

## Assignment 01: Determine the wordcount

The comments/sections provided are your cues to perform the assignment. You don't need to limit yourself to the number of rows/cells provided. You can add additional rows in each section to add more lines of code.

If at any point in time you need help on solving this assignment, view our demo video to understand the different steps of the code.

Happy coding!

## Determine the word count of the given Amazon dataset:

### DESCRIPTION

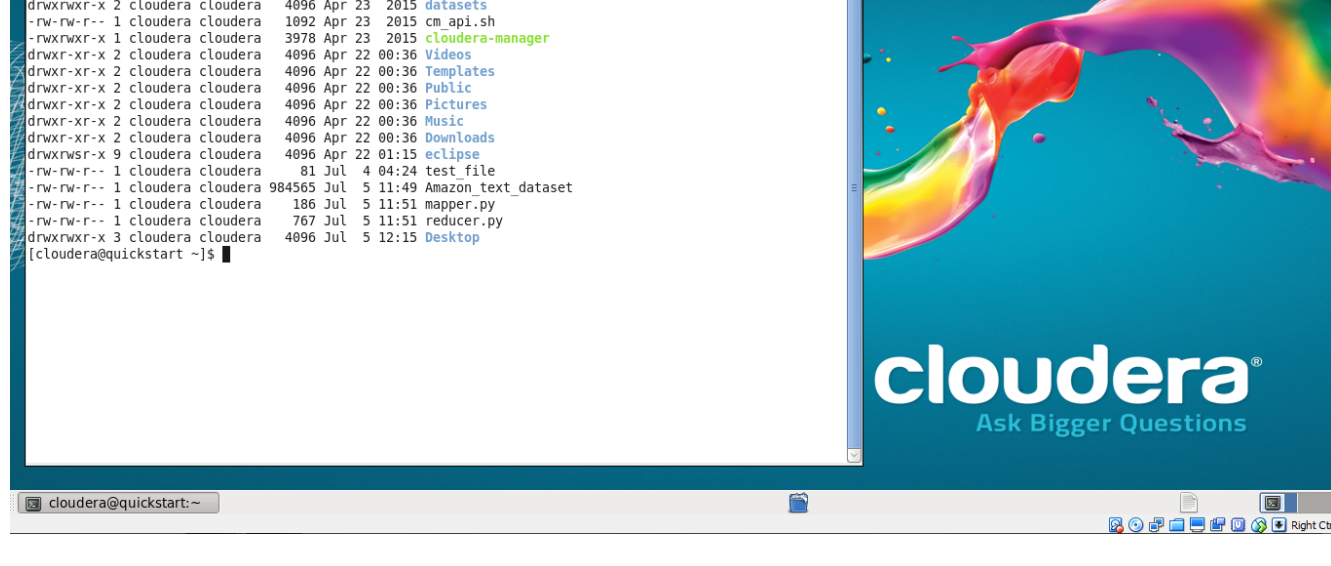
**Problem:**

- Create a MapReduce program to determine the word count of the Amazon dataset.
- Submit the MapReduce task to HDFS and run it.
- Verify the output.

Create a MapReduce program to determine the word count of the Amazon dataset

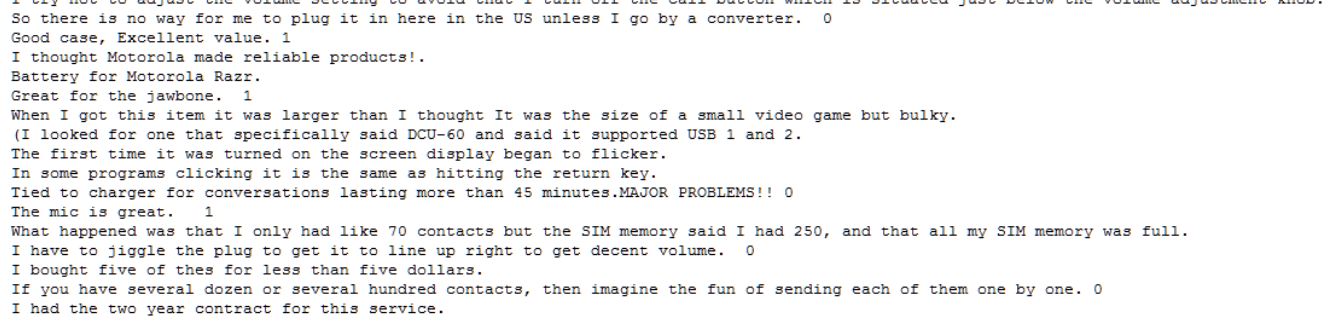
Copy the Dataset to Home Directory

- Copy the Amzon\_text\_dataset to cloudera Home Directory.
- Open the Terminal Application.
- Enter "ls -lrt" to list the folders and files in Home directory.



View the Dataset

- View the dataset using "cat Amazon\_text\_data" command.



Create the Mapper code

- Create mapper python program in Vi editor using "Vi mapper.py" command.
- Use "cat mapper.py" to view the contents of the file.

```
[cloudera@quickstart ~]$ cat mapper.py
#!/usr/bin/env python

import sys

for line in sys.stdin:
    #remove loading and trailing spaces
    line = line.strip()
    words = line.split()

    for word in words:
        print("%s\t%s"%(word,1))
```

Create the Reducer code

- Create reducer python program in Vi editor using "Vi reducer.py" command.
- Use "cat reducer.py" to view the contents of the file.

```
[cloudera@quickstart ~]$ cat reducer.py
import sys
from operator import itemgetter

current_word = None
current_count = 0
word = None

#input comes from console
for line in sys.stdin:
    #remove leading and trailing whitespace
    line = line.strip()

    #parse the input we got from mapper.py
    word, count = line.split('\t',1)

    #convert count (currently a string) to int
    try:
        count= int(count)
    except ValueError:
        #ignore if the count is not a number
        continue

    #use the sorted map output
    if current_word == word:
        current_count+=count
    else:
        if current_word:
            #write result in console
            print("%s\t%s" %(current_word, current_count))
        current_count = count
        current_word = word

#finally output the word to the console
if current_word == word:
    print("%s\t%s" %(current_word,current_count))
```

Execute the Mapper program over the Amazon\_text\_dataset

- Use "cat Amazon\_text\_dataset | python mapper.py" comment to execute the Mapper program over the Amazon\_text\_dataset.

```
[cloudera@quickstart ~]$ cat Amazon_text_dataset | python mapper.py
I 1
try 1
not 1
to 1
adjust 1
the 1
volume 1
setting 1
to 1
avoid 1
that 1
I 1
turn 1
off 1
the 1
call 1
button 1
which 1
is 1
situated 1
-----
```

Submit the MapReduce task to HDFS and run it.

Copy the Amazon\_text\_dataset file to HDFS location

- Use "hdfs dfs -put Amazon\_text\_dataset /user/cloudera" to copy the "Amazon\_text\_dataset" to HDFS location.
- Use "hdfs dfs -ls /user/cloudera" to list the folders and files in HDFS location.

```
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera
Found 7 items
-rw-r--r-- 1 cloudera cloudera 984565 2021-07-05 11:55 /user/cloudera/Amazon_text_dataset
-rw-r--r-- 1 cloudera cloudera 81 2021-07-04 03:34 /user/cloudera/test_file
drwxr-xr-x - cloudera cloudera 0 2021-07-04 03:46 /user/cloudera/wc_output01
drwxr-xr-x - cloudera cloudera 0 2021-07-04 04:20 /user/cloudera/wc_output02
drwxr-xr-x - cloudera cloudera 0 2021-07-04 04:31 /user/cloudera/wc_output03
drwxr-xr-x - cloudera cloudera 0 2021-07-05 12:11 /user/cloudera/wc_output04
drwxr-xr-x - cloudera cloudera 0 2021-07-05 21:14 /user/cloudera/wc_output05
```

Submit the MapReduce Task

- Using below command, submit the MapReduce task, add the below line in bash.
- Once the task is submitted, it runs the mapper first and then runs the reducer program over the dataset.

**hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-jar -file /home/cloudera/mapper.py /home/cloudera/reducer.py -mapper "python mapper.py" -reducer "python reducer.py" -input /user/cloudera/Amazon\_text\_dataset -output /user/cloudera/wc\_output06**

```
[cloudera@quickstart ~]$ hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-jar -file /home/cloudera/mapper.py /home/cloudera/reducer.py
-mapper "python mapper.py" -reducer "python reducer.py" -input /user/cloudera/Amazon_text_dataset -output /user/cloudera/wc_output06
21/07/05 21:27:06 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/home/cloudera/mapper.py, /home/cloudera/reducer.py] [/usr/lib/hadoop-mapreduce/hadoop-streaming-2.6.0-cdh5.4.0.jar] /tmp/streamjob1614669751
21/07/05 21:27:15 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.18032
21/07/05 21:27:19 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.18032
21/07/05 21:27:23 INFO mapred.FileInputFormat: Total input paths to process : 1
21/07/05 21:27:23 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1625542583719_0002
21/07/05 21:27:27 INFO impl.YarnClientImpl: Submitted application application_1625542583719_0002/
21/07/05 21:27:27 INFO mapreduce.Job: Running job: job_1625542583719_0002
21/07/05 21:28:07 INFO mapreduce.Job: Job job_1625542583719_0002 running in uber mode : false
21/07/05 21:28:07 INFO mapreduce.Job: map 0% reduce 0%
21/07/05 21:28:07 INFO mapreduce.Job: map 16% reduce 0%
21/07/05 21:28:11 INFO mapreduce.Job: map 44% reduce 0%
21/07/05 21:28:15 INFO mapreduce.Job: map 60% reduce 0%
21/07/05 21:28:18 INFO mapreduce.Job: map 67% reduce 0%
21/07/05 21:28:22 INFO mapreduce.Job: map 100% reduce 0%
21/07/05 21:28:27 INFO mapreduce.Job: map 100% reduce 77%
21/07/05 21:28:30 INFO mapreduce.Job: map 100% reduce 100%
21/07/05 21:29:31 INFO mapreduce.Job: Job job_1625542583719_0002 completed successfully
21/07/05 21:29:32 INFO mapreduce.Job: Counters: 49
File System Counters
  FILE: Number of bytes read=1693573
  FILE: Number of bytes written=3728548
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=191772
  HDFS: Number of bytes written=191772
  HDFS: Number of read operations=0
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=0
  MapReduce Counters
    input split bytes=204
    Combine input records=0
    Combine output records=0
    Reduce input groups=18818
    Reduce shuffle bytes=1693579
    Reduce input records=180753
    Reduce output records=18818
    Spilled Records=361506
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time spent (ms)=923
    CPU time spent (ms)=9430
    Physical memory (bytes) snapshot=550928384
    Virtual memory (bytes) snapshot=4511481856
    Total committed heap usage (bytes)=312680448
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=987899
File Output Format Counters
  Bytes Written=191772
21/07/05 21:29:32 INFO streaming.StreamJob: Output directory: /user/cloudera/wc_output06
```

Verify the output

- Verify the output folder present in HDFS location using below command.
  - hdfs dfs -ls /user/cloudera

```
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera
Found 8 items
-rw-r--r-- 1 cloudera cloudera 984565 2021-07-05 11:55 /user/cloudera/Amazon_text_dataset
-rw-r--r-- 1 cloudera cloudera 81 2021-07-04 03:34 /user/cloudera/test_file
drwxr-xr-x - cloudera cloudera 0 2021-07-04 03:46 /user/cloudera/wc_output01
drwxr-xr-x - cloudera cloudera 0 2021-07-04 04:20 /user/cloudera/wc_output02
drwxr-xr-x - cloudera cloudera 0 2021-07-04 04:31 /user/cloudera/wc_output03
drwxr-xr-x - cloudera cloudera 0 2021-07-05 12:11 /user/cloudera/wc_output04
drwxr-xr-x - cloudera cloudera 0 2021-07-05 21:14 /user/cloudera/wc_output05
drwxr-xr-x - cloudera cloudera 0 2021-07-05 21:29 /user/cloudera/wc_output06
```

- To check the Output on HDFS use the below command.
  - hdfs dfs -ls /user/cloudera/wc\_output06

```
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/wc_output06
Found 2 items
-rw-r--r-- 1 cloudera cloudera 0 2021-07-05 21:29 /user/cloudera/wc_output06/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 191772 2021-07-05 21:29 /user/cloudera/wc_output06/part-00000
```

- To View the part file use below command
  - hdfs dfs -cat /user/cloudera/wc\_output06/part\*

```
[cloudera@quickstart ~]$ hdfs dfs -cat /user/cloudera/wc_output06/part*
! 17
!! 4
!!! 2
!!!! 1
!!!!. 1
!!!!. 1
!!.. 1
!$$$ 1
!.. 1
!..I 1
!...ridiculous, 1
!1) 1
!1. 1
!2. 1
!At 1
!BURN 1
!Do 1
!GONNA 1
!Gonna 1
!Hope 1
!I 2
!I've 1
!If 1
```