hadoop

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities ▾

Browse Directory

/user/sruthy/output Go!

Show 25 entries Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	-rw-r--r--	sruthysanthosh	supergroup	0 B	Jun 07 20:46	1	128 MB	_SUCCESS
<input type="checkbox"/>	-rw-r--r--	sruthysanthosh	supergroup	416.61 KB	Jun 07 20:46	1	128 MB	part-00000

Showing 1 to 2 of 2 entries Previous 1 Next

Hadoop, 2022.

Datanode Volume Failures Snapshot Startup Progress Utilities ▾

File information - part-00000

[Download](#) [Head the file \(first 32K\)](#) [Tail the file \(last 32K\)](#)

Block information -- Block 0

Block ID: 1073741853
Block Pool ID: BP-545430982-127.0.0.1-1654625867185
Generation Stamp: 1029
Size: 233636
Availability:

- localhost

File contents

```
"Access" 1
"Acid" 3
"Ad-hocracy" 1
"Advanced" 1
"Agents" 1
"AI" 3
"All" 3
"American" 1
```

- Exercise 2: Analysis of Airport efficiency with Map Reduce

The required csv file was downloaded containing 450017 lines. The file was transferred to the hadoop directory using the -put command. I have used one mapper and one reducer for implementing each of the tasks.

- 1) Computing the maximum, minimum, and average departure delay for each airport.

Mapper file:

```
import sys

#for each line in input
for line in sys.stdin:

    # Removing leading, trailing whitespace and quotes
    line = line.strip().replace('\\"', '')
    # getting each attribute by splitting using , as it is a csv file
    line = line.split(",")

    if len(line) >= 2:
        #Getting the required column values
        dep_airport = line[3]
        dep_delay = line[6]

        #If column value is empty, assigning 0
        if dep_delay == '':
            dep_delay = 0

        #Forming the output of mapper using tab character
        string = '%s\t%s' % (dep_airport, dep_delay)
        print(string)
```

In the mapper, we get each attribute after splitting. Then the required attributes which are the departure airport and departure delay are selected and send as output, separated by a tab character.

The sample output of mapper:

PHL	-10.00
CLT	-8.00
RDU	-7.00
CLT	-7.00
PHL	-3.00
FLL	-7.00
PHL	-5.00
RDU	-8.00
PHX	-4.00
ORD	-6.00
IND	-13.00
ALB	-14.00
CLT	-7.00
BWI	-6.00
CLT	28.00
MCO	-7.00
PHL	-8.00
CLT	-7.00
PIT	-9.00
CLT	-5.00
PHL	13.00
RSW	11.00
CLT	4.00
PBI	23.00
PHL	-10.00
ORD	-5.00
PHL	-1.00
PHX	-7.00
PHL	-8.00
PHX	-16.00

The different airports will be sorted accordingly and send as input to reducer.
Sample input of reducer:

XNA	79.00
XNA	8.00
XNA	80.00
XNA	80.00
XNA	83.00
XNA	83.00
XNA	840.00
XNA	85.00
XNA	88.00
XNA	89.00
XNA	9.00
XNA	92.00
XNA	93.00
YAK	-10.00
YAK	-11.00
YAK	-11.00
YAK	-12.00
YAK	-13.00
YAK	-14.00
YAK	-18.00
YAK	-18.00
YAK	-19.00
YAK	-19.00
YAK	-2.00
YAK	-20.00
YAK	-20.00
YAK	-21.00

Reducer File:

```

#for each line in the input of reducer
# ( mapper output after sorting based on key)
for line in sys.stdin:
    #Removing trailing white spaces
    line = line.strip()
    #Getting the key and values send by mapper by splitting
    #according to tab character
    dep_airport, dep_delay= line.split('\t')
    #print(line)
    #Converting string to float as we need the dep delay in float
    try:
        dep_delay = float(dep_delay)
    except ValueError:
        continue

    #Forming dictionary for each unique key, with values as list
    if dep_airport in dict:
        dict[dep_airport].append(float(dep_delay))
    else:
        dict[dep_airport] = []
        dict[dep_airport].append(float(dep_delay))

    #For each departure airport
for dep_airport in dict.keys():
    #each value in dict is a list of delays for that airport
    #average delay calc
    avg_delay = sum(dict[dep_airport])*1.0 / len(dict[dep_airport])
    #max delay
    max_delay = max(dict[dep_airport])
    #min delay
    min_delay = min(dict[dep_airport])

    #Forming output of reducer by using tab as separator
    #Format - Departure Airport, Average delay, Maximum delay, Minimum delay
    string = '%s\t%s\t%s\t%s' % (dep_airport, avg_delay,max_delay,min_delay)
    #printing to output
    print(string)

```

Each key and value received from mapper is appended to a dictionary. Then the calculations are done on the dictionary. Each key in dictionary corresponds to a departure airport and its value will be a list of all the delays of that airport. Average, Maximum and minimum delays are found for each airport as shown in code. The output of the Mapreduce is stored in the output file in Hadoop.

The screenshot shows the Hadoop File Explorer interface at localhost:9870/explorer.html#/user/sruthy. The top navigation bar includes links for Overview, Datanodes, Datanode Volume Failures, Snapshot, Startup Progress, and Utilities. The main content area displays a list of files in the directory /user/sruthy, showing columns for Permission, Owner, Group, Size, Last Modified, Replication, Block Size, and Name. The files listed are T_ONTIME_REPORTING.csv, output, output1, and text.txt.

Browse Directory

This screenshot shows the detailed file information for 'part-00000' within the directory /user/sruthy/output1. A modal window titled 'File information - part-00000' is open, displaying the following details:

- Download**, **Head the file (first 32K)**, **Tail the file (last 32K)**
- Block information**: Block 0
- Block ID**: 1073741834
- Block Pool ID**: BP-545430982-127.0.0.1-1654625867185
- Generation Stamp**: 1010
- Size**: 10024
- Availability**: localhost

The main browser view shows the directory listing again, with 'Showing 1 to 2 of 2 entries'.

Hadoop, 2022.

This screenshot shows the detailed file information for 'part-00000' within the directory /user/sruthy/output1. A modal window titled 'File information - part-00000' is open, displaying the following details:

- Download**, **Head the file (first 32K)**, **Tail the file (last 32K)**
- Block information**: Block 0
- Block ID**: 1073741834
- Block Pool ID**: BP-545430982-127.0.0.1-1654625867185
- Generation Stamp**: 1010
- Size**: 10024
- Availability**: localhost

The main browser view shows the directory listing again, with 'Showing 1 to 2 of 2 entries'.

2) Computing a ranking list that contains top 10 airports by their average Arrival delay.

Mapper File: It is similar to the mapper file for the above question. Here we take the arrival airport and arrival delay attributes and then pass them as the output. The format of the output is also same.

```

#for each line in input
for line in sys.stdin:

    # Removing leading, trailing whitespace and quotes
    line = line.strip().replace('"', '')
    # getting each attribute by splitting using , as it is a csv file
    line = line.split(",")

    if len(line) >= 2:
        #Getting the required column values
        arr_airport = line[4]
        arr_delay = line[8]

        #If column value is empty, assigning 0
        if arr_delay == '':
            arr_delay = 0

        #Forming the output of mapper using tab character
        string = '%s\t%s' % (arr_airport, arr_delay)
        print(string)

```

Reducer File: In Reducer, like the previous question, the keys and values are taken and stored as dictionary. Hence each key would be an arrival airport and its value would be a list of delays corresponding to that airport. Then the average delay for each key is calculated and this is stored in another dictionary. This dictionary is then sorted in descending order, to get the airports having highest average rating. The output is written to Hadoop directory

```
#for each line in the input of reducer
# ( mapper output after sorting based on key)
for line in sys.stdin:
    #Removing trailing white spaces
    line = line.strip()
    #Getting the key and values send by mapper by splitting
    #according to tab character
    arr_airport, arr_delay= line.split('\t')

    #Converting string to float as we need the arr delay in float
    try:
        arr_delay = float(arr_delay)
    except ValueError:
        continue

    #Forming dictionary for each unique key, with values as list
    if arr_airport in dict:
        dict[arr_airport].append(float(arr_delay))
    else:
        dict[arr_airport] = []
        dict[arr_airport].append(float(arr_delay))

#For each arrival airport
for arr_airport in dict.keys():
    #each value in dict is a list of delays for that airport
    #average delay calc
    avg_delay = sum(dict[arr_airport])*1.0 / len(dict[arr_airport])

    #adding the average delay of each airport to a new dictionary
    new_dict[arr_airport]= avg_delay
```

```
count = 10

#Sorting in descending order based on the values
sorted_keys = sorted(new_dict, key=new_dict.get, reverse=True)

#Getting the 10 airports having highest average delay
for i in sorted_keys:
    #Forming output
    string = '%s\t%s' %(i, new_dict[i])
    print(string)

    count = count-1
    if(count==0):
        break
```

The output: The top 10 airports are ELM, BPT, GGG, BMI, ABI, LAW, LWS ,GRB, CHA and ACT.

The screenshot shows a browser window with the URL localhost:9870/explorer.html#/user/sruthy/output2. The main content is a "Browse Directory" view for the path `/user/sruthy/output2`, showing two entries: `_SUCCESS` and `part-00000`. A modal dialog box titled "File information - part-00000" is open. It contains tabs for "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Block information" tab is selected, showing details for "Block 0": Block ID: 1073741835, Block Pool ID: BP-545430982-127.0.0.1-1654625867185, Generation Stamp: 1011, Size: 202, and Availability: localhost. The "File contents" tab shows the following data:

```

ELM 81.76923076923077
BPT 47.857142857142854
GGG 46.375
BMI 37.584070796460175
ABI 34.142857142857146
LAW 29.733333333333334
LWS 29.06
GRB 24.98032786885246

```

The output files of all questions are attached.

- Exercise 3: Analysis of Movie dataset using Map and Reduce

Two files were downloaded for this question: `movies.dat` and `ratings.dat`. These files were merged on the Movie ID column to get the file `merge_data5.csv`. These files were then transferred to the hadoop directory using the `-put` command. In this exercise, time taken by each mapper X reducer also needs to be found, hence I have added time also to the hadoop streaming command. The number of reducers and mappers can be set using `-D mapred.job.maps=x -D mapred.job.reduces=x`.

1. Find the movie title which has the maximum average rating?

Mapper file: Here I have used the merged file for finding the movie title having the maximum average rating. First each line is read and then

split accordingly. Then the required attributes are send as output of the mapper. Here I have collected the movie id, movie title and rating and then send them separated by a tab character. Since I used merge function on movield, the file would be already sorted based on movie Id.

```
import sys
#For each line in input
for line in sys.stdin:

    # Removing leading, trailing whitespace and quotes
    line = line.strip().replace('"', '')
    # getting each attribute by splitting using , as it is a csv file
    line = line.split(",")

    if len(line) >= 2:
        #Getting the required column values
        movie_id = line[0]
        movie_name = line[1]
        rating = line[4]

        #Forming the output of mapper using tab character
        string = '%s\t%s\t%s' % (movie_id, movie_name,rating)
        print(string)
```

Sample Output of mapper file:

1580	Men in Black (1997)	3.5
1580	Men in Black (1997)	5.0
1580	Men in Black (1997)	4.0
1580	Men in Black (1997)	3.0
1580	Men in Black (1997)	1.0
1580	Men in Black (1997)	0.5
1580	Men in Black (1997)	3.0
1580	Men in Black (1997)	3.5
1580	Men in Black (1997)	1.0
1580	Men in Black (1997)	2.0
1580	Men in Black (1997)	4.0
1580	Men in Black (1997)	3.0
1580	Men in Black (1997)	5.0
1580	Men in Black (1997)	4.0
1580	Men in Black (1997)	3.5
1580	Men in Black (1997)	4.0
1580	Men in Black (1997)	4.0
1580	Men in Black (1997)	5.0
1580	Men in Black (1997)	3.0
1580	Men in Black (1997)	4.5
1580	Men in Black (1997)	2.5
1580	Men in Black (1997)	4.5
1580	Men in Black (1997)	2.5
1580	Men in Black (1997)	4.5

Reducer File: In the reducer file, the values send as output by mapper are collected. A dictionary is created with key as movie name and value as a list of ratings given for that movie. Then the average rating for each key is found and appended to a new dictionary.

```

#for each line in the input of reducer
# ( mapper output after sorting based on key)
for line in sys.stdin:
    #Removing trailing white spaces
    line = line.strip()
    #Getting the key and values send by mapper by splitting
    #according to tab character
    movie_id, movie_name, rating= line.split('\t')

    #Converting string to float as we need the rating in float
    try:
        rating= float(rating)
    except ValueError:
        continue

    #Forming dictionary for each unique key, with values as list
    if movie_name in dict:
        dict[movie_name].append(float(rating))
    else:
        dict[movie_name] = []
        dict[movie_name].append(float(rating))

#For each movie name
for name in dict.keys():
    #each movie name is a key and its values contain a list of all
    # the ratings it got
    #Average rating calculation
    avg_rating = sum(dict[name])*1.0 / len(dict[name])

    #New dictionary with key as movie name and average rating as va
    new_dict[name]= avg_rating

```

Then the new dictionary is sorted in descending order and we see that the highest average rating is 5 and 5 movies have that. These are written separated by tab to the output file in Hadoop directory.

```

max = 5.0

#Sorting the new dictionary is descending order
sorted_keys = sorted(new_dict, key=new_dict.get, reverse=True)

#for each movie name, if average rating is equal to the maximum rating,
#printing the same
for i in sorted_keys:
    if(new_dict[i] == max):
        #Forming output of the movie names with highest average rating
        string = '%s\t%s' %(i, new_dict[i])
        print(string)

```

Terminal:

```

Merged Map Outputs=0
GC time elapsed (ms)=930
Total committed heap usage (bytes)=6759645184
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=700272367
File Output Format Counters
  Bytes Written=204
2022-06-10 15:28:04,248 INFO streaming.StreamJob: Output directory: /user/sruthy/movieoutput1
$ hadoop jar -input /user/sruthy/merge_data5.csv -output -mapper -reducer    75.68s user 8.07s system 139% cpu 1:00.18 total
[sudo] sruthysanthosh@Sruthys-MacBook-Air ~ %

```

Hadoop Interface:

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities ▾

Browse Directory

/user/sruthy Go!

Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	-rw-r--r--	sruthysanthosh	supergroup	25.81 MB	Jun 10 11:12	1	128 MB	T_ONTIME_REPORTING.csv	
<input type="checkbox"/>	-rw-r--r--	sruthysanthosh	supergroup	667.81 MB	Jun 10 15:21	1	128 MB	merge_data5.csv	
<input type="checkbox"/>	drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 15:25	0	0 B	movieoutput1	
<input type="checkbox"/>	drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 07 20:46	0	0 B	output	
<input type="checkbox"/>	drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 11:16	0	0 B	output1	
<input type="checkbox"/>	drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 11:22	0	0 B	output2	
<input type="checkbox"/>	-rw-r--r--	sruthysanthosh	supergroup	252.82 MB	Jun 10 15:21	1	128 MB	ratings.dat	
<input type="checkbox"/>	-rw-r--r--	sruthysanthosh	supergroup	649.47 KB	Jun 07 20:30	1	128 MB	text.txt	

Showing 1 to 8 of 8 entries Previous 1 Next

Output:

The screenshot shows the Apache Hadoop Web UI. In the top navigation bar, there are several tabs: Overview, Datanodes, Datanode Volume Failures, Snapshot, Startup Progress, and Utilities. Below the navigation bar, there is a sub-navigation bar with tabs: Hadoop, Overview, Datanodes, Datanode Volume Failures, Snapshot, Startup Progress, and Utilities. A modal window titled "File information - part-00000" is open. It contains three tabs: "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Block information" tab is selected, showing details such as Block ID: 1073741844, Block Pool ID: BP-545430982-127.0.0.1-1654625867185, Generation Stamp: 1020, Size: 204, and Availability: localhost. The "File contents" tab shows a list of movie titles and their average ratings:

Title	Average Rating
Satan's Tango (Sátántangó) (1994)	5.0
Shadows of Forgotten Ancestors (1964)	5.0
Fighting Elegy (Kenka erejii) (1966)	5.0
Sun Alley (Sonnenallee) (1999)	5.0
Blue Light. The (Das Blaue Licht) (1932)	5.0

The movie titles having the maximum average rating of 5 are Satan's Tango, Shadows of Forgotten Ancestors, Fighting Elegy, Sun Alley, Blue Light, The.

I have analysed the time for different number of mappers and reducers.

Number of Mappers	Number of Reducers	Time in seconds
1	1	75.68s
2	2	76.81s
4	4	76.29s

Number of Mappers	Number of Reducers	Time in seconds
2	1	74.55s

4	1	75.15s
8	1	77.53s

Number of Mappers	Number of Reducers	Time in seconds
1	2	76.94s
1	4	76.30s
1	8	74.85s

From the above tables we can see that, on increasing both mappers and reducers, the time taken increases slightly (<1s). This can be due to computation or overhead of having many mappers and reducers.

When increasing only mappers, we see that initially the time decreases and then increases after 4 mappers. Also for increase in only reducers time also decreased slightly. But overall the time changes is very very less when changing the number of mappers and reducers. This can be because for this operation, increasing mappers and reducers may not help much.

But still we can say that increasing reducers help in decreasing time taken. Since reducing function does the main operation, it makes sense to have more reducers.

2. Find the user who has assigned lowest average rating among all the users who rated more than 40 times?

Here I have only used the ratings.dat file as input.

Mapper File: As done in previous question, the mapper splits and gets the required attributes. Here we collect the user Id and the rating. These values are the output of the mapper file.

```
#For each line in input file
for line in sys.stdin:

    # Removing leading, trailing whitespace and quotes
    line = line.strip().replace('"', '')
    # getting each attribute by splitting using :: as it is a dat file
    line = line.split("::")

    if len(line) >=2 :
        #Getting the required column values
        user_id = line[0]
        rating = line[2]

        #Forming the output of mapper using tab character
        string = '%s\t%s' % (user_id, rating)
        print(string)
```

Then the output of mapper is sorted according to user ids and send as the input to Reducer. Hence all same userIds will be together.

Reducer File: In the reducer file, the inputs are split and then appended to a dictionary. The userIds will be keys and the values will be the list of ratings which were given by that user.

Then the users having more than 40 ratings in their value is selected and the average of the ratings for each such user is calculated and stored in a new dictionary. This dictionary is then sorted in ascending order. The first key in the sorted dictionary will contain the user who assigned the lowest average rating among users who have rated more than 40 times.

```

#for each line in the input of reducer
# ( mapper output after sorting based on key)
for line in sys.stdin:
    #Removing trailing white spaces
    line = line.strip()
    #Getting the key and values send by mapper by splitting
    #according to tab character
    user_id, rating= line.split('\t')

    #Converting string to float as we need the rating in float
try:
    rating= float(rating)
except ValueError:
    continue

#Forming dictionary for each unique key, with values as list
if user_id in dict:
    dict[user_id].append(float(rating))
else:
    dict[user_id] = []
    dict[user_id].append(float(rating))

#For each user in dictionary
#each userid is a key and its values contain a list of all
# the ratings given by that user
for user in dict.keys():
    #If the number of ratings given by user is greater than or equal to 40
    if(len(dict[user])>=40):
        #Average rating calculation
        avg_rating = sum(dict[user])*1.0 / len(dict[user])

        #New dictionary with userid as key and the average rating as value
        new_dict[user]= avg_rating

```

```

min = 1.0
#Sorting dictionary in ascending order
sorted_keys = sorted(new_dict, key=new_dict.get)

for i in sorted_keys:
    #Finding the user with the least average rating
    if(new_dict[i] == min):
        #Forming the output as :userid  average rating
        string = '%s\t%s' %(i, new_dict[i])
        print(string)

```

The output of the reducer would be the user Id and the average rating fulfilling the above conditions.

Hadoop Interface:

The screenshot shows the Hadoop interface at localhost:9870/explorer.html#/user/sruthy/movieoutput2_map1. A modal window titled "File information - part-00000" is open, displaying details about Block 0. The modal has tabs for "Download", "Head the file (first 32K)", and "Tail the file (last 32K)". The "Block information" tab is selected, showing the following details:

- Block ID: 1073741847
- Block Pool ID: BP-545430982-127.0.0.1-1654625867185
- Generation Stamp: 1023
- Size: 10
- Availability: • localhost

Below the modal, the main interface shows a "Browse Directory" view for the path /user/sruthy/movieoutput2_map1. It lists two entries, both owned by sruthysar and have permission -rw-r--r--. The interface also shows a footer with "Hadoop, 2022."

The output is userId 24178 having an average rating of 1.0

Number of Mappers	Number of Reducers	Time in seconds
1	1	65.23s
2	2	65.53s
4	4	66.31s

Number of Mappers	Number of Reducers	Time in seconds
1	2	64.39s
1	4	64.05s
1	8	65.44s

Number of Mappers	Number of Reducers	Time in seconds
2	1	64.44s
4	1	65.10
8	1	64.60

From the tables we can see that the trend is similar to that of above question. On increasing both mapper and reducer, the time increases very slightly. While increasing only reducers, the time decreases and then increases for 8 reducers. In the case where we increase only mappers, we see that time slightly increases and then decreases. The change is very very slight and can be due to CPU usage at the time.

Thus we can say the overall the time increases with increase in number of mappers and reducers as the computation overhead will be higher.

3. Find the highest average rated movie genre?

Here I used the merged csv file.

Mapper File: As done in previous question, the mapper splits and gets the required attributes. Here we collect the genre and the rating. These values are the output of the mapper file. The values are sorted according to genre before sending as input to reducer.

```
#For each line in input file
for line in sys.stdin:

    # Removing leading, trailing whitespace and quotes
    line = line.strip().replace('"', '')

    # getting each attribute by splitting using , as it is a csv file
    line = line.split(",")

    if len(line) >=2 :
        #Getting the required column values
        genre = line[2]
        rating = line[4]

        #Forming the output of mapper using tab character
        string = '%s\t%s' % (genre, rating)
        print(string)
```

Reducer File: In the reducer file, the inputs are split and then appended to a dictionary. The genres will be keys and the values will be the list of ratings which were given for each combination of genre.

Then the average of the ratings for each such genre list is calculated and stored in a new dictionary. This dictionary is then sorted in descending order. The first key in the sorted dictionary will contain the genre having the highest average rating. The output of reducer would be the genre and the average rating separated by a tab character. This will be written to the output directory mentioned in the hadoop command.

```

#for each line in the input of reducer
# ( mapper output after sorting based on key)
for line in sys.stdin:
    #Removing trailing white spaces
    line = line.strip()
    #Getting the key and values send by mapper by splitting
    #according to tab character
    genre, rating= line.split('\t')

    #Converting string to float as we need the rating in float
    try:
        rating= float(rating)
    except ValueError:
        continue

    #Forming dictionary for each unique key, with values as list
    if genre in dict:
        dict[genre].append(float(rating))
    else:
        dict[genre] = []
        dict[genre].append(float(rating))

#For each genre categories corresponding to movies in dictionary
#each genre list is a key and its values contain a list of all
# the ratings movies belonging to that genre list got
for genre in dict.keys():
    #Average rating calculation
    avg_rating = sum(dict[genre])*1.0 / len(dict[genre])

#New dictionary with genre list as key and the average rating as value
new_dict[genre]= avg_rating

```

```

count = 1

##Sorting dictionary in descending order to get the highest average rating
sorted_keys = sorted(new_dict, key=new_dict.get, reverse=True)

for i in sorted_keys:
    if(count>0):
        #The genre list having the highest average rated movies is found
        #Output is formed as : Genre list  Average rating
        string = '%s\t%s' %(i, new_dict[i])
        print(string)
        count = count - 1

```

Hadoop Interface:

localhost 9870/explorer.html#/user/sruthy

Info & Services for... IELTS Writing Task... IBM Garage Meth... M.Sc. Computer S... The interview: Ge... Certifying docum... myUNSW Portal... Welcome to the H... Academic Excelle...

Browse Directory

/user/sruthy

Show 25 entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	sruthysanthosh	supergroup	25.81 MB	Jun 10 11:12	1	128 MB	T_ONTIME_REPORTING.csv
-rw-r--r--	sruthysanthosh	supergroup	667.81 MB	Jun 10 15:21	1	128 MB	merge_data5.csv
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 15:25	0	0 B	movieoutput1
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 16:14	0	0 B	movieoutput2_map1
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 15:43	0	0 B	movieoutput_map2
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 15:46	0	0 B	movieoutput_map3
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 07 20:46	0	0 B	output
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 11:16	0	0 B	output1
drwxr-xr-x	sruthysanthosh	supergroup	0 B	Jun 10 11:22	0	0 B	output2
-rw-r--r--	sruthysanthosh	supergroup	252.82 MB	Jun 10 15:21	1	128 MB	ratings.dat
-rw-r--r--	sruthysanthosh	supergroup	649.47 KB	Jun 07 20:30	1	128 MB	text.txt

Showing 1 to 11 of 11 entries Previous 1 Next

localhost:9870/explorer.html#/user/sruthy/movieoutput3_new

IELTS Writing Task... IBM Garage Meth... M.Sc. Computer S... The interview: Ge... Certifying docum... myUNSW Portal...

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

File information - part-00000

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073741867
 Block Pool ID: BP-545430982-127.0.0.1-1654625867185
 Generation Stamp: 1043
 Size: 27
 Availability:
 • localhost

File contents

Animation|IMAX|Sci-Fi 4.75

Showing 1 to 2 of 2 entries

Hadoop, 2022.

The genre having the highest rating of 4.75 is Animation|IMAX|Sci-Fi

Sample Terminal:

```

Failed Shuffles=0
Merged Map outputs=6
GC time elapsed (ms)=451
Total committed heap usage (bytes)=6268387328
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=700272367
File Output Format Counters
Bytes Written=238
022-06-10 16:25:04,465 INFO streaming.StreamJob: Output directory: /user/sruthy/movieoutput3_map1
hadoop jar -D mapred.job.maps=1 -D mapred.job.reduces=1 -input -output      60.64s user 5.15s system 138% cpu 47.458 total
base) sruthysanthosh@Sruthys-MacBook-Air ~ %

```

Number of Mappers	Number of Reducers	Time in seconds
1	1	60.64s
2	2	62.52s
4	4	60.68s

Number of Mappers	Number of Reducers	Time in seconds
1	2	61.11s
1	4	60.21s
1	8	62.51

Number of Mappers	Number of Reducers	Time in seconds
2	1	61.62s
4	1	60.66s
8	1	60.62s

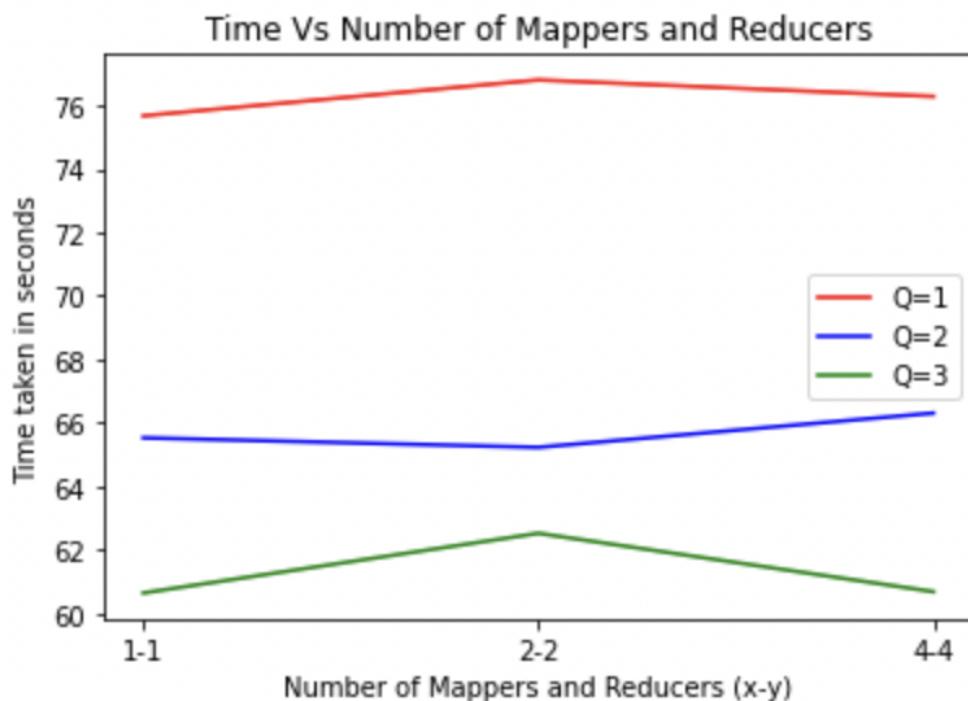
From the tables we can see that here, when we increase both mappers and reducers, the time increases and then decreases. This can be due to randomness and

the CPU usage. While increasing only reducers , we see a slight decrease initially, then increases. When we increase only mappers, we see that there is a slight decrease in time.

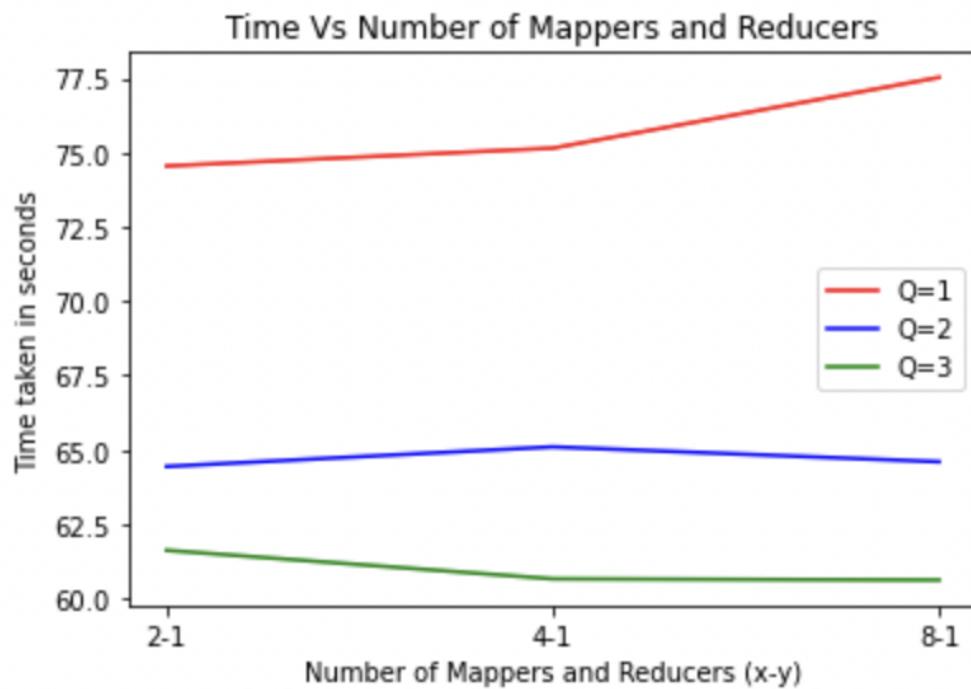
Thus overall we can say that increasing only mappers or reducers show a slightly increase in trend initially. But this is still very very slight and can be considered as randomness also (not sure).

I have written the code for merging of datasets and the code for plotting graphs in Jupyter notebook which is also attached.

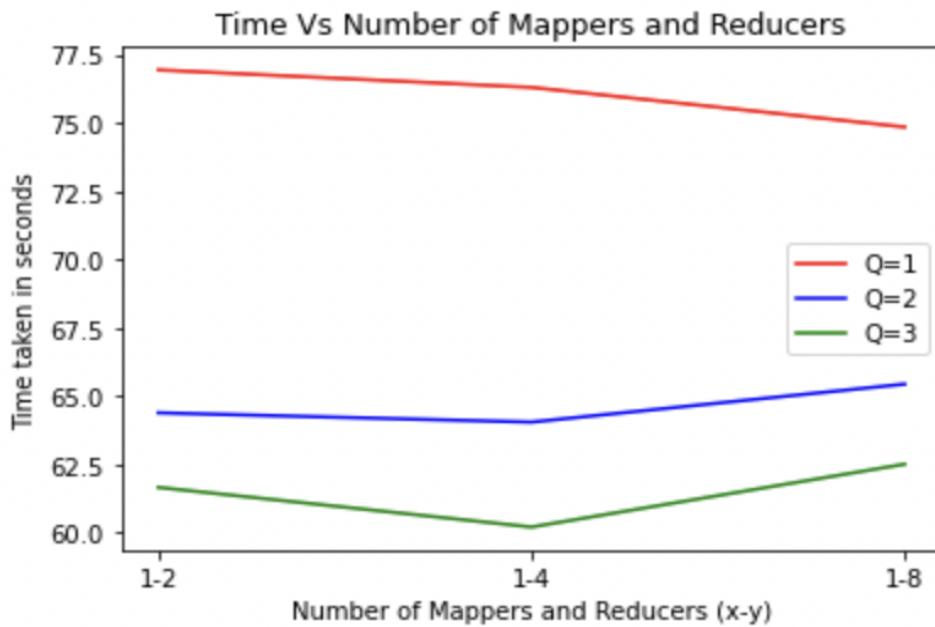
Graphs



We see that for Question 1 and 3, time decreases after initially increasing



We see that time varies very slightly.



We can see that the time decreases initially and then increases for Question 2 and 3.

Also the time taken by Question 1 > Question 2 > Question 3

References:

- <https://stackoverflow.com/questions/20577840/python-dictionary-sorting-in-descending-order-based-on-values>
- <https://stackoverflow.com/questions/22346086/measure-the-shell-script-execution-time-in-milliseconds-on-mac-os>
- <https://stackoverflow.com/questions/37697195/how-to-merge-two-data-frames-based-on-particular-column-in-pandas-python>
- <https://stackoverflow.com/questions/16923281/writing-a-pandas-dataframe-to-csv-file>
- <https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.replace.html>
- <https://stackoverflow.com/questions/6885441/setting-the-number-of-map-tasks-and-reduce-tasks>
- <https://codewithharjun.medium.com/>
- <https://blog.contactsunny.com/data-science/installing-hadoop-on-the-new-m1-pro-and-m1-max-macbook-pro>
- <https://www.geeksforgeeks.org/hadoop-streaming-using-python-word-count-problem/>
-