# BIG DATA ANALYTICS

## Practical 6

Analyse impact of different number of mapper and reducer
on same definition as practical 4.

**Prepared By**

18BCE151 - Parth Kapadia

# Creating Mapper and Reducer

```python
mapper.py > ...
1   #!/usr/bin/env python
2
3   # import sys because we need to read and write data to STDIN and STDOUT
4   import sys
5
6   # reading entire line from STDIN (standard input)
7   for line in sys.stdin:
8       # to remove leading and trailing whitespace
9       line = line.strip()
10      # split the line into words
11      words = line.split()
12
13      # we are looping over the words array and printing the word
14      # with the count of 1 to the STDOUT
15      for word in words:
16          # write the results to STDOUT (standard output);
17          # what we output here will be the input for the
18          # Reduce step, i.e. the input for reducer.py
19          print ('%s\t%s' % (word, 1))
```

```python
reducer.py > ...
1   #!/usr/bin/env python
2
3   from operator import itemgetter
4   import sys
5
6   current_num = None
7   current_count = 0
8   num = None
9
10  g_largest=0
11  g_avg = 0.0
12  g_count = 0
13  dict = {}
14
15
16  # read the entire line from STDIN
17  for line in sys.stdin:
18      # remove leading and trailing whitespace
19      line = line.strip()
20      # slpiting the data on the basis of tab we have provided in mapper.py
21      num, count = line.split('\t', 1)
22      # convert count (currently a string) to int
```

```python
    try:
        count = int(count)
        num = int(num)
    except ValueError:
        # count was not a number, so silently
        # ignore/discard this line
        continue

    # this IF-switch only works because Hadoop sorts map output
    # by key (here: num) before it is passed to the reducer

    if current_num == num:
        current_count += count
    else:
        if current_num:
            # write result to STDOUT
            #print ('%s\t%s' % (current_num, current_count))
            dict[current_num] = current_count
        current_count = count
        current_num = num

# do not forget to output the last num if needed!
if current_num == num:
    #print ('%s\t%s' % (current_num, current_count))
    dict[current_num] = current_count

# Find avg, largest

for key, value in dict.items():
    g_avg += key*value
    g_count += value
    g_largest = max(g_largest, key)

g_avg/=g_count


print("The largest Integer is: ", g_largest)
print("The avg of Integers is: ", g_avg)
print("Total Distinct integers are:", len(dict))
print("All integers are:")

for key,value in dict.items():
    print(key)
```
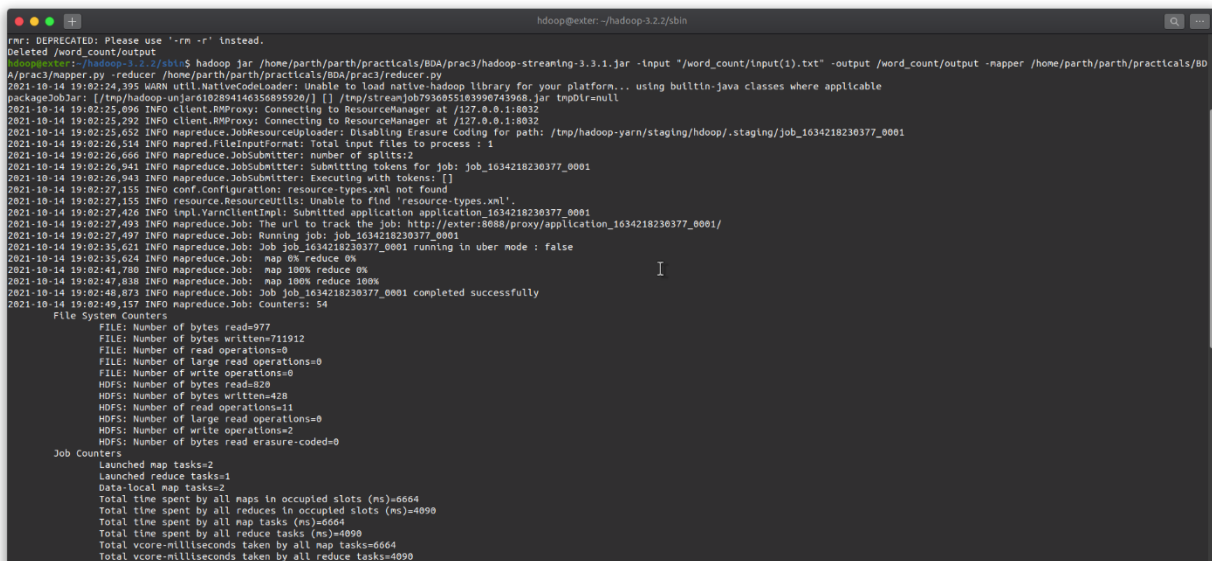
# Running the Map Reduce program

hadoop jar /home/parth/parth/practicals/BDA/prac3/hadoop-streaming-3.3.1.jar

-input "/word_count/input(1).txt"

-output /word_count/output

-mapper /home/parth/parth/practicals/BDA/prac3/mapper.py

-reducer /home/parth/parth/practicals/BDA/prac3/reducer.py

```
                Total time spent by all maps in occupied slots (ms)=6664
                Total time spent by all reduces in occupied slots (ms)=4090
                Total time spent by all map tasks (ms)=6664
                Total time spent by all reduce tasks (ms)=4090
                Total vcore-milliseconds taken by all map tasks=6664
                Total vcore-milliseconds taken by all reduce tasks=4090
                Total megabyte-milliseconds taken by all map tasks=6823936
                Total megabyte-milliseconds taken by all reduce tasks=4188160
        Map-Reduce Framework
                Map input records=7
                Map output records=140
                Map output bytes=691
                Map output materialized bytes=983
                Input split bytes=194
                Combine input records=0
                Combine output records=0
                Reduce input groups=74
                Reduce shuffle bytes=983
                Reduce input records=140
                Reduce output records=78
                Spilled Records=280
                Shuffled Maps =2
                Failed Shuffles=0
                Merged Map outputs=2
                GC time elapsed (ms)=240
                CPU time spent (ms)=3300
                Physical memory (bytes) snapshot=806744064
                Virtual memory (bytes) snapshot=7675355136
                Total committed heap usage (bytes)=835715072
                Peak Map Physical memory (bytes)=303910912
                Peak Map Virtual memory (bytes)=2556661760
                Peak Reduce Physical memory (bytes)=200216576
                Peak Reduce Virtual memory (bytes)=2563387392
        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=626
        File Output Format Counters
                Bytes Written=428
2021-10-14 19:02:49,158 INFO streaming.StreamJob: Output directory: /word_count/output
hdoop@exter:~/hadoop-3.2.2/sbin$
```

## Output

hdfs dfs -cat /word_count/output/part-00000



## Changing number of mappers and reducers

-Dmapreduce.job.reduces=3 -Dmapreduce.job.maps=2

hadoop jar /home/parth/parth/practicals/BDA/prac3/hadoop-streaming-3.3.1.jar -Dmapreduce.job.reduces=3 -Dmapreduce.job.maps=2 -input /word_count/Essay700kB.txt -output /word_count/output -mapper

/home/parth/parth/practicals/BDA/prac6/mapper.py -reducer
/home/parth/parth/practicals/BDA/prac6/reducer.py

## Execution time of Job (in ns)

### wrt to number of mappers



## Execution time of Job (in ns)

### wrt to reducers