**System Info:**

| | |
|---|---|
| Processor | i7-5500U , 2.40GHz |
| Cores | 4 |
| Operating system | Windows 64 Bit |
| Ram | 8GB |
| Programming Language | Python 3.7.7 |

# Q1:

Making directory in HDFS and copied all the files I will use in exercise 1, 2, 3

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop fs -mkdir /lab6
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop fs -put E:\DDA\Exercise_6\files\* /lab6
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop fs -ls /lab6
Found 5 items
-rw-r--r--   3 fahad supergroup    1058582 2020-06-21 23:59 /lab6/1400-0.txt
-rw-r--r--   3 fahad supergroup     927445 2020-06-21 23:59 /lab6/158-0.txt
-rw-r--r--   3 fahad supergroup     234089 2020-06-21 23:59 /lab6/219-0.txt
-rw-r--r--   3 fahad supergroup     560162 2020-06-21 23:59 /lab6/2591-0.txt
-rw-r--r--   3 fahad supergroup    1586488 2020-06-21 23:59 /lab6/4300-0.txt
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib>
```

Running the prebuild World count program to count occurrence of words.

I have not attached full output because output was very lengthy instead I have taken a screenshot of running the command and end of output.

```
Windows PowerShell                                                                          —  □  ×
PS E:\hadoop-2.7.0\hadoop-2.7.0> bin\yarn jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\mapreduce\hadoop-mapreduce-examples-2.7.0.jar
ordcount /lab6/2591-0.txt WordCount
20/06/22 00:08:56 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
20/06/22 00:08:56 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
20/06/22 00:08:57 INFO input.FileInputFormat: Total input paths to process : 1
20/06/22 00:08:57 INFO mapreduce.JobSubmitter: number of splits:1
20/06/22 00:08:57 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1993384500_0001
20/06/22 00:08:57 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
20/06/22 00:08:57 INFO mapreduce.Job: Running job: job_local1993384500_0001
20/06/22 00:08:57 INFO mapred.LocalJobRunner: OutputCommitter set in config null
20/06/22 00:08:57 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
20/06/22 00:08:57 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
20/06/22 00:08:57 INFO mapred.LocalJobRunner: Waiting for map tasks
20/06/22 00:08:57 INFO mapred.LocalJobRunner: Starting task: attempt_local1993384500_0001_m_000000_0
20/06/22 00:08:57 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
20/06/22 00:08:57 INFO util.ProcfsBasedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux.
20/06/22 00:08:58 INFO mapred.Task:  Using ResourceCalculatorProcessTree : org.apache.hadoop.yarn.util.WindowsBasedProcessTree@260cac7

20/06/22 00:08:58 INFO mapred.MapTask: Processing split: hdfs://0.0.0.0:9000/lab6/2591-0.txt:0+560162
20/06/22 00:08:58 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
20/06/22 00:08:58 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
20/06/22 00:08:58 INFO mapred.MapTask: soft limit at 83886080
20/06/22 00:08:58 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
20/06/22 00:08:58 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
20/06/22 00:08:58 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
20/06/22 00:08:58 INFO input.LineRecordReader: Found UTF-8 BOM and skipped it
20/06/22 00:08:58 INFO mapred.LocalJobRunner:
20/06/22 00:08:58 INFO mapred.MapTask: Starting flush of map output
20/06/22 00:08:58 INFO mapred.MapTask: Spilling map output
20/06/22 00:08:58 INFO mapred.MapTask: bufstart = 0; bufend = 964786; bufvoid = 104857600
20/06/22 00:08:58 INFO mapred.MapTask: kvstart = 26214396(104857584); kvend = 25797812(103191248); length = 416585/6553600
20/06/22 00:08:58 INFO mapreduce.Job: Job job_local1993384500_0001 running in uber mode : false
20/06/22 00:08:58 INFO mapreduce.Job:  map 0% reduce 0%
20/06/22 00:08:58 INFO mapred.MapTask: Finished spill 0
20/06/22 00:08:58 INFO mapred.Task: Task:attempt_local1993384500_0001_m_000000_0 is done. And is in the process of committing
20/06/22 00:08:58 INFO mapred.LocalJobRunner: map
20/06/22 00:08:58 INFO mapred.Task: Task 'attempt_local1993384500_0001_m_000000_0' done.          Activate Windows
20/06/22 00:08:58 INFO mapred.LocalJobRunner: Finishing task: attempt_local1993384500_0001_m_000000_0   Go to Settings to activate Windows.
20/06/22 00:08:58 INFO mapred.LocalJobRunner: map task executor complete.
20/06/22 00:08:58 INFO mapred.LocalJobRunner: Waiting for reduce tasks
20/06/22 00:08:58 INFO mapred.LocalJobRunner: Starting task: attempt_local1993384500_0001_r_000000_0
```



```
Windows PowerShell
        File System Counters
                FILE: Number of bytes read=853146
                FILE: Number of bytes written=1564004
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=1120324
                HDFS: Number of bytes written=110408
                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=9571
                Map output records=104147
                Map output bytes=964786
                Map output materialized bytes=152958
                Input split bytes=100
                Combine input records=104147
                Combine output records=10949
                Reduce input groups=10949
                Reduce shuffle bytes=152958
                Reduce input records=10949
                Reduce output records=10949
                Spilled Records=21898
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=5
                Total committed heap usage (bytes)=544210944
        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=560162
        File Output Format Counters
                Bytes Written=110408
PS E:\hadoop-2.7.0\hadoop-2.7.0>
```

Printing output of prebuild word count program



```
PS E:\hadoop-2.7.0\hadoop-2.7.0> hadoop dfs -cat WordCount/*
```

```
Windows PowerShell
┌Ÿthere          3
┌Ÿthey  6
┌Ÿthis  4
┌Ÿthose          1
┌Ÿthou  2
┌Ÿthy   1
┌Ÿtis   1
┌Ÿto    5
┌Ÿtry   1
┌Ÿtwas  1
┌Ÿupon  1
┌Ÿwas   1
┌Ÿwe    6
┌Ÿwhat  27
┌ÿwhat┌ös         1
┌Ÿwhen  4
┌Ÿwho   3
┌Ÿwhose          1
┌Ÿwhy   6
┌Ÿwith  2
┌Ÿwould          2
┌Ÿyonder         1
┌Ÿyou   24
┌ÿyour┌öll        1
┌£Ah,   1
┌£Defects,┌¥      1
┌£Good  1
┌£Heads          1
┌£Here  2
┌£I     2
┌£Information    1
┌£Iron  1
┌£It    1
┌£Jip!┌¥┌ö        1
┌£Merrily        1
┌£Plain          2
┌£Project        5
┌£Right          1
┌£Under          1
┌£what  1
```

## Running my word count program and printing its output

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop fs -mkdir /lab6
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop fs -put E:\DDA\Exercise_6\files\* /lab6
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop fs -ls /lab6
Found 5 items
-rw-r--r--   3 fahad supergroup   1058582 2020-06-21 23:59 /lab6/1400-0.txt
-rw-r--r--   3 fahad supergroup    927445 2020-06-21 23:59 /lab6/158-0.txt
-rw-r--r--   3 fahad supergroup    234089 2020-06-21 23:59 /lab6/219-0.txt
-rw-r--r--   3 fahad supergroup    560162 2020-06-21 23:59 /lab6/2591-0.txt
-rw-r--r--   3 fahad supergroup   1586488 2020-06-21 23:59 /lab6/4300-0.txt
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-stream
-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_1\mapper.py -mapper "python mapper.py" -file E:\DDA\Exercise_6\Exercise_1\reducer.py -r
cer "python reducer.py" -input /lab6/2591-0.txt -output /lab6/output1.txt
20/06/22 00:01:19 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [E:\DDA\Exercise_6\Exercise_1\mapper.py, E:\DDA\Exercise_6\Exercise_1\reducer.py] [] C:\Users\fahad\AppData\Local\Tem
treamjob4716391358864535160.jar tmpDir=null
20/06/22 00:01:20 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
20/06/22 00:01:20 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
20/06/22 00:01:20 INFO jvm.JvmMetrics: Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
20/06/22 00:01:20 INFO mapred.FileInputFormat: Total input paths to process : 1
20/06/22 00:01:21 INFO mapreduce.JobSubmitter: number of splits:1
20/06/22 00:01:21 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1785863152_0001
20/06/22 00:01:22 INFO mapred.LocalDistributedCacheManager: Localized file:/E:/DDA/Exercise_6/Exercise_1/mapper.py as file:/tmp/hado
fahad/mapred/local/1592809281624/mapper.py
20/06/22 00:01:22 INFO mapred.LocalDistributedCacheManager: Localized file:/E:/DDA/Exercise_6/Exercise_1/reducer.py as file:/tmp/had
-fahad/mapred/local/1592809281625/reducer.py
20/06/22 00:01:22 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
20/06/22 00:01:22 INFO mapreduce.Job: Running job: job_local1785863152_0001
20/06/22 00:01:22 INFO mapred.LocalJobRunner: OutputCommitter set in config null
20/06/22 00:01:22 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
20/06/22 00:01:22 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
20/06/22 00:01:22 INFO mapred.LocalJobRunner: Waiting for map tasks
20/06/22 00:01:22 INFO mapred.LocalJobRunner: Starting task: attempt_local1785863152_0001_m_000000_0
20/06/22 00:01:22 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
20/06/22 00:01:22 INFO util.ProcfsBasedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux.
20/06/22 00:01:22 INFO mapred.Task:  Using ResourceCalculatorProcessTree : org.apache.hadoop.yarn.util.WindowsBasedProcessTree@1f1fc
```

```
20/06/22 00:01:27 INFO mapreduce.Job: Counters: 35
        File System Counters
                FILE: Number of bytes read=3181580
                FILE: Number of bytes written=5338674
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=1120324
                HDFS: Number of bytes written=147112
                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=9571
                Map output records=104147
                Map output bytes=1381494
                Map output materialized bytes=1589794
                Input split bytes=87
                Combine input records=0
                Combine output records=0
                Reduce input groups=10371
                Reduce shuffle bytes=1589794
                Reduce input records=104147
                Reduce output records=1
                Spilled Records=208294
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=0
                Total committed heap usage (bytes)=545259520
        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=560162
        File Output Format Counters
                Bytes Written=147112
```

```
20/06/22 00:01:27 INFO streaming.StreamJob: Output directory: /lab6/output1.txt
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop dfs -cat /lab6/output1.txt/*
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
{'#2591]': 2, '$5,000)': 2, ($1': 2, '(1785-1863)': 2, '(1786-1859),': 2, '(801)': 2, '(a)': 2, '(and': 4, '(any': 2, '(available': 2,
'(b)': 2, '(but': 2, '(c)': 2, '(does': 2, '(for': 3, '(if': 2, '(or': 4, '(rapunzel),': 2, '(trademark/copyright)': 2, '(who': 3, '(w
hose': 2, '(www.gutenberg.org),': 2, '(ΓÇfthe': 2, '***': 7, '*****': 4, '-': 8, '1.': 4, '1.a.': 2, '1.b.': 2, '1.c': 2, '1.c.': 2, '1
.d.': 2, '1.e': 2, '1.e.': 2, '1.e.1': 4, '1.e.1.': 3, '1.e.2.': 2, '1.e.3.': 2, '1.e.4.': 2, '1.e.5.': 2, '1.e.6.': 2, '1.e.7': 3, '1.
e.7.': 2, '1.e.8': 3, '1.e.8.': 3, '1.e.9.': 4, '1.f.': 2, '1.f.1.': 2, '1.f.2.': 2, '1.f.3.': 4, '1.f.3.': 2, '1.f.4.': 2, '1.f.5.': 2
, '1.f.6.': 2, '14,': 2, '1500': 2, '1812': 2, '1814.': 2, '1823,': 2, '2.': 4, '20%': 2, '2001': 2, '2001,': 2, '2008': 2, '2016': 2,
'2591-0.txt': 2, '2591-0.zip': 2, '3': 3, '3.': 3, '30': 2, '4': 2, '4.': 2, '4.': 2, '4557': 2, '5.': 2, '50': 2, '501(c)(3)': 3, '596
-1887.': 2, '60': 2, '64-6221541.': 3, '7': 2, '809': 2, '84116,': 2, '90': 3, '99712.': 2, '[*]': 2, [ebook': 2, '[little': 3, '_ju
g.': 2, '_my_': 2, 'a': 1961, 'a-fishing;': 2, 'a-hunting': 2, 'abashed,': 2, 'abide': 4, 'able': 27, 'abode': 3, 'about': 143, 'about,
': 22, 'about.': 9, 'about;': 3, 'about;': 3, 'about?': 3, 'about?ΓÇÖ': 3, 'above': 11, 'above,': 5, 'absence': 2, 'abundance.ΓÇÖ':
2, 'accept': 2, 'accepted': 2, 'accepted,': 2, 'accepting': 2, 'access': 11, 'accessed,': 2, 'accessible': 2, 'accident': 2, 'accidenta
lly': 2, 'accomplish.ΓÇÖ': 2, 'accomplished': 2, 'accord,': 2, 'accordance': 3, 'according': 5, 'accordingly': 2, 'accordingly,': 2, 'a
ccount!ΓÇÖ': 2, 'account': 6, 'account.': 2, 'accursed': 2, 'accused': 4, 'acknowledge': 2, 'acquaintance': 3, 'across': 13, 'across,':
 4, 'act': 2, 'active': 3, 'actual,': 2, 'actually': 3, 'actually,': 2, 'add': 3, 'add:': 2, 'added': 2, 'added:': 2, 'addition': 2, 'a
dditional': 5, 'additions': 2, 'address': 3, 'addresses.': 2, 'admiring': 2, 'admitted': 2, 'ado.': 2, 'adopt': 2, 'adrift': 2, 'adrift
,': 3, 'adrift.': 2, 'advanced': 2, 'advantage.': 2, 'adventure': 2, 'adventures': 3, 'advice': 3, 'advice,': 3, 'advice.': 5, 'advice
.ΓÇÖ': 2, 'advice:': 2, 'advice;': 2, 'afar': 7, 'affair.ΓÇÖ': 2, 'affairs,': 2, 'affected': 2, 'afraid': 14, 'afraid,': 6, 'afraid.':
2, 'afraid;': 4, 'after': 139, 'after,': 5, 'afternoon': 2, 'afterwards': 21, 'afterwards,': 8, 'again!': 7, 'again!ΓÇ\udc9d': 2, 'agai
n!ΓÇÖ': 2, 'again': 118, 'again,': 66, 'again,ΓÇÖ': 5, 'again.': 43, 'again.ΓÇÖ': 27, 'again:': 7, 'again;': 14, 'against': 27, 'againΓ
ÇÖ': 2, 'age.': 2, 'aged': 4, 'agent': 2, 'ago': 2, 'ago,': 2, 'ago.ΓÇÖ': 2, 'agree': 13, 'agreed': 25, 'agreed,': 3, 'a
greed.': 2, 'agreement': 10, 'agreement,': 7, 'agreement.': 4, 'ah!': 3, 'ah,': 3, 'ahead': 2, 'aid': 4, 'aid,': 2, 'ailed': 2, 'ails':
3, 'aim,': 2, 'aimed': 3, 'air!ΓÇÖ': 2, 'air': 8, 'air,': 10, 'air.': 7, 'air;': 3, 'ak,': 2, 'alarm': 2, 'alarmed,': 5, 'alas!': 13,
'alas!ΓÇÖ': 2, 'ale': 9, 'ale,': 4, 'ale.ΓÇÖ': 3, 'ale;': 2, 'alight': 2, 'alighted': 5, 'alighted;': 2, 'alike,': 3, 'alike:': 2, 'ali
ve!': 2, 'alive!ΓÇÖ': 2, 'alive': 8, 'alive,': 7, 'alive,ΓÇÖ': 2, 'alive.': 2, 'alive;': 2, 'alive?ΓÇÖ': 2, 'all!': 2, 'all!ΓÇÖ': 2, 'a
ll': 538, 'all,': 21, 'all,ΓÇÖ': 7, 'all--how': 2, 'all--skin,': 2, 'all-gone,': 2, 'all-gone;': 2, 'all.': 7, 'all.ΓÇÖ': 2, 'all;':
 4, 'all;': 13, 'all?ΓÇÖ': 3, 'allow': 6, 'allowed': 9, 'allowed,': 2, 'allowed;': 2, 'almost': 11, 'alms': 2, 'alone': 21, 'alone,':
6, 'alone,ΓÇÖ': 3, 'alone.': 9, 'alone;': 2, 'along': 25, 'along,': 8, 'along.': 2, 'along.ΓÇÖ': 2, 'aloud': 3, 'aloud,': 2, 'aloud':
3, 'already!ΓÇÖ': 3, 'already': 29, 'already,': 3, 'already.ΓÇÖ': 5, 'also': 30, 'also,': 9, 'also--that': 2, 'also.
': 3, 'also.ΓÇÖ': 2, 'also;': 2, 'altar': 2, 'alter': 2, 'alteration,': 2, 'alternate': 3, 'although': 10, 'altogether': 4, 'always': 3
8, 'always,': 2, 'always;': 2, 'am!ΓÇÖ': 11, 'am': 140, 'am,': 2, 'am,ΓÇÖ': 2, 'amazed': 2, 'amazed.': 2, 'ambassador': 2, 'ambush,': 2
'amends': 2, 'amicably': 2, 'amid': 3, 'amidst': 2, 'amiss,': 2, 'among': 21, 'amongst': 8, 'amuse': 3, 'amused': 2, 'amusement': 3,
'an': 138, 'and': 5489, 'and,': 15, 'angel': 3, 'angels': 4, 'anger': 4, 'anger,': 3, 'angered': 2, 'angrily': 4, 'angry': 9, 'angry,':
13, 'angry;': 2, 'animal': 3, 'animal,': 2, 'animals': 5, 'animals,': 3, 'animals.': 3, 'announced': 7, 'announced,': 2, 'announced.'
2, 'anointed': 2, 'another!ΓÇÖ': 3, 'another': 56, 'another,': 12, 'another.': 4, 'another;': 2, 'answer': 4, 'answer,': 3, 'answer,ΓÇ
Ö': 2, 'answer.': 2, 'answer;': 3, 'answered': 69, 'answered,': 16, 'answered.': 5, 'answered:': 27, 'answered;': 2, 'answers': 2, 'an
t-hill.': 2, 'ant-king': 4, 'ante-chamber': 2, 'ante-chamber,': 3, 'ante-chamber,ΓÇÖ': 2, 'ante-chamber.': 2, 'ants': 3, 'ants,': 2, 'a
```

```
: 2, ┌ćÿcoxcomb.┌ćö: 2, ┌ćÿcreep: 2, ┌ćÿcut: 2, ┌ćÿdear: 17, ┌ćÿdearest: 2, ┌ćÿdescend: 2, ┌ćÿdid: 2, ┌ćÿdo: 22, ┌ćÿd
one!┌ćö: 2, ┌ćÿdon┌ćöt: 6, ┌ćÿdost: 2, ┌ćÿdrink: 2, ┌ćÿeach: 2, ┌ćÿearly: 3, ┌ćÿelsie: 2, ┌ćÿelsie.: 2, ┌ćÿenough: 2,
┌ćÿeven: 2, ┌ćÿeverything: 2, ┌ćÿfair: 2, ┌ćÿfaith.┌ćö: 2, ┌ćÿfalada.: 4, ┌ćÿfalse: 2, ┌ćÿfather: 8, ┌ćÿfather.┌ćö: 2
┌ćÿfear: 2, ┌ćÿfeel.┌ćö: 2, ┌ćÿfine: 4, ┌ćÿfirst: 2, ┌ćÿfirst.: 2, ┌ćÿfive.┌ćö: 2, ┌ćÿfool!┌ćö: 2, ┌ćÿfool: 2, ┌ćÿf
r: 10, ┌ćÿfourthly: 2, ┌ćÿfox: 2, ┌ćÿfrederick: 3, ┌ćÿfrederick: 3, ┌ćÿfrom: 2, ┌ćÿfull: 2, ┌ćÿfundevogel: 3, ┌ćÿga
ve: 4, ┌ćÿgently!: 2, ┌ćÿget: 6, ┌ćÿgive: 6, ┌ćÿgo: 18, ┌ćÿgo.: 2, ┌ćÿgod: 5, ┌ćÿgood: 41, ┌ćÿgood.┌ćö: 3, ┌ćÿgoodbye
: 17, ┌ćÿgoodness: 2, ┌ćÿgracious: 2, ┌ćÿgrete.: 4, ┌ćÿgretel: 2, ┌ćÿgrowler: 2, ┌ćÿha!┌ćö: 2, ┌ćÿha: 2, ┌ćÿhadst
: 2, ┌ćÿhalf-done!: 2, ┌ćÿhalf-done.┌ćö: 2, ┌ćÿhans: 2, ┌ćÿhansel: 5, ┌ćÿhark: 6, ┌ćÿhas: 5, ┌ćÿhave: 10, ┌ćÿhe: 21,
┌ćÿheads: 3, ┌ćÿhearken: 2, ┌ćÿheaven: 6, ┌ćÿhere: 9, ┌ćÿhere: 3, ┌ćÿhere.┌ćö: 2, ┌ćÿhercös: 2, ┌ćÿhi!: 3, ┌ćÿhick
: 2, ┌ćÿhis: 2, ┌ćÿhither: 3, ┌ćÿhome!: 2, ┌ćÿhow: 36, ┌ćÿhow?: 2, ┌ćÿhowever: 2, ┌ćÿhullo!┌ćö: 2, ┌ćÿhurrah!: 2, ┌ćÿhu
rry: 2, ┌ćÿhusband.┌ćö: 3, ┌ćÿhusband.┌ćö: 5, ┌ćÿhush!: 2, ┌ćÿi: 285, ┌ćÿi: 2, ┌ćÿiff: 77, ┌ćÿin: 11, ┌ćÿiron: 2, ┌ćÿis
: 17, ┌ćÿit: 45, ┌ćÿit┌ćös: 3, ┌ćÿir┌ćöll: 8, ┌ćÿircöm: 3, ┌ćÿjip!: 2, ┌ćÿjip!┌ćö: 2, ┌ćÿjust: 14, ┌ćÿjustice: 2, ┌ćÿka
te!: 2, ┌ćÿkate: 4, ┌ćÿkeep: 5, ┌ćÿlaces: 2, ┌ćÿlate: 2, ┌ćÿlay: 2, ┌ćÿlearn: 2, ┌ćÿleave: 3, ┌ćÿlet: 19, ┌ćÿlie: 2
┌ćÿlift: 3, ┌ćÿlisten: 2, ┌ćÿlisten: 5, ┌ćÿlittle: 6, ┌ćÿlook: 5, ┌ćÿlook: 4, ┌ćÿlord: 3, ┌ćÿlord: 3, ┌ćÿmadam: 3,
: 3, ┌ćÿmake: 2, ┌ćÿmaster!: 2, ┌ćÿmaster: 2, ┌ćÿmay: 5, ┌ćÿmind: 2, ┌ćÿmine: 2, ┌ćÿmiserable: 2, ┌ćÿmistress: 2, ┌ćÿm
oney: 2, ┌ćÿmore: 3, ┌ćÿmost: 2, ┌ćÿmother: 3, ┌ćÿmother.┌ćö: 3, ┌ćÿmother: 2, ┌ćÿmurder!┌ćö: 2, ┌ćÿmust: 2, ┌ćÿmy: 35, ┌ćÿnay.┌ćö
: 5, ┌ćÿneither: 5, ┌ćÿnever: 13, ┌ćÿnibble: 2, ┌ćÿno!┌ćö: 2, ┌ćÿno: 4, ┌ćÿno: 16, ┌ćÿno.┌ćö: 26, ┌ćÿnobody: 2, ┌ćÿno
nsense!┌ćö: 2, ┌ćÿnot: 13, ┌ćÿnothing: 3, ┌ćÿnothing.: 2, ┌ćÿnow: 30, ┌ćÿnow: 15, ┌ćÿnow.┌ćö: 4, ┌ćÿnursery: 2, ┌ćÿo:
: 12, ┌ćÿo.: 2, ┌ćÿof: 3, ┌ćÿoh!: 4, ┌ćÿoh!┌ćö: 5, ┌ćÿoh: 3, ┌ćÿoh.: 35, ┌ćÿoh.┌ćö: 5, ┌ćÿoho!┌ćö: 2, ┌ćÿold: 3, ┌ćÿo
ne: 9, ┌ćÿonly: 2, ┌ćÿopen: 10, ┌ćÿor: 3, ┌ćÿour: 3, ┌ćÿout: 2, ┌ćÿorcöer: 2, ┌ćÿpearls: 2, ┌ćÿpeasant: 2, ┌ćÿperhap
s: 3, ┌ćÿpoor: 3, ┌ćÿpray: 12, ┌ćÿpray.: 2, ┌ćÿpray.┌ćö: 2, ┌ćÿprince: 4, ┌ćÿprincess: 2, ┌ćÿprithee: 2, ┌ćÿpromise
: 2, ┌ćÿpull: 2, ┌ćÿput: 3, ┌ćÿquick: 2, ┌ćÿrapunzel: 5, ┌ćÿreadily.┌ćö: 2, ┌ćÿright: 2, ┌ćÿriver: 2, ┌ćÿround: 3,
┌ćÿrun: 3, ┌ćÿsay: 2, ┌ćÿscour: 2, ┌ćÿsee: 2, ┌ćÿsee: 5, ┌ćÿsee.┌ćö: 4, ┌ćÿseek: 2, ┌ćÿseven: 3, ┌ćÿsew: 2, ┌ćÿshake
: 3, ┌ćÿshake: 4, ┌ćÿshall: 2, ┌ćÿshe: 20, ┌ćÿshow: 3, ┌ćÿsilly: 2, ┌ćÿsince: 3, ┌ćÿsing: 2, ┌ćÿsir: 2, ┌ćÿsit: 2,
┌ćÿsitting: 2, ┌ćÿsnow-white: 2, ┌ćÿsnow-white.: 2, ┌ćÿsnowdrop: 2, ┌ćÿso: 5, ┌ćÿsoftly: 4, ┌ćÿsome: 2, ┌ćÿsomeone: 2
, ┌ćÿsomething: 3, ┌ćÿsooner: 2, ┌ćÿspinning.┌ćö: 2, ┌ćÿstand: 2, ┌ćÿstanding: 2, ┌ćÿstay: 2, ┌ćÿstir: 2, ┌ćÿstop!: 2
┌ćÿstop!┌ćö: 2, ┌ćÿstop.: 2, ┌ćÿstrew: 2, ┌ćÿstrike: 2, ┌ćÿstuck: 3, ┌ćÿsuch: 3, ┌ćÿsuppose: 3, ┌ćÿsurely: 2, ┌ćÿsure
ly.┌ćö: 2, ┌ćÿsweep: 2, ┌ćÿtake: 18, ┌ćÿtell: 7, ┌ćÿthan: 2, ┌ćÿthank: 3, ┌ćÿthat: 67, ┌ćÿthat: 2, ┌ćÿthat┌ćös: 3, ┌ćÿ
ćÿthe: 52, ┌ćÿthen: 20, ┌ćÿthen.┌ćö: 3, ┌ćÿthere!: 2, ┌ćÿthere: 12, ┌ćÿthere: 7, ┌ćÿthere.┌ćö: 2, ┌ćÿthere┌ćös: 3, ┌ćÿ
these: 3, ┌ćÿthey: 13, ┌ćÿthirdly: 2, ┌ćÿthis: 20, ┌ćÿthose: 2, ┌ćÿthou: 11, ┌ćÿthou.: 7, ┌ćÿthree: 2, ┌ćÿthrow: 2,
┌ćÿthy: 2, ┌ćÿtill: 2, ┌ćÿtis: 3, ┌ćÿto: 16, ┌ćÿtomorrow: 3, ┌ćÿtomorrow.: 2, ┌ćÿtonight: 2, ┌ćÿtook: 4, ┌ćÿtop: 3, ┌
ćÿtop-off!: 2, ┌ćÿtread: 2, ┌ćÿtry: 2, ┌ćÿtu: 2, ┌ćÿturn: 4, ┌ćÿtwas: 2, ┌ćÿtwo: 2, ┌ćÿuncouth: 2, ┌ćÿunlucky: 4, ┌ćÿ
upon: 3, ┌ćÿvery: 8, ┌ćÿwait: 2, ┌ćÿwait: 2, ┌ćÿwait.┌ćö: 3, ┌ćÿwallface.┌ćö: 2, ┌ćÿwas: 2, ┌ćÿwe: 31, ┌ćÿwell: 4,
┌ćÿwell: 19, ┌ćÿwell!┌ćö: 14, ┌ćÿwhat!┌ćö: 2, ┌ćÿwhat: 128, ┌ćÿwhat: 2, ┌ćÿwhat?┌ćö: 2, ┌ćÿwhatever: 2, ┌ćÿwhat┌ćös: 5
, ┌ćÿwhen: 13, ┌ćÿwhere: 23, ┌ćÿwhere┌ćös: 2, ┌ćÿwhither: 8, ┌ćÿwho: 22, ┌ćÿwhoever: 2, ┌ćÿwhose: 4, ┌ćÿwhy: 19, ┌ćÿwhy
: 8, ┌ćÿwhy.┌ćö: 3, ┌ćÿwhy?┌ćö: 2, ┌ćÿwife: 8, ┌ćÿwife.┌ćö: 4, ┌ćÿwill: 4, ┌ćÿwillingly.┌ćö: 2, ┌ćÿwish: 2, ┌ćÿwith
: 17, ┌ćÿwoman: 2, ┌ćÿwould: 5, ┌ćÿyellow: 2, ┌ćÿyes: 17, ┌ćÿyes.┌ćö: 22, ┌ćÿyes: 2, ┌ćÿyes!┌ćö: 2, ┌ćÿyet.┌ćö: 2, ┌ć
yyonder: 2, ┌ćÿyou: 88, ┌ćÿyou.: 2, ┌ćÿyou.┌ćö: 2, ┌ćÿyour: 6, ┌ćÿyourcöll: 2, ┌ćÿyourcöre: 2, ┌ćfdefects, ┌ć
\udc9d: 2, ┌ćfgood: 2, ┌ćfheads: 2, ┌ćfhere: 3, ┌ćfi: 3, ┌ćfinformation: 2, ┌ćfiron: 2, ┌ćfito: 2, ┌ćfjip!┌ć\udc9d┌ćö:
2, ┌ćfmerrily: 2, ┌ćfplain: 3, ┌ćfproject: 6, ┌ćfright: 2, ┌ćfunder: 2, ┌ćfwhat: 2}
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib>
```

## Mapper:

It will read data from STDIN, remove any leading or trailing whitespaces, split rows into words and output tuple containing word and its count to STDOUT

```python
import sys

for line in sys.stdin:
    line = line.strip()
    words = line.split()  # split the line into words
    for word in words:
        word = (word, 1)
        print(word)
```

## Reducer:

It will read the results of mapper.py from STDIN (the output format of mapper.py and the expected input format of reducer.py must match).

Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple.

```python
dic = {}
for line in sys.stdin:
    key = literal_eval(line)[0]
```
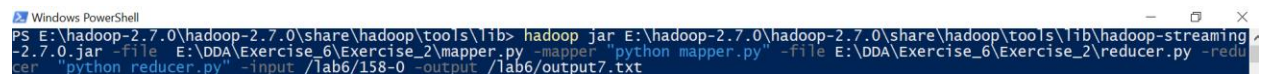
Then I check if word exists in dictionary or not. If word already exist then take its previous count and add 1 to it and save new count. If word does not exist in dictionary then I simply add word as key and count value as 1 in dictionary.

```python
    dic[key] = dic.get(key, 0) + 1
```

# Q2:

## Mapper:

It will read data from STDIN and remove any leading or trailing whitespaces. Then it read line character by character and skips any punctuations and number and concatenates other characters. Then split the line in words, check if that word exists in stop word dictionary or not. If not then output specific word with count to STDOUT

```python
punctuations = '''!()-[]{};:'"\,<>./?@#$%^&*_~'''
numbers='0123456789'
stop_words=['a','able','about','across','after','all','almost','also','am','among','a
n','and','any','are','as','at','be','because','been','but',

'by','can','cannot','could','dear','did','do','does','either','else','ever','every','
for','from','get','got','had','has','have','he',

'her','hers','him','his','how','however','i','if','in','into','is','it','its','just',
'least','let','like','likely','may','me','might',

'most','must','my','neither','no','nor','not','of','off','often','on','only','or','ot
her','our','own','rather','said','say','says',

'she','should','since','so','some','than','that','the','their','them','then','there',
'these','they','this','tis','to','too','twas',

'us','wants','was','we','were','what','when','where','which','while','who','whom','wh
y','will','with','would','yet','you','your']

no_punct_and_number = ""

for line in sys.stdin:
    line = line.strip()
    for char in line:
        if char not in punctuations and numbers:
            no_punct_and_number = no_punct_and_number + char
    words = no_punct_and_number.split()# split the line into words
    for word in words:
        if word not in stop_words:
            word = (word.lower(), 1)
            print(word)
```

Following line uses build in function of translation. maketrans first two parameters tells translate function to translate nothing to nothing and translate any punctuation or numbers to None (i.e. remove them). This function worked faster so I used this in my code. My implementation of removing punctuation and number was taking a lot of time.

```python
line = line.translate(str.maketrans('','',punctuation))
line = line.translate(str.maketrans('','','1234567890'))
line = line.strip()# split the line into words
```

## Reducer:

It will read the results of mapper.py from STDIN (the output format of mapper.py and the expected input format of reducer.py must match).

Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple.

```
dic = {}
for line in sys.stdin:
    key = literal_eval(line)[0]
```

Then I check if word exists in dictionary or not. If word already exists then take its previous count and add 1 to it and save new count. If word does not exist in dictionary then I simply add word as key and count value as 1 in dictionary.

```
    dic[key] = dic.get(key, 0) + 1
```

## Show the final results:

```
for k,v in dic.items():
    print('{} {}'.format(k,v))
```

# Running commands: All 5 files



```
Windows PowerShell                                                           —  □  ×
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming
-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_2\mapper.py -mapper "python mapper.py"  -file E:\DDA\Exercise_6\Exercise_2\reducer.py -redu
cer  "python reducer.py" -input /lab6/158-0 -output /lab6/output7.txt
```

                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=1854890
                HDFS: Number of bytes written=137641
                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=16633
                Map output records=82014
                Map output bytes=1233434
                Map output materialized bytes=1397468
                Input split bytes=86
                Combine input records=0
                Combine output records=0
                Reduce input groups=11142
                Reduce shuffle bytes=1397468
                Reduce input records=82014
                Reduce output records=11142
                Spilled Records=164028
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=0
                Total committed heap usage (bytes)=544210944
        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=927445
        File Output Format Counters
                Bytes Written=137641
20/06/22 09:37:59 INFO streaming.StreamJob: Output directory: /lab6/output7.txt
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop dfs -cat /lab6/output7.txt/*

## Windows PowerShell

```
ГÇ£well 33
ГÇ£wellas 1
ГÇ£wellbut 1
ГÇ£wellif 3
ГÇ£wellГÇ¥ 11
ГÇ£were 3
ГÇ£what 28
ГÇ£whatever 1
ГÇ£whatГÇ¥ 1
ГÇ£when 7
ГÇ£whenever 2
ГÇ£where 2
ГÇ£wheremay 1
ГÇ£which 3
ГÇ£while 1
ГÇ£who 4
ГÇ£whoever 3
ГÇ£whom 1
ГÇ£why 9
ГÇ£will 8
ГÇ£with 3
ГÇ£without 1
ГÇ£women 1
ГÇ£worse 1
ГÇ£would 4
ГÇ£writes 1
ГÇ£wrong 1
ГÇ£yes 40
ГÇ£yesa 1
ГÇ£yesentirely 1
ГÇ£yesi 3
ГÇ£yesit 1
ГÇ£yesrather 1
ГÇ£yesratheri 1
ГÇ£yeswhat 1
ГÇ£yesГÇ¥ 14
ГÇ£york 1
ГÇ£you 101
ГÇ£your 8
ГÇ£yours 1
```

```
ГÇ£yours 1
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming
-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_2\mapper.py -mapper "python mapper.py"  -file E:\DDA\Exercise_6\Exercise_2\reducer.py -redu
cer  "python reducer.py" -input /lab6/219-0.txt -output /lab6/output8.txt
```

```
                FILE: Number of bytes read=768544
                FILE: Number of bytes written=1713954
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=468178
                HDFS: Number of bytes written=73162
                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=3709
                Map output records=22869
                Map output bytes=335346
                Map output materialized bytes=381090
                Input split bytes=86
                Combine input records=0
                Combine output records=0
                Reduce input groups=6442
                Reduce shuffle bytes=381090
                Reduce input records=22869
                Reduce output records=6442
                Spilled Records=45738
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=117
                Total committed heap usage (bytes)=824180736
        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=234089
        File Output Format Counters
                Bytes Written=73162
20/06/22 09:39:40 INFO streaming.StreamJob: Output directory: /lab6/output8.txt
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib>
```

```
Γç£kurtz 1
Γç£love 1
Γç£mindΓç¥ 1
Γç£mistah 1
Γç£my 1
Γç£near 1
Γç£next 1
Γç£no 4
Γç£now 2
Γç£of 1
Γç£oh 1
Γç£on 1
Γç£one 5
Γç£plain 2
Γç£poor 1
Γç£project 5
Γç£repeat 1
Γç£right 1
Γç£she 6
Γç£some 2
Γç£sometimes 1
Γç£suddenly 2
Γç£the 16
Γç£there 2
Γç£they 3
Γç£this 3
Γç£through 1
Γç£thus 1
Γç£to 1
Γç£towards 1
Γç£true 2
Γç£try 1
Γç£two 1
Γç£unsound 1
Γç£we 7
Γç£what 1
Γç£when 3
Γç£yes 1
Γç£yet 1
Γç£you 5
```

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming
-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_2\mapper.py -mapper "python mapper.py" -file E:\DDA\Exercise_6\Exercise_2\reducer.py -redu
cer  "python reducer.py" -input /lab6/1400-0.txt -output /lab6/output8.txt
```

```
ΓÇ£wretchesΓÇ¥ 1
ΓÇ£yah 1
ΓÇ£yahΓÇ¥ 2
ΓÇ£ye 1
ΓÇ£yes 76
ΓÇ£yesΓÇ¥ 32
ΓÇ£yet 3
ΓÇ£yetΓÇ¥ 1
ΓÇ£yonderΓÇ¥ 1
ΓÇ£you 191
ΓÇ£young 4
ΓÇ£your 11
ΓÇ£yours 1
ΓÇ£yoursΓÇ¥ 1
ΓÇ£youΓÇ¥ 1
ΓÇ£youΓÇöd 1
ΓÇ£youΓÇöll 2
ΓÇ£youΓÇöre 11
ΓÇ£youΓÇöve 3
ΓÇ£ΓÇöas 1
ΓÇ£ΓÇöat 1
ΓÇ£ΓÇöby 1
ΓÇ£ΓÇöhad 1
ΓÇ£ΓÇöinvest 1
ΓÇ£ΓÇöthat 1
ΓÇ£ΓÇöthen 1
ΓÇ£ΓÇöthereΓÇös 1
ΓÇ£ΓÇöwhich 2
ΓÇ£ΓÇöyes 1
ΓÇ£ΓÇÿaccount 1
ΓÇ£ΓÇÿeat 1
ΓÇ£ΓÇÿgod 1
ΓÇ£ΓÇÿhe 1
ΓÇ£ΓÇÿi 1
ΓÇ£ΓÇÿjoseph 1
ΓÇ£ΓÇÿluck 1
ΓÇ£ΓÇÿshe 1
ΓÇ£ΓÇÿto 1
ΓÇ£ΓÇÿwhat 1
ΓÇ£ΓÇÿyes 1
```

```
ΓÇ£ΓÇÿyes 1
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming
-2.7.0.jar -file E:\DDA\Exercise_6\Exercise_2\mapper.py -mapper "python mapper.py" -file E:\DDA\Exercise_6\Exercise_2\reducer.py -redu
cer "python reducer.py" -input /lab6/4300-0.txt -output /lab6/output9.txt
```

```
Windows PowerShell
Çöyesterday 1
Çöyou 45
Çöyour 4
ÇöyourÇöre 10
Çözinfandel 1
ÇörÇö 24
ÇörÇörÇörÇörÇö 2
ÇörÇölÌdo 1
ÇörÇötis 2
ÇÖ 2
ÇÖem 2
ÇÖmid 1
ÇÖneath 1
ÇÖpon 1
ÇÖs 1
ÇÖslife 1
ÇÖtis 13
ÇÖtwas 14
ÇÖtwere 1
ÇÖtwixt 1
Çf come 1
ÇfjÇ¥ 1
ÇfviatorÇ¥ 1
Çfyou 1
ÇÓ 1
f 1
ú 31
⣿ 3
Ⲁ 9
rç 1
Ⳁangus 3
rça 1
rê 1
rëlus 1
f£bermensch 2
ⳇ 3
r-clat 1
r-lite 3
r-tat 1
rülΣüvΣütΣür 1
```



```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming
-2.7.0.jar -file E:\DDA\Exercise_6\Exercise_2\mapper.py -mapper "python mapper.py" -file E:\DDA\Exercise_6\Exercise_2\reducer.py -redu
cer "python reducer.py" -input /lab6/2591-0.txt -output /lab6/output5.txt
```



```
Windows PowerShell
rçÝwho 21
rçÝwhoever 1
rçÝwhose 3
rçÝwhy 25
rçÝwhyrçö 3
rçÝwife 7
rçÝwiferçö 3
rçÝwill 3
rçÝwillinglyrçö 1
rçÝwish 1
rçÝwith 16
rçÝwoman 1
rçÝwould 4
rçÝyellow 1
rçÝyes 17
rçÝyesrçö 22
rçÝyetrçö 1
rçÝyonder 1
rçÝyou 88
rçÝyour 5
rçÝyourçö 1
rçÝyourçöll 1
rçÝyourçöre 1
rçfah 1
rçfdefectsrç¥ 1
rçfgood 1
rçfheads 1
rçfhere 2
rçfi 2
rçfinformation 1
rçfiron 1
rçfit 1
rçfjiprç¥rçö 1
rçfmerrily 1
rçfplain 2
rçfproject 5
rçfright 1
rçfthe 1
rçfunder 1
rçfwhat 1
```

# Q3:

## Mapper:

It will read data from STDIN and get file names from environment variable "`map_input_file`". These are input files names which are passes as command line argument. Then it remove any punctuations or number from line. Then remove any leading or trailing whitespaces, split the line in words, check if that word exists in stop word dictionary or not. If not then output specific word with count to STDOUT

```python
for line in sys.stdin:
    filename = os.environ["map_input_file"]
    line = line.translate(str.maketrans('','',punctuation))
    line = line.translate(str.maketrans('','','1234567890'))
    line = line.strip()# remove leading and trailing whitespace
    words = line.split() # split the line into words
    for word in words:
        if word not in stop_words:
            print((filename,word,1))
```

## Reducer:

It will read the results of mapper.py from STDIN (the output format of mapper.py and the expected input format of reducer.py must match).

Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple. After converting it to tuple we only take first two arguments (filename, word) of tuple as third argument is just 1 in all cases.

```python
import sys
from ast import literal_eval
dic = {}
for line in sys.stdin:
    key = literal_eval(line)[:2]
```

Then I check if word exists in dictionary or not. If word already exists then take its previous count and add 1 to it and save new count. If word does not exist in dictionary then I simply add word as key and count value as 1 in dictionary.

```python
    dic[key] = dic.get(key , 0) + 1
```

Then we output the tuple as following format (filename, word, count of words in that filename).

```python
for k,v in dic.items():
    print((k[0],k[1],v))
```

## Mapper2:

It will read data from STDIN (output of reducer1). Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple. Then this tuple (filename, word, count of words in that filename) is passed as output to STDOUT

```python
for line in sys.stdin:
    filename_word_count = literal_eval(line)
    print(filename_word_count)
```

## Reduce2:

It will read the results of mapper.py from STDIN (the output format of mapper.py and the expected input format of reducer.py must match).

Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple. After converting it to tuple we calculate total words in each file. Also we save tuple that we received from STDIN which will help use later.

```python
for line in sys.stdin:

    filename, word, count = literal_eval(line)
    count = int(count)
    if prev_filename == filename:
        N = N + count
    else:
        if prev_filename != None:
            file_name_with_total_wordCount[prev_filename] = N
            # print(file_name_with_total_wordCount[prev_filename])

        N = 0
        prev_filename = filename
    saved_previous_data.append(line)
file_name_with_total_wordCount[prev_filename] = N #saved last file count
```

Then we map the list of tuple that we saved previously to total words in the specific files from which the word in tuple belong. Finally we output the tuple as following format (filename, word, count of words in that filename, total_words_count_in_specific file).

```
for line in saved_previous_data:
    filename, word, count = literal_eval(line)
    for k in file_name_with_total_wordCount.keys():
        if filename == k:
            print((word, filename, count, file_name_with_total_wordCount[k]))
```

## Mapper3:

It will read data from STDIN (output of reducer1). Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple. Then this tuple (filename, word, count of words in that filename, total_words_count_in_specific file,1)  is passed as output to STDOUT

```
for line in sys.stdin:
    word_filename_count_totalcount = literal_eval(line)

print((word_filename_count_totalcount[0],word_filename_count_totalcount[1],word_filen
ame_count_totalcount[2],word_filename_count_totalcount[3],1))
```

## Reduce3:

It will read the results of mapper.py from STDIN (the output format of mapper.py and the expected input format of reducer.py must match).

Since the value received from STDIN is in string format so first we need to convert this string tuple to tuple. After converting it to tuple we calculate number of documents that contains token t. Also we save tuple (word, filename) that which will help use later.

```
for line in sys.stdin:
    word, filename, wordCount, total_WordCount, count = literal_eval(line)
    if prev_word == word:
        total_count = total_count + int(count)
    else:
        if prev_word != None:
            df[prev_word] = (wordCount, total_WordCount, total_count)
            word_filename = (prev_word, filename)
            saved_previous_data.append(word_filename)
        total_count = 1
        prev_word = word

df[prev_word] = (wordCount, total_WordCount, total_count)
```

```
word_filename = (prev_word, filename)
saved_previous_data.append(word_filename)
```

Now we have a list of tuple which contain word and filename. Then we have dictionary that has word as key and tuple (wordCount, total_WordCount, total_count_of_tokenK_i n document) as value.

Then first I take token (word)  and get its  data from dictionary.Then I used this data to find tfifd of token using formula given in exercise

```
for line in saved_previous_data:
    word, filename = line
    for k in df.keys():
        if word == k:
            wordCount, total_WordCount, total_count_of_word_in_different_document =
df[k]
            tfidf = (wordCount / total_WordCount) * log10(5 /
total_count_of_word_in_different_document)
            print('{}: {}'.format(word, tfidf))
```

Running commands:

First mapper, reducer:

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming
-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_3\mapper.py -mapper "python mapper.py" -file E:\DDA\Exercise_6\Exercise_3\reducer.py wsredu
cer "python reducer.py" -input /lab6/* -output /lab6/output1.txt
```

```
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwife', 2)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwill', 8)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwine', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwise', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwith', 4)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwithout', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwoa', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwoke', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwonder', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöwould', 3)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöyes', 76)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöyesterday', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöyou', 45)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöyour', 4)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöyourçore', 10)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçözinfandel', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöhome', 2)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçörçö', 24)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçörçörçörçörçö', 2)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçörçöTis', 2)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçörçölldo', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçö', 2)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöslife', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöTis', 9)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöTwas', 9)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöTwixt', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöem', 2)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçömid', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöneath', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçöpon', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçös', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçötis', 4)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçötwas', 5)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçötwere', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçfcome', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçfJrç\udc9d', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçfviatorrç\udc9d', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçfyou', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rçö', 1)
('hdfs://0.0.0.0:9000/lab6/4300-0.txt', 'rf', 1)
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib>
```

## Second mapper reducer:

```
>> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_3\mapper
2.py -mapper "python mapper2.py" -file E:\DDA\Exercise_6\Exercise_3\reducer2.py -reducer  "python reducer2.py" -input /lab6/output1.txt
/part-00000 -output /lab6/output2.txt
```

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop dfs -cat /lab6/output2.txt/*
```

## Third mapper reducer:

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib>
>> hadoop jar E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib\hadoop-streaming-2.7.0.jar -file  E:\DDA\Exercise_6\Exercise_3\mappe
3.py -mapper "python mapper3.py" -file E:\DDA\Exercise_6\Exercise_3\reducer3.py -reducer  "python reducer3.py" -input /lab6/output2.t
/part-00000 -output /lab6/output6.txt
20/06/22 09:15:40 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [E:\DDA\Exercise_6\Exercise_3\mapper3.py, E:\DDA\Exercise_6\Exercise_3\reducer3.py] [] C:\Users\fahad\AppData\Local\Te
\streamjob6658525403061605237.jar tmpDir=null
20/06/22 09:15:41 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
20/06/22 09:15:41 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
20/06/22 09:15:41 INFO jvm.JvmMetrics: Cannot initialize JVM Metrics with processName=JobTracker, sessionId= - already initialized
20/06/22 09:15:42 INFO mapred.FileInputFormat: Total input paths to process : 1
20/06/22 09:15:42 INFO mapreduce.JobSubmitter: number of splits:1
20/06/22 09:15:43 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1126490717_0001
20/06/22 09:15:43 INFO mapred.LocalDistributedCacheManager: Localized file:/E:/DDA/Exercise_6/Exercise_3/mapper3.py as file:/tmp/hado
-fahad/mapred/local/1592842543264/mapper3.py
20/06/22 09:15:43 INFO mapred.LocalDistributedCacheManager: Localized file:/E:/DDA/Exercise_6/Exercise_3/reducer3.py as file:/tmp/had
p-fahad/mapred/local/1592842543265/reducer3.py
20/06/22 09:15:43 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
20/06/22 09:15:43 INFO mapreduce.Job: Running job: job_local1126490717_0001
20/06/22 09:15:43 INFO mapred.LocalJobRunner: OutputCommitter set in config null
20/06/22 09:15:43 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
20/06/22 09:15:43 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
20/06/22 09:15:43 INFO mapred.LocalJobRunner: Waiting for map tasks
20/06/22 09:15:43 INFO mapred.LocalJobRunner: Starting task: attempt_local1126490717_0001_m_000000_0
20/06/22 09:15:43 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
20/06/22 09:15:43 INFO util.ProcfsBasedProcessTree: ProcfsBasedProcessTree currently is supported only on Linux.
20/06/22 09:15:44 INFO mapred.Task:  Using ResourceCalculatorProcessTree : org.apache.hadoop.yarn.util.WindowsBasedProcessTree@4de7224

20/06/22 09:15:44 INFO mapred.MapTask: Processing split: hdfs://0.0.0.0:9000/lab6/output2.txt/part-00000:0+4767116
20/06/22 09:15:44 INFO mapred.MapTask: numReduceTasks: 1
20/06/22 09:15:44 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
20/06/22 09:15:44 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
20/06/22 09:15:44 INFO mapred.MapTask: soft limit at 83886080
20/06/22 09:15:44 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
20/06/22 09:15:44 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
20/06/22 09:15:44 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
20/06/22 09:15:44 INFO streaming.PipeMapRed: PipeMapRed exec [python, mapper3.py]
20/06/22 09:15:44 INFO Configuration.deprecation: mapred.work.output.dir is deprecated. Instead, use mapreduce.task.output.dir
20/06/22 09:15:44 INFO Configuration.deprecation: map.input.start is deprecated. Instead, use mapreduce.map.input.start
```

```
                FILE:  Number of bytes read=10294656
                FILE:  Number of bytes written=16007034
                FILE:  Number of read operations=0
                FILE:  Number of large read operations=0
                FILE:  Number of write operations=0
                HDFS:  Number of bytes read=9534232
                HDFS:  Number of bytes written=1644263
                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=75568
                Map output records=75568
                Map output bytes=4993822
                Map output materialized bytes=5144966
                Input split bytes=99
                Combine input records=0
                Combine output records=0
                Reduce input groups=75568
                Reduce shuffle bytes=5144966
                Reduce input records=75568
                Reduce output records=50422
                Spilled Records=151136
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=0
                Total committed heap usage (bytes)=546308096
        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=4767116
        File Output Format Counters
                Bytes Written=1644263
20/06/22 09:20:27 INFO streaming.StreamJob: Output directory: /lab6/output6.txt
```

```
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib> hadoop dfs -cat /lab6/output6.txt/*
```

Windows PowerShell
ΓÇ£whose: 7.019251090450 9876e-06
ΓÇ£why: 7.0192510904509876e-06
ΓÇ£wife: 2.1057753271352963e-05
ΓÇ£will: 1.9981120952813223e-05
ΓÇ£with: 7.0192510904509876e-06
ΓÇ£without: 3.996224190562645e-06
ΓÇ£would: 3.996224190562645e-06
ΓÇ£wouldnΓÇöt: 8.532348685742418e-06
ΓÇ£writes: 7.0192510904509876e-06
ΓÇ£yes: 3.062300128525822e-05
ΓÇ£yet: 0.0002877892947 0849046
ΓÇ£you: 3.996224190562645e-06
ΓÇ£young: 1.4038502180901975e-05
ΓÇ£your: 3.996224190562645e-06
ΓÇ£yourΓÇöll: 2.1057753271352963e-05
ΓÇ£yourΓÇöre: 7.0192510904509876e-06
ΓÇ£yourΓÇöve: 7.0192510904509876e-06
ΓÇ£ΓÇöAt: 7.0192510904509876e-06
ΓÇ£ΓÇöBy: 7.0192510904509876e-06
ΓÇ£ΓÇöHad: 7.0192510904509876e-06
ΓÇ£ΓÇöInvest: 7.0192510904509876e-06
ΓÇ£ΓÇöThat: 7.0192510904509876e-06
ΓÇ£ΓÇöThen: 1.4038502180901975e-05
ΓÇ£ΓÇöwhich: 7.0192510904509876e-06
ΓÇ£ΓÇöYes: 7.0192510904509876e-06
ΓÇ£ΓÇöas: 7.0192510904509876e-06
ΓÇ£ΓÇöthereΓÇös: 7.0192510904509876e-06
ΓÇ£ΓÇÿEat: 7.0192510904509876e-06
ΓÇ£ΓÇÿGod: 7.0192510904509876e-06
ΓÇ£ΓÇÿI: 7.0192510904509876e-06
ΓÇ£ΓÇÿJoseph: 7.0192510904509876e-06
ΓÇ£ΓÇÿLuck: 7.0192510904509876e-06
ΓÇ£ΓÇÿShe: 7.0192510904509876e-06
ΓÇ£ΓÇÿTo: 7.0192510904509876e-06
ΓÇ£ΓÇÿwhat: 7.0192510904509876e-06
ΓÇ£ΓÇÿYes: 7.0192510904509876e-06
ΓÇ£ΓÇÿaccount: 7.0192510904509876e-06
ΓÇ£ΓÇÿhe: 4.285241365303498e-06
ΓÇö: 4.285241365303498e-06
Γ£: 4.285241365303498e-06
PS E:\hadoop-2.7.0\hadoop-2.7.0\share\hadoop\tools\lib>