

**Followed instructions the below link to install hadoop-2.8.0**

**<https://github.com/MuhammadBilalYar/Hadoop-On-Window/wiki/Step-by-step-Hadoop-2.8.0-installation-on-Window-10>**

### **Exercise 1.2: Basic hadoop operations**

1.

```
C:\hadoop-2.8.0\sbin>hadoop version
Hadoop 2.8.0
Subversion https://git-wip-us.apache.org/repos/asf/hadoop.git -r 91f2b7a13d1e97be65db92ddabc627cc29ac0009
Compiled by jdu on 2017-03-17T04:12Z
Compiled with protoc 2.5.0
From source with checksum 60125541c2b3e266cbf3becc5bda666
This command was run using /C:/hadoop-2.8.0/share/hadoop/common/hadoop-common-2.8.0.jar
```

2.

```
C:\hadoop-2.8.0\sbin>hadoop fs -ls /
Found 1 items
drwxr-xr-x  - manna supergroup          0 2019-07-03 14:48 /hadoopdemo
```

3.

```
C:\hadoop-2.8.0\sbin>hadoop fs -mkdir /hadoopdemo
```

**Hadoop** Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities -

### Browse Directory

Show 25 entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
drwxr-xr-x	manna	supergroup	0 B	Jul 03 14:48	0	0 B	hadoopdemo	

Showing 1 to 1 of 1 entries Previous **1** Next

4.

```
C:\hadoop-2.8.0\sbin> hadoop fs -mkdir /hadoopdemo/text_files
C:\hadoop-2.8.0\sbin> hadoop fs -mkdir /hadoopdemo/raw_data
```

**Hadoop** Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities ▾

## Browse Directory

Show  ▾ entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<a href="#">drwxr-xr-x</a>	<a href="#">manna</a>	<a href="#">supergroup</a>	0 B	Jul 03 14:48	<a href="#">0</a>	0 B	<a href="#">raw_data</a>
<a href="#">drwxr-xr-x</a>	<a href="#">manna</a>	<a href="#">supergroup</a>	0 B	Jul 03 14:48	<a href="#">0</a>	0 B	<a href="#">text_files</a>

Showing 1 to 2 of 2 entries Previous **1** Next

Hadoop, 2017.

5.

```
C:\hadoop-2.8.0\sbin>hadoop fs -put C:\Users\manna\desktop\file.txt /hadoopdemo/text_files
```

**Hadoop** Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities ▾

## Browse Directory

Show  ▾ entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<a href="#">-rw-r--r--</a>	<a href="#">manna</a>	<a href="#">supergroup</a>	778 B	Jul 03 14:57	<a href="#">1</a>	128 MB	<a href="#">file.txt</a>

Showing 1 to 1 of 1 entries Previous **1** Next

Hadoop, 2017.

7.

```
C:\hadoop-2.8.0\sbin>hadoop fs -rm -r /hadoopdemo/text_files
Deleted /hadoopdemo/text_files
```

[Hadoop](#) [Overview](#) [Datanodes](#) [Datanode Volume Failures](#) [Snapshot](#) [Startup Progress](#) [Utilities](#)

## Browse Directory

Show  entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<a href="#">drwxr-xr-x</a>	<a href="#">manna</a>	<a href="#">supergroup</a>	0 B	Jul 03 14:48	<a href="#">0</a>	0 B	<a href="#">raw_data</a>

Showing 1 to 1 of 1 entries Previous **1** Next

Hadoop, 2017.

8 & 9.

BEFORE

```
C:\hadoop-2.8.0\sbin>hadoop fs -cat /hadoopdemo/text_files/file.txt
https://www.google.com/maps/dir/Bromberger+Str.+44,+31141+Hildesheim,+Germany/PSV+Green+White+Hildesheim+e.V.,+Marienburger+Stra%C3%9Fe,+Hildesheim/@52.1340656,9.9692939,16z/am=t/data=!4m2!3m1!1s0x47baae3f6062ee71:0x65d9561fd127fad0!2m2!1d9.9785292!2d52.129663!2m3!6e0!7e2!8j1561021320!3e3
9,50 euro

https://www.google.com/maps/dir/Bromberger+Str.+44,+31141+Hildesheim,+Germany/DJK+Blau-Wei%C3%9F+Hildesheim,+Lucienv%C3%86rder+Allee+8,+31139+Hildesheim/@52.1434274,9.9471582,15z/data=!4m2!1d9.9737509!2d52.138727!1m5!1m1!1s0x47baae3f6062ee71:0x4be7444a9d97324!2m2!1d9.9413713!2d52.1415661!2m3!6e0!7e2!8j1561021320!3e3
8,00 euro
```

AFTER

```
C:\hadoop-2.8.0\sbin>hadoop fs -cat /hadoopdemo/text_files/file.txt
https://www.google.com/maps/dir/Bromberger+Str.+44,+31141+Hildesheim,+Germany/PSV+Green+White+Hildesheim+e.V.,+Marienburger+Stra%C3%9Fe,+Hildesheim/@52.1340656,9.9692939,16z/am=t/data=!4m18!4m17!1m5!1m1!1s0x47baae3f6062ee71:0x65d9561fd127fad0!2m2!1d9.9785292!2d52.129663!2m3!6e0!7e2!8j1561021320!3e3
9,50 euro

https://www.google.com/maps/dir/Bromberger+Str.+44,+31141+Hildesheim,+Germany/DJK+Blau-Wei%C3%9F+Hildesheim,+Lucienv%C3%86rder+Allee+8,+31139+Hildesheim/@52.1434274,9.9471582,15z/data=!4m18!4m17!1m5!1m1!1s0x47baae3f6062ee71:0x4be7444a9d97324!2m2!1d9.9413713!2d52.1415661!2m3!6e0!7e2!8j1561021320!3e3
8,00 euro

hello world!!
C:\hadoop-2.8.0\sbin>
```

## Exercise 1.3: Word Count Map Reduce example

1. Created directory to store the text file in “input” directory after word count get the output in “output folder”.

```
C:\hadoop-2.8.0>hadoop fs -mkdir /input  
C:\hadoop-2.8.0>hadoop fs -mkdir /output
```

2. Transfer and store the text file “big.txt” from local system to hadoop system.

```
C:\hadoop-2.8.0>hadoop fs -put C:\Users\manna\Desktop\big.txt /input
```

3. Run the word count MapReduce job provided in

```
%HADOOP_HOME%>bin\yarn jar share\hadoop\mapreduce\hadoop-mapreduce-examples-  
2.8.0.jar /input output
```

Show 20 ▾ entries															Search: <input type="text"/>		
ID ▾	User ▾	Name ▾	Application Type ▾	Queue ▾	Application Priority ▾	StartTime ▾	FinishTime ▾	State ▾	FinalStatus ▾	Running Containers ▾	Allocated CPU Vcores ▾	Allocated Memory MB ▾	% of Queue ▾	% of Cluster ▾	Progress ▾	Tracking UI ▾	Blacklisted Nodes ▾
application_1562185904484_0001	manna	word count	MAPREDUCE	default	0	Wed Jul 3 22:45:24 +0200 2019	Wed Jul 3 22:45:52 +0200 2019	FINISHED	SUCCEEDED	N/A	N/A	N/A	0.0	0.0	<div></div>	<a href="#">History</a>	0
Showing 1 to 1 of 1 entries																	
															First Previous 1 Next Last		

First few output are given

```
stories! 1  
stories!" 1  
stories) 1  
stories, 6  
stories," 1  
stories. 2  
storing 2  
storm 23  
storm, 2  
storm-tossed 1  
storm. 5  
storm; 1  
storm?" 1  
stormcloud 2  
stormed 8  
stormed! 1  
storming 1  
storms 3  
storms. 1  
stormy 7  
stormy, 1  
story 83  
story! 2  
story), 1  
story, 13  
story--they 1  
story-teller. 1  
story. 21  
story." 2  
story: 1  
story?" 1  
stout 44  
stout, 15  
stout-built, 1  
stouter 3  
stouter!" 1  
stouter!..." 1  
stouter?" 1  
stoutest 2  
stoutly 3  
stoutness, 2  
stove 3  
stove. 3
```

## Exercise 2: Analysis of Airport efficiency with Map Reduce (10 points)

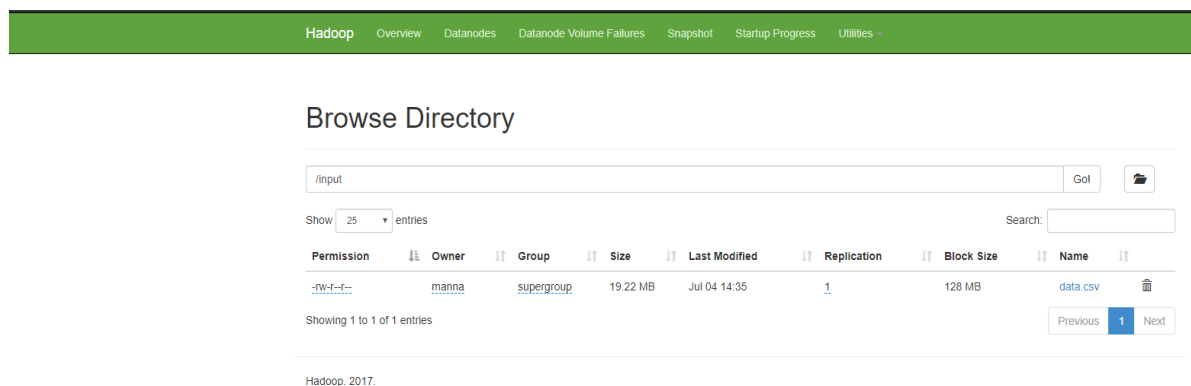
The following are the steps which are being followed while doing map reduce on Airport efficiency.

1. Two directories were created in hdfs i.e input(to store input data) and output(to get output data)

```
C:\hadoop-2.8.0\sbin>hadoop fs -mkdir /input  
C:\hadoop-2.8.0\sbin>hadoop fs -mkdir /output
```


2. Uploaded the .csv file in hdfs input folder

```
C:\hadoop-2.8.0\bin>hadoop fs -put C:/Users/manna/desktop/data.csv /input
```



Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

### Browse Directory

/input  

Show 25 entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	manna	supergroup	19.22 MB	Jul 04 14:35	1	128 MB	data.csv

Showing 1 to 1 of 1 entries

Hadoop, 2017.

3. Write Mapper and Reducer in python

In Hadoop, there is a Java program called Hadoop streaming-jar. This program internally read (stdin) and print out (stdout) line by line. Therefore, Python can read each line as a string and parse it by using functions like strip and split (" , "). For example, the first line would be parsed like this.

```
[2017-01-01, 20409, 264, SEA, JFK, 2102, 2.00, 0506, 1.00] -> ['2017-01-01','20409','264','SEA','JFK','2102','2.00','0506','1.00']
```

So, need to only pick the fourth element (airport) and the ninth element (arrival delay) in each array, and print out (stdout) them as Key-Value pair.

Here is Mapper.py code

```
1#!/usr/bin/python
2
3import sys
4
5# input comes from STDIN (standard input)
6for line in sys.stdin:
7    line = line.strip()
8    line = line.split(",")
9
10    if len(line) >=8:
11        airport = line[3]
12        arrdelay = line[8]
13
14    print ('%s\t%s' % (airport, arrdelay))
```

In reducer, the string is parsed which is coming from mapper.py as Key-Value pair, and keep them in dictionary like {airport:[arrival delay,.....]}

```
1#!/usr/bin/python
2
3#Reducer.py
4import sys
5
6
7(last_airport,max_arrdelay)=(None,0)
8(last_airport,min_arrdelay)=(None,0)
9airport_arrdelay = {}
10ave_delay_rank={}
11
12
13
14for line in sys.stdin:
15    data = line.strip().split('\t')
16
17    if len(data) != 2:
18        # Something has gone wrong. Skip this line.
19        continue
20    airport,arrdelay = data
21
22    if airport in airport_arrdelay:
23        airport_arrdelay[airport].append(int(arrdelay))
24    else:
25        airport_arrdelay[airport] = []
26        airport_arrdelay[airport].append(int(arrdelay))
27
28    if last_airport and last_airport!=airport:
29        print('maximum arrival delay:', '%s\t%s' % (last_airport, max_arrdelay))
30        print('minimum arrival delay:', '%s\t%s' % (last_airport, min_arrdelay))
31
32        (last_airport,max_arrdelay)=(airport,int(arrdelay))
33        (last_airport,min_arrdelay)=(airport,int(arrdelay))
34    else:
35        (last_airport,max_arrdelay)=(airport,max(max_arrdelay,int(arrdelay)))
36        (last_airport,min_arrdelay)=(airport,min(min_arrdelay,int(arrdelay)))
37
38
39for airport in airport_arrdelay.keys():
40    ave_delay = sum(airport_arrdelay[airport])*1.0 / len(airport_arrdelay[airport])
41    ave_delay_rank[airport] = []
42    ave_delay_rank[airport].append(float(ave_delay))
43    print ('average arrival delay', '%s\t%s' % (airport, ave_delay))
44    sorted_x = sorted(ave_delay_rank.items(), key=lambda kv: kv[1])
45    print("Top 10 airports by their average Arrival delay")
46    count=0
47    for keys,values in sorted_x:
48        if count<10:
49            print('%s\t%s' % (keys,values[0]))
50            count+=1
51
```

### 3.1 Explaining the reducer code

Collecting data line by line by strip and split function as well as performed a quick check if I have data of length less than 2. If the data is correct then the 'data' is unpacked and saved in 'airport' and 'arrdelay'

```
1 for line in sys.stdin:
2     data = line.strip().split('\t')
3
4     if len(data) != 2:
5         # Something has gone wrong. Skip this line.
6         continue
7     airport, arrdelay = data
```

This if else is performed to create a dict of airport and arrival delay time so that it can be used later on to find the average arrival delay time for each airport.

```
1 if airport in airport_arrdelay:
2     airport_arrdelay[airport].append(int(arrdelay))
3 else:
4     airport_arrdelay[airport] = []
5     airport_arrdelay[airport].append(int(arrdelay))
```

This where the maximum and minimum is time is found for each airport. First the data is sorted after coming out from the mapper. Then in every iteration the currentkey(airport) is compared with the oldkey(last\_airport) to find the minimum and maximum for each airport.

```
1 if last_airport and last_airport!=airport:
2     print('maximum arrival delay:', '%s\t%s'% (last_airport, max_arrdelay))
3     print('minimum arrival delay:', '%s\t%s'% (last_airport, min_arrdelay))
4
5     (last_airport, max_arrdelay) = (airport, int(arrdelay))
6     (last_airport, min_arrdelay) = (airport, int(arrdelay))
7 else:
8     (last_airport, max_arrdelay) = (airport, max(max_arrdelay, int(arrdelay)))
9     (last_airport, min_arrdelay) = (airport, min(min_arrdelay, int(arrdelay)))
```

At the end average arrival delay is calculated for every airport using the airport\_arrdelay dictionary which was created earlier. On the other hand, the ave\_delay\_rank dict holds the key (airport) and value (average arrival delay). Then it is sorted to find the top 10 airports by their average Arrival delay.

```
for airport in airport_arrdelay.keys():
    ave_delay = sum(airport_arrdelay[airport])*1.0 / len(airport_arrdelay[airport])
    ave_delay_rank[airport] = []
    ave_delay_rank[airport].append(float(ave_delay))
    print ('average arrival delay', '%s\t%s'% (airport, ave_delay))
sorted_x = sorted(ave_delay_rank.items(), key=lambda kv: kv[1])
print("Top 10 airports by their average Arrival delay")
count=0
for keys, values in sorted_x:
    if count<10:
        print('%s\t%s'% (keys, values[0]))
        count+=1
```

#### 4. Change access permission on Mapper and Reducer

Changed access permission on your mapper and reducer as below.

```
chmod +x mapper.py
chmod+x reducer.py
```

#### 5. Tested mapper.py and reducer.py locally

simulated process-flow by just using "cat" command in the terminal

```
manafee@Mannafee:/mnt/c/Users/manafee/desktop$ cat data.csv | python3 mapper.py | sort | python3 reducer.py
```



Showing output of minimum and maximum for the first few airports.

```
mannafee@mannafee:/mnt/c/Users/manna/desktop$ cat data.csv | python3 mapper.py | sort | python3 reducer.py
maximum arrival delay: ABE 786
minimum arrival delay: ABE -30
maximum arrival delay: ABI 261
minimum arrival delay: ABI -20
maximum arrival delay: ABQ 908
minimum arrival delay: ABQ -47
maximum arrival delay: ABR 1245
minimum arrival delay: ABR -29
maximum arrival delay: ABY 307
minimum arrival delay: ABY -42
maximum arrival delay: ACT 192
minimum arrival delay: ACT -27
maximum arrival delay: ACV 565
minimum arrival delay: ACV -26
maximum arrival delay: ACY 649
minimum arrival delay: ACY -43
maximum arrival delay: ADK 41
minimum arrival delay: ADK -31
maximum arrival delay: ADQ 71
minimum arrival delay: ADQ -31
maximum arrival delay: AEX 351
minimum arrival delay: AEX -16
```

Showing output of average arrival delay for the first few airports.

```
average arrival delay ABE 15.43010752688172
average arrival delay ABI 28.962962962962962
average arrival delay ABQ 5.056603773584905
average arrival delay ABR 35.033333333333333
average arrival delay ABY 3.654320987654321
average arrival delay ACT 11.142857142857142
average arrival delay ACV 10.952941176470588
average arrival delay ACY 2.139751552795031
average arrival delay ADK 2.3333333333333335
average arrival delay ADQ -5.225
average arrival delay AEX 9.93061224489796
average arrival delay AGS 16.647058823529413
average arrival delay ALB 2.7625994694960214
average arrival delay AMA 0.8148148148148148
average arrival delay ANC 4.973176865046102
average arrival delay APN 20.904761904761905
average arrival delay ASE 32.85406301824212
average arrival delay ATL 9.904612978889757
average arrival delay ATW 25.306532663316585
average arrival delay AUS 6.0779642663779105
average arrival delay AVL 8.51231527093596
average arrival delay AVP 28.821052631578947
average arrival delay AZO 4.115107913669065
average arrival delay BDL -3.306301050175029
```

Showing output for the top 10 airports by their average Arrival delay

```
Top 10 airports by their average Arrival delay
LSE -22.3
YAK -19.3125
PPG -12.5
CDV -11.26923076923077
ISN -9.619047619047619
LBE -8.7
EAU -8.018518518518519
ORH -7.649122807017544
INL -7.530612244897959
BFL -6.788571428571428
```

## 6. Testing mapper and reducer script in Hadoop platform

Using Hadoop streaming-2.8.0 jar to execute the mapper and reducer script

Command “**hadoop jar C:\hadoop-2.8.0\share\hadoop\tools\lib\hadoop-streaming-2.8.0.jar -file mapper.py -mapper "python mapper.py" -file reducer.py -reducer "python reducer.py" -input /input/data.csv -output 14**”

```
C:\hadoop-2.8.0>hadoop jar C:\hadoop-2.8.0\share\hadoop\tools\lib\hadoop-streaming-2.8.0.jar -file mapper.py -mapper "python mapper.py" -file reducer.py -reducer "python reducer.py" -input /input/data.csv -output 14
19/07/05 11:33:31 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [mapper.py, reducer.py] [] C:\Users\manna\AppData\Local\Temp\streamjob2848630007336944374.jar tmpDir=null
19/07/05 11:33:32 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
19/07/05 11:33:32 INFO Jvm2Metrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
19/07/05 11:33:32 INFO Jvm2Metrics: Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
19/07/05 11:33:33 INFO mapred.FileInputFormat: Total input files to process : 1
19/07/05 11:33:33 INFO mapreduce.JobSubmitter: number of splits:1
19/07/05 11:33:33 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local822758103_0001
19/07/05 11:33:34 INFO mapred.LocalDistributedCacheManager: Localized file: C:\hadoop-2.8.0\mapper.py as file: tmp/hadoop-manna/mapred/local/1562319214005/mapper.py
19/07/05 11:33:34 INFO mapred.LocalDistributedCacheManager: Localized file: C:\hadoop-2.8.0\reducer.py as file: tmp/hadoop-manna/mapred/local/1562319214006/reducer.py
19/07/05 11:33:34 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
19/07/05 11:33:34 INFO mapred.LocalJobRunner: OutputCommiter set in config null
19/07/05 11:33:34 INFO mapreduce.Job: Running job: job_local822758103_0001
19/07/05 11:33:34 INFO mapred.LocalJobRunner: OutputCommiter is org.apache.hadoop.mapred.FileOutputCommitter
19/07/05 11:33:34 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
19/07/05 11:33:34 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
19/07/05 11:33:34 INFO mapred.LocalJobRunner: Waiting for map tasks
19/07/05 11:33:34 INFO mapred.LocalJobRunner: Starting task: attempt_local822758103_0001_m_000000_0
19/07/05 11:33:34 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
19/07/05 11:33:34 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
19/07/05 11:33:34 INFO util.ProcfsBasedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux.
19/07/05 11:33:34 INFO mapred.Task: Using ResourceCalculatorProcessTree : org.apache.hadoop.yarn.util.WindowsBasedProcessTree@7bc2ab3
19/07/05 11:33:34 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/input/data.csv:0+20151734
19/07/05 11:33:34 INFO mapred.MapTask: numReduceTasks: 1
19/07/05 11:33:34 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
19/07/05 11:33:34 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
19/07/05 11:33:34 INFO mapred.MapTask: soft limit at 83886080
19/07/05 11:33:34 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
19/07/05 11:33:34 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
19/07/05 11:33:34 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
```

The map reduce operation was completed successfully

```
19/07/05 11:33:42 INFO mapred.LocalJobRunner: reduce task executor complete.
19/07/05 11:33:43 INFO mapreduce.Job: map 100% reduce 100%
19/07/05 11:33:43 INFO mapreduce.Job: Job job_local822758103_0001 completed successfully
19/07/05 11:33:43 INFO mapreduce.Job: Counters: 35
File System Counters
  FILE: Number of bytes read=8303786
  FILE: Number of bytes written=13110521
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=40303468
  HDFS: Number of bytes written=31489
  HDFS: Number of read operations=13
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=4
Map-Reduce Framework
  Map input records=450017
  Map output records=450017
  Map output bytes=3248855
  Map output materialized bytes=4148895
  Input split bytes=88
  Combine input records=0
  Combine output records=0
  Reduce input groups=298
  Reduce shuffle bytes=4148895
  Reduce input records=450017
  Reduce output records=900
  Spilled Records=900034
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=6
  Total committed heap usage (bytes)=662700032
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=20151734
File Output Format Counters
  Bytes Written=31489
19/07/05 11:33:43 INFO streaming.StreamJob: Output directory: 14
```

The output is stored in output directory named "14".  
Sample output is shown

```
C:\hadoop-2.8.0>hadoop fs -cat /user/manna/14/part-00000
maximum arrival delay: ABE      786
minimum arrival delay: ABE      -30
maximum arrival delay: ABI      261
minimum arrival delay: ABI      -20
maximum arrival delay: ABQ      908
minimum arrival delay: ABQ      -47
maximum arrival delay: ABR     1245
minimum arrival delay: ABR      -29
maximum arrival delay: ABY      307
minimum arrival delay: ABY      -42
maximum arrival delay: ACT      192
minimum arrival delay: ACT      -27
maximum arrival delay: ACV      565
minimum arrival delay: ACV      -26
maximum arrival delay: ACY      649
minimum arrival delay: ACY      -43
maximum arrival delay: ADK       41
minimum arrival delay: ADK      -31
maximum arrival delay: ADQ       71
minimum arrival delay: ADQ      -31
maximum arrival delay: AEX      351
minimum arrival delay: AEX      -40

average arrival delay ABE      15.43010752688172
average arrival delay ABI      28.962962962962962
average arrival delay ABQ       5.056603773584905
average arrival delay ABR      35.033333333333333
average arrival delay ABY       3.654320987654321
average arrival delay ACT      11.142857142857142
average arrival delay ACV      10.952941176470588
average arrival delay ACY       2.139751552795031
average arrival delay ADK       2.3333333333333335
average arrival delay ADQ       -5.225
average arrival delay AEX       9.93061224489796
average arrival delay AGS      16.647058823529413
average arrival delay ALB       2.7625994694960214
average arrival delay AMA       0.8148148148148148
average arrival delay ANC       4.973176865046102
average arrival delay APN      20.904761904761905
average arrival delay ASE      32.85406301824212
average arrival delay ATL       9.904612978889757
average arrival delay ATW      25.306532663316585
average arrival delay AUS       6.0779642663779105

Top 10 airports by their average Arrival delay
LSE      -22.3
YAK      -19.3125
PPG      -12.5
CDV      -11.26923076923077
ISN      -9.619047619047619
LBE      -8.7
EAU      -8.018518518518519
ORH      -7.649122807017544
INL      -7.530612244897959
BFL      -6.788571428571428
```

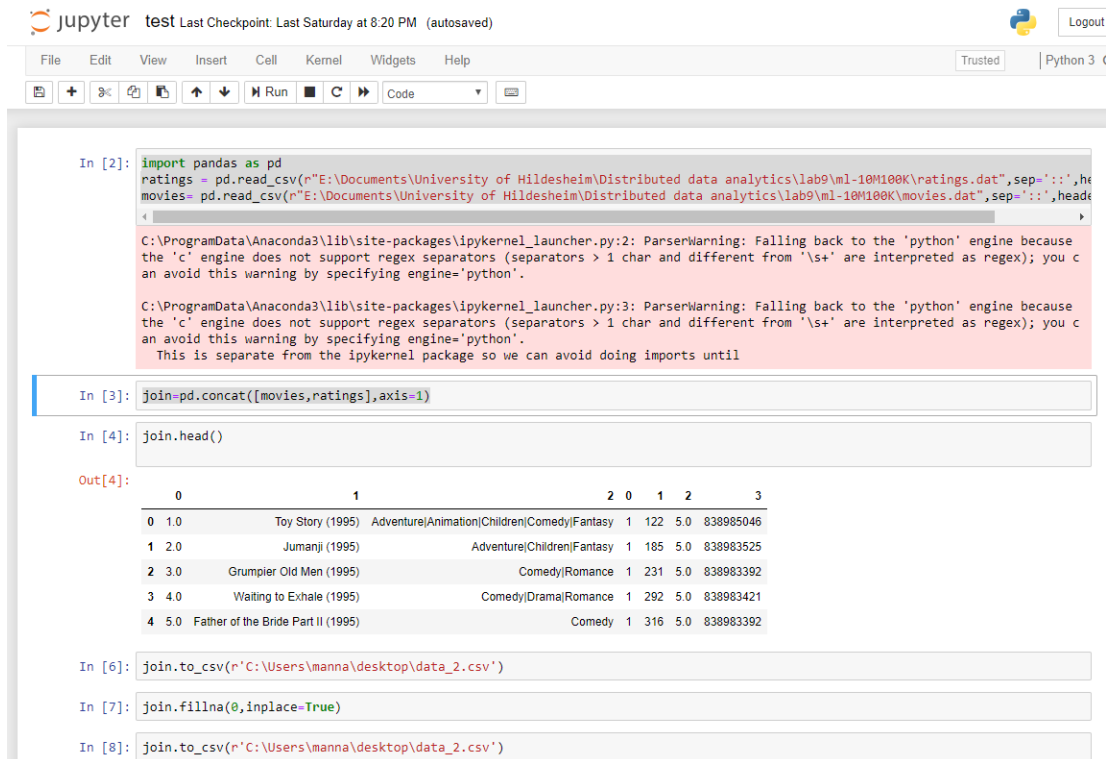
## Exercise 3: Analysis of Movie dataset using Map and Reduce

The following are the steps which are being followed while doing map reduce on Airport efficiency.

1. Two directories were created in hdfs i.e input(to store input data) and output(to get output data).

```
C:\hadoop-2.8.0\sbin>hadoop fs -mkdir /input
C:\hadoop-2.8.0\sbin>hadoop fs -mkdir /output
```

2.The two dat files(movies.dat and rating.dat) are merged to create one data.csv file



The screenshot shows a Jupyter Notebook interface with the following code and output:

```
In [2]: import pandas as pd
ratings = pd.read_csv(r"E:\Documents\University of Hildesheim\Distributed data analytics\lab9\ml-10M100K\ratings.dat", sep='::', header=0)
movies = pd.read_csv(r"E:\Documents\University of Hildesheim\Distributed data analytics\lab9\ml-10M100K\movies.dat", sep='::', header=0)
```

Warning messages are displayed for the CSV reads:

```
C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:2: ParserWarning: Falling back to the 'python' engine because the 'c' engine does not support regex separators (separators > 1 char and different from '\s+' are interpreted as regex); you can avoid this warning by specifying engine='python'.
C:\ProgramData\Anaconda3\lib\site-packages\ipykernel_launcher.py:3: ParserWarning: Falling back to the 'python' engine because the 'c' engine does not support regex separators (separators > 1 char and different from '\s+' are interpreted as regex); you can avoid this warning by specifying engine='python'.
This is separate from the ipykernel package so we can avoid doing imports until
```

```
In [3]: join = pd.concat([movies, ratings], axis=1)

In [4]: join.head()
```

Output [4] shows the first five rows of the merged data:

	0	1	2	0	1	2	3
0	1.0	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1	122	5.0	838985046
1	2.0	Jumanji (1995)	Adventure Children Fantasy	1	185	5.0	838983525
2	3.0	Grumpier Old Men (1995)	Comedy Romance	1	231	5.0	838983392
3	4.0	Waiting to Exhale (1995)	Comedy Drama Romance	1	292	5.0	838983421
4	5.0	Father of the Bride Part II (1995)	Comedy	1	316	5.0	838983392

```
In [6]: join.to_csv(r'C:\Users\manna\desktop\data_2.csv')

In [7]: join.fillna(0, inplace=True)

In [8]: join.to_csv(r'C:\Users\manna\desktop\data_2.csv')
```

### 3. Uploaded the .csv file in hdfs input folder

```
C:\hadoop-2.8.0\bin>hadoop fs -put C:/Users/manna/desktop/data.csv /input
```

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

#### Browse Directory

/input

Show 25 entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	man	supergroup	19.22 MB	Jul 04 14:35	1	128 MB	data.csv

Showing 1 to 1 of 1 entries

Hadoop, 2017.

### 4. Write Mapper and Reducer in python

mapper.py code

```
1 import sys
2
3 # input comes from STDIN (standard input)
4 movieid_name={}
5 for line in sys.stdin:
6     line = line.strip()
7     line = line.split(",")
8
9     movie_id_map=line[0]
10    movie_name=line[1]
11    genre=line[2]
12    user_id = line[3]
13    movie_id= line[4]
14    rating=line[5]
15
```

## Reducer.py

```
1
2
3 #Reducer.py
4 import sys
5 import operator
6
7
8
9 movieid_name={}
10 movieid_rating={}
11 userid_rating={}
12 genre_rating={}
13 ave_movie_rating={}
14 ave_user_rating={}
15 ave_genre_rating={}
16
17 for line in sys.stdin:
18     data = line.strip().split('\t')
19
20     if len(data) != 4:
21         # Something has gone wrong. Skip this line.
22         continue
23     movie_id_map,movie_name,genre,user_id,movie_id,rating= data
24
25     if movie_id_map in movieid_name:
26         movieid_name[movie_id_map].append(int(movie_name))
27     else:
28         movieid_name[movie_id_map] = []
29         movieid_name[movie_id_map].append(int(movie_name))
30
31     if movie_id in movieid_rating:
32         movieid_rating[movie_id].append(int(rating))
33     else:
34         movieid_rating[movie_id] = []
35         movieid_rating[movie_id].append(int(rating))
36
37     if user_id in userid_rating:
38         userid_rating[user_id].append(int(rating))
39     else:
40         userid_rating[user_id] = []
41         userid_rating[user_id].append(int(rating))
42
43     if genre in genre_rating:
44         genre_rating[genre].append(int(rating))
45     else:
46         genre_rating[genre] = []
47         genre_rating[genre].append(int(rating))
48
49
50 |
51 for movie_id in movieid_rating.keys():
52     ave_movie_rating[movieid_name[movie_id]].append(sum(movieid_rating[movie_id])*1.0 / len(movieid_rating[movie_id]))
53
54
55 for user_id in userid_rating.keys():
56     ave_user_rating[user_id].append(sum(userid_rating[user_id])*1.0 / len(userid_rating[user_id]))
57     if userid_rating[user_id]>40:
58         target_user=user_id
59
60 for genre in genre_rating.keys():
61     ave_genre_rating[genre_rating[genre]].append(sum(genre_rating[genre])*1.0 / len(genre_rating[genre]))
62
63
64 print(max(ave_movie_rating.iteritems(), key=operator.itemgetter(1))[0])
65 print(min(ave_user_rating[user_id].iteritems(), key=operator.itemgetter(1))[0])
66 print(max(ave_genre_rating[genre].iteritems(), key=operator.itemgetter(1))[0])
67
```

## 4.1 Explaining the reducer.py code

Collecting data line by line by strip and split function as well as performed a quick check if I have data of length is not equals to 6. If the data is correct then the 'data' is unpacked and saved in movie\_id\_map, movie\_name, genre, user\_id, movie\_id, rating respectively

```
for line in sys.stdin:
    data = line.strip().split('\t')

    if len(data) != 6:
        # Something has gone wrong. Skip this line.
        continue
    movie_id_map, movie_name, genre, user_id, movie_id, rating = data
```

This if else is used to store the movie name according to the movie\_id in a dictionary 'movieid\_name'

```
if movie_id_map in movieid_name:
    movieid_name[movie_id_map].append(int(movie_name))
else:
    movieid_name[movie_id_map] = []
    movieid_name[movie_id_map].append(int(movie_name))
```

This if else is used to store rating of movie according to their respective movie\_id in a dictionary 'movieid\_rating'

```
if movie_id in movieid_rating:
    movieid_rating[movie_id].append(int(rating))
else:
    movieid_rating[movie_id] = []
    movieid_rating[movie_id].append(int(rating))
```

This if else is used to store rating of movie according to the user\_id who rated in a dictionary 'userid\_rating'

```
if user_id in userid_rating:
    userid_rating[user_id].append(int(rating))
else:
    userid_rating[user_id] = []
    userid_rating[user_id].append(int(rating))
```

This if else is used to store rating according to genre in a dictionary 'genre\_rating'

```
if genre in genre_rating:
    genre_rating[genre].append(int(rating))
else:
    genre_rating[genre] = []
    genre_rating[genre].append(int(rating))
```

This for loop find the average rating for each movie id and in the print function the maximum value is collected then the key(movie\_name) is printed which has the maximum average rating.

```
for movie_id in movieid_rating.keys():
    ave_movie_rating[movieid_name[movie_id]].append(sum(movieid_rating[movie_id])*1.0 / len(movieid_rating[movie_id]))
print(max(ave_movie_rating.iteritems(), key=operator.itemgetter(1))[0])
```

This for loop finds the average rating according to each user but the user has to rate at least 40 times hence a condition check is done whether the length is greater than 40 or not. If it is greater the key and value is stored in ave\_user\_rating dictionary. At the end the user\_id is printed who has given the minimum rating for at least 40 times.

```
for user_id in userid_rating.keys():
    if(len(userid_rating[user_id]>40)):
        ave_user_rating[user_id].append(sum(userid_rating[user_id])*1.0 / len(userid_rating[user_id]))
print(min(ave_user_rating.iteritems(), key=operator.itemgetter(1))[0])
```

This loop is used to find the average rating for each genre and later the genre is printed which has the maximum rating.

```
for genre in genre_rating.keys():
    ave_genre_rating[genre_rating[genre]].append(sum(genre_rating[genre])*1.0 / len(genre_rating[genre]))
print(max(ave_genre_rating[genre].iteritems(), key=operator.itemgetter(1))[0])
```



## 6. Testing mapper and reducer script in Hadoop platform

Using Hadoop streaming-2.8.0 jar to execute the mapper and reducer script

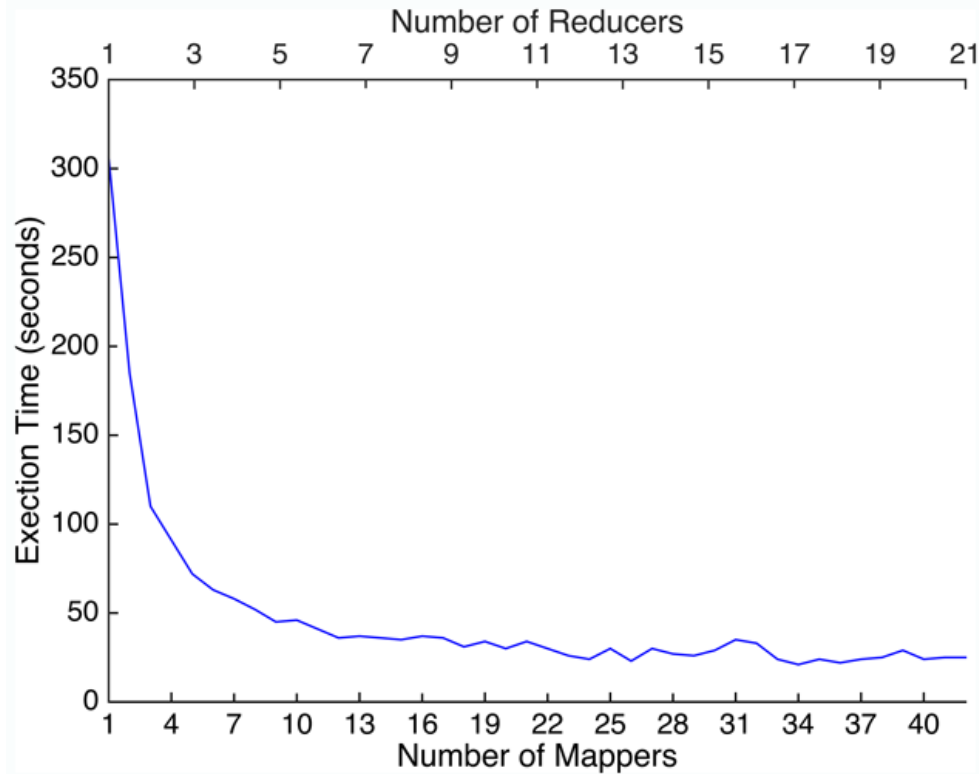
Command “**hadoop jar C:\hadoop-2.8.0\share\hadoop\tools\lib\hadoop-streaming-2.8.0.jar -file mapper.py -mapper "python mapper.py" -file reducer.py -reducer "python reducer.py" -input /input/data\_2.csv -output 14**”

```
C:\hadoop-2.8.0>hadoop jar C:\hadoop-2.8.0\share\hadoop\tools\lib\hadoop-streaming-2.8.0.jar -file mapper.py -mapper "python mapper.py" -file reducer.py -reducer "python reducer.py" -input /input/data.csv -output 14
19/07/05 11:33:31 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [mapper.py, reducer.py] [] C:\Users\manna\AppData\Local\Temp\streamjob2848638087336944374.jar tmpDir=null
19/07/05 11:33:32 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
19/07/05 11:33:32 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
19/07/05 11:33:32 INFO jvm.JvmMetrics: Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
19/07/05 11:33:33 INFO mapred.FileInputFormat: Total input files to process : 1
19/07/05 11:33:33 INFO mapreduce.JobSubmitter: number of splits:1
19/07/05 11:33:33 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local822758103_0001
19/07/05 11:33:34 INFO mapred.LocalDistributedCacheManager: Localized file:/C:/hadoop-2.8.0/mapper.py as file:/tmp/hadoop-manna/mapred/local/1562319214005/mapper.py
19/07/05 11:33:34 INFO mapred.LocalDistributedCacheManager: Localized file:/C:/hadoop-2.8.0/reducer.py as file:/tmp/hadoop-manna/mapred/local/1562319214006/reducer.py
19/07/05 11:33:34 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
19/07/05 11:33:34 INFO mapred.LocalJobRunner: OutputCommitter set in config null
19/07/05 11:33:34 INFO mapreduce.Job: Running job: job_local822758103_0001
19/07/05 11:33:34 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
19/07/05 11:33:34 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
19/07/05 11:33:34 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
19/07/05 11:33:34 INFO mapred.LocalJobRunner: Waiting for map tasks
19/07/05 11:33:34 INFO mapred.LocalJobRunner: Starting task: attempt_local822758103_0001_m_000000_0
19/07/05 11:33:34 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
19/07/05 11:33:34 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
19/07/05 11:33:34 INFO util.ProcfsBasedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux.
19/07/05 11:33:34 INFO mapred.Task: Using ResourceCalculatorProcessTree : org.apache.hadoop.yarn.util.WindowsBasedProcessTree@7bc2ab3
19/07/05 11:33:34 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/input/data.csv:0+20151734
19/07/05 11:33:34 INFO mapred.MapTask: numReduceTasks: 1
19/07/05 11:33:34 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
19/07/05 11:33:34 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
19/07/05 11:33:34 INFO mapred.MapTask: soft limit at 83886080
19/07/05 11:33:34 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
19/07/05 11:33:34 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
19/07/05 11:33:34 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
```

The map reduce operation was completed successfully

```
19/07/05 11:33:42 INFO mapred.LocalJobRunner: reduce task executor complete.
19/07/05 11:33:43 INFO mapreduce.Job: map 100% reduce 100%
19/07/05 11:33:43 INFO mapreduce.Job: Job job_local822758103_0001 completed successfully
19/07/05 11:33:43 INFO mapreduce.Job: Counters: 35
  File System Counters
    FILE: Number of bytes read=8303786
    FILE: Number of bytes written=13110521
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=40303468
    HDFS: Number of bytes written=31489
    HDFS: Number of read operations=13
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=4
  Map-Reduce Framework
    Map input records=450017
    Map output records=450017
    Map output bytes=3248855
    Map output materialized bytes=4148895
    Input split bytes=88
    Combine input records=0
    Combine output records=0
    Reduce input groups=298
    Reduce shuffle bytes=4148895
    Reduce input records=450017
    Reduce output records=900
    Spilled Records=900034
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=6
    Total committed heap usage (bytes)=662700032
  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=20151734
  File Output Format Counters
    Bytes Written=31489
19/07/05 11:33:43 INFO streaming.StreamJob: Output directory: 14
```

The graph of execution time vs the number of mappers × reducers for the overall task



The output is stored in output directory named "14".  
Sample output is shown.

```
Follow the Bitch(1998)
('3598',1.0153846153846153)
('Animation|Comedy|Thriller',4.473837209302325)
```