

Name: Xinrun Zhang

Date: 15/10/2018

## Task1

1. Test mapper.py and reducer.py:

```
zhangxinrun@ubuntu:~/handson-2/task1$ cat pg20417.txt | ./mapper.py | sort | ./reducer.py
('the', 9028)
('of', 5530)
('and', 2886)
('a', 2715)
('in', 2336)
('to', 2247)
('is', 2120)
('it', 1383)
('that', 1349)
('are', 947)
```

2. Make a new directory in Hadoop cluster:

```
zhangxinrun@ubuntu:~$ hdfs dfs -mkdir /user/bigdata/handson-2/task1
```

3. Copy the pg20417.txt file into Hadoop cluster:

```
zhangxinrun@ubuntu:~$ hdfs dfs -copyFromLocal /home/zhangxinrun/handson-2/task1/pg20417.txt /user/bigdata/handson-2/task1
```

4. Run the Hadoop job and save the output into a new folder:

```
zhangxinrun@ubuntu:~$ hadoop jar /usr/local/hadoop/hadoop-streaming-2.9.1.jar -mapper /home/zhangxinrun/handson-2/task1/mapper.py -reducer /home/zhangxinrun/handson-2/task1/reducer.py -input /user/bigdata/handson-2/task1/pg20417.txt -output /user/bigdata/handson-2/task1/output
```

5. Result:

The Hadoop job is running

```
18/10/15 20:35:07 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1539659905077_0001
18/10/15 20:35:08 INFO impl.YarnClientImpl: Submitted application application_1539659905077_0001
18/10/15 20:35:08 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1539659905077_0001/
18/10/15 20:35:08 INFO mapreduce.Job: Running job: job_1539659905077_0001
18/10/15 20:35:25 INFO mapreduce.Job: Job job_1539659905077_0001 running in uber mode : false
18/10/15 20:35:25 INFO mapreduce.Job:  map 0% reduce 0%
18/10/15 20:35:47 INFO mapreduce.Job:  map 100% reduce 0%
18/10/15 20:35:59 INFO mapreduce.Job:  map 100% reduce 100%
18/10/15 20:36:02 INFO mapreduce.Job: Job job_1539659905077_0001 completed successfully
```

Take a look at Hadoop output folder and the outputs are showed:

```
zhangxinrun@ubuntu:~$ hdfs dfs -ls /user/bigdata/handson-2/task1/output

Found 2 items
-rw-r--r--  1 zhangxinrun supergroup          0 2018-10-15 20:35 /user/bigdata/handson-2/task1/output/_SUCCESS
-rw-r--r--  1 zhangxinrun supergroup       143 2018-10-15 20:35 /user/bigdata/handson-2/task1/output/part-00000
```

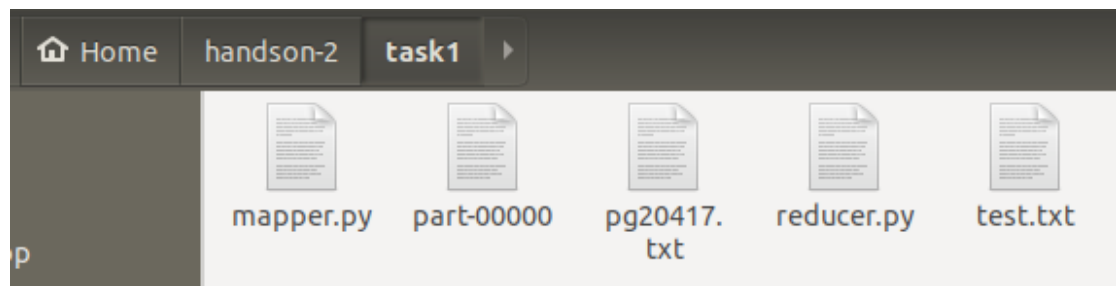
Use `-cat` command to execute output. The outputs are exactly same as test results:

```
zhangxinrun@ubuntu:~$ hdfs dfs -cat /user/bigdata/handson-2/task1/output/part-000000
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/usr/local/hadoop/share/hadoop/common/lib/hadoop-auth-2.9.1.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
('the', 9028)
('of', 5530)
('and', 2886)
('a', 2715)
('in', 2336)
('to', 2247)
('is', 2120)
('it', 1383)
('that', 1349)
('are', 947)
```

Download outputs into local file system.

```
zhangxinrun@ubuntu:~$ hdfs dfs -get /user/bigdata/handson-2/task1/output/part-000000 /home/zhangxinrun/handson-2/task1
```

the file in local file system, it's a .txt file:



## 6. Code:

Mapper.py:

```
1  #!/usr/bin/python3
2
3  import sys
4
5  # This function is used to handle double quote, single quote,
6  #or other non-alphabet character in the prefix or suffix.
7  #At first i tried to handle words
8  #however, i found that it can't handle some specific alphas like "--" and "://"
9  #def handle_word(word):
10     #word = word.replace('"','')
11     #word = word.replace("'",'')
12     #word = word.replace(',')
13     #x = word.find('--')
14     #if x != -1:
15         #newWord = word.split('--')
16         #for word in newWord:
17             #print('%s\t%s' % (word, 1))
18     #else:
19         #print('%s\t%s' % (word, 1))
20
21 #so, i turned to try to handle lines
22 #and it works well
23 def handle_line(line):
24     line = line.replace('"',' ')
25     line = line.replace("'",' ')
26     line = line.replace(',')
27     line = line.replace('://',' ')
28     line = line.replace(':', ' ')
29     line = line.replace('--',' ')
30     line = line.replace('-', ' ')
31     line = line.replace('@',' ')
32     line = line.replace('(',' ')
33     line = line.replace(')',' ')
34     line = line.replace('$',' ')
35     line = line.replace('","',' ')
36     line = line.replace('!', ' ')
37     line = line.replace(';', ' ')
38     line = line.replace('_', ' ')
39     line = line.replace('?', ' ')
40     line = line.replace("[", ' ')
41     line = line.replace("]", ' ')
42     line = line.replace("/", ' ')
43     line = line.replace("{", ' ')
44     line = line.replace("}", ' ')
45     line = line.replace("^", ' ')
46     return line
47
48 # input comes from STDIN (standard input)
49 for line in sys.stdin:
50     # remove leading and trailing whitespace
51     line = line.strip()
52     line = handle_line(line)
53     # split the line into words
54     words = line.split()
55     # increase counters
56     for word in words:
57         # write the results to STDOUT (standard output);
58         # what we output here will be the input for the
59         # Reduce step, i.e. the input for reducer.py
60         #
61         # tab-delimited; the trivial word count is 1
62         word = word.lower()
63         print('%s\t%s' % (word, 1))
64
```

I defined a `handle_line()` function to handle the alphas in each line. It replaced alphas with a blank. Then the words in each line were separated naturally and clearly.

Reducer.py:

```
1  #!/usr/bin/python3
2  """reducer.py"""
3
4  from operator import itemgetter
5  import sys
6
7  current_word = None
8  current_count = 0
9  word = None
10 newDict = {}
11
12 # input comes from STDIN
13 for line in sys.stdin:
14     # remove leading and trailing whitespace
15     line = line.strip()
16
17     # parse the input we got from mapper.py
18     word, count = line.split('\t', 1)
19
20     # convert count (currently a string) to int
21     try:
22         count = int(count)
23     except ValueError:
24         # count was not a number, so silently
25         # ignore/discard this line
26         continue
27
28     # by key (here: word) before it is passed to the reducer
29     if current_word == word:
30         current_count += count
31     else:
32         if current_word:
33             # write result to newDict{}
34             newDict[current_word] = current_count
35         current_count = count
36         current_word = word
37
38 # put the last word into newDict{}
39 if current_word == word:
40     newDict[current_word] = current_count
41
42 #sort the newDict{} and put the result into sorted_Dict{}
43 #output top 10 values and keys
44 sorted_Dict = sorted(newDict.items(),key = lambda item:item[1],reverse = 1)
45 m = 0
46 while m < 10:
47     print(sorted_Dict[m])
48     m += 1
```

To display top 10 frequency words, I saved the outputs of reducer in a dictionary newDict{} and sorted it. After that, I displayed the first ten keys and values of dictionary.

## Task2

1. Test mapper.py and reducer.py

```
zhangxinrun@ubuntu:~/handson-2/task2$ cat purchases.txt | ./mapper.py | sort | ./reducer.py
Albuquerque 149.8976161356743
Anaheim 146.42173162239143
Anchorage 140.5401075683234
Arlington 141.86862165937328
Atlanta 135.36741044462065
Aurora 147.24716460673767
Austin 136.57540721909703
Bakersfield 164.6138453204633
Baltimore 159.90480134406775
Baton Rouge 124.41992734361764
Birmingham 159.64722908498777
Boise 137.5190310123468
Boston 152.80345090576012
```

2. Make a new directory in Hadoop cluster:

```
zhangxinrun@ubuntu:~/handson-2/task2$ cd
zhangxinrun@ubuntu:~$ hdfs dfs -mkdir /user/bigdata/handson-2/task2
```

3. Copy the purchases.txt into Hadoop cluster:

```
zhangxinrun@ubuntu:~$ hdfs dfs -copyFromLocal /home/zhangxinrun/handson-2/task2/purchases.txt /user/bigdata/handson-2/task2
```

4. Run the Hadoop job and save the output into a new folder:

```
zhangxinrun@ubuntu:~$ hadoop jar /usr/local/hadoop/hadoop-streaming-2.9.1.jar -mapper /home/zhangxinrun/handson-2/task2/mapper.py -reducer /home/zhangxinrun/handson-2/task2/reducer.py -input /user/bigdata/handson-2/task2/purchases.txt -output /user/bigdata/handson-2/task2/output
```

5. Result:

```
18/10/29 13:21:56 INFO mapreduce.Job: Job job_1540842225502_0002 running in uber mode : false
18/10/29 13:21:56 INFO mapreduce.Job: map 0% reduce 0%
18/10/29 13:22:11 INFO mapreduce.Job: map 100% reduce 0%
18/10/29 13:22:20 INFO mapreduce.Job: map 100% reduce 100%
18/10/29 13:22:21 INFO mapreduce.Job: Job job_1540842225502_0002 completed successfully
18/10/29 13:22:22 INFO mapreduce.Job: Counters: 49
```

```
zhangxinrun@ubuntu:~$ hdfs dfs -ls /user/bigdata/handson-2/task2/output
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/usr/local/hadoop/share/hadoop/common/lib/hadoop-auth-2.9.1.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
Found 2 items
-rw-r--r-- 1 zhangxinrun supergroup 0 2018-10-29 13:22 /user/bigdata/handson-2/task2/output/_SUCCESS
-rw-r--r-- 1 zhangxinrun supergroup 2911 2018-10-29 13:22 /user/bigdata/handson-2/task2/output/part-00000
```

```
zhangxinrun@ubuntu:~$ hdfs dfs -cat /user/bigdata/handson-2/task2/output/part-00000
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.security.authentication.util.KerberosUtil (file:/usr/local/hadoop/share/hadoop/common/lib/hadoop-auth-2.9.1.jar) to method sun.security.krb5.Config.getInstance()
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.security.authentication.util.KerberosUtil
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
Albuquerque 142.63462259996976
Anaheim 144.69636714877205
Anchorage 139.47109810213223
Arlington 141.46898374025895
Atlanta 134.70971367024174
Aurora 141.16462827405482
Austin 130.57552668958868
Bakersfield 156.2696982218905
```

## 6. Code:

mapper.py:

```
mapper.py
~/handson-2/task2

#!/usr/bin/python3
import sys
for line in sys.stdin:
    line = line.strip()
    words = line.split('\t')
    store = words[2]
    sale = words[4]
    print('%s\t%s' % (store, sale))
```

reducer.py:

```
#!/usr/bin/python3
import sys
import math
delta1 = int(0)
delta2 = int(0)
M2 = int(0)
variance = int(0)
(current_store, tot_sales, count) = (None, 0.0, int(0))
for line in sys.stdin:
    (store, sale) = line.split('\t')
    if current_store and store != current_store:
        print('%s\t%s' % (current_store, SD))
        delta1 = int(0)
        delta2 = int(0)
        M2 = int(0)
        variance = int(0)
        (current_store, tot_sales) = (store, float(sale))
        count = 1
    else:
        mean = tot_sales / (count + 1)
        delta1 = float(sale) - mean
        mean = mean + delta1 / (count + 1)
        delta2 = float(sale) - mean
        M2 = M2 + delta1 * delta2
        variance = M2 / (count + 1)
        SD = variance ** 0.5
        (current_store, tot_sales, count) = (store, tot_sales + float(sale), count + 1)
        # print('%s\t%s' % (current_store, SD))
        # print('%s\t%s\t%s\t%s' % (current_store, tot_sales, SD, count))
if current_store:
    print('%s\t%s' % (current_store, SD))
```

### Task3

1. Test mapper.py and reducer.py

```
zhangxinrun@ubuntu:~/handson-2/task3$ cat test.txt | ./mapper.py | sort | ./reducer.py
148.224.17.21 1
160.69.35.16 2
198.31.73.63 5
204.255.215.49 7
gandalf.slip.cc.uq.oz.au 3
ix-al12-19.ix.netcom.com 1
ngriffin.itc.gu.edu.au 11
ns.okc.com 44
sydney-ts-38.nstn.ca 1
xslip48.csr.v.uidaho.edu 12
```

2. Make a new directory in Hadoop cluster

```
zhangxinrun@ubuntu:~/handson-2/task3$ hdfs dfs -mkdir /user/bigdata/handson-2/task3
```

3. Copy the epa-http.txt into the directory:

```
zhangxinrun@ubuntu:~/handson-2/task3$ cd
zhangxinrun@ubuntu:~$ hdfs dfs -copyFromLocal /home/zhangxinrun/handson-2/task3/epa-http.txt /user/bigdata/handson-2/task3
```

4. Run the Hadoop job and save the output into a new folder:

```
zhangxinrun@ubuntu:~$ hadoop jar /usr/local/hadoop/hadoop-streaming-2.9.1.jar -mapper /home/zhangxinrun/handson-2/task3/mapper.py -reducer /home/zhangxinrun/handson-2/task3/reducer.py -input /user/bigdata/handson-2/task3/epa-http.txt -output /user/bigdata/handson-2/task3/output
```

5. Result:

```
18/10/29 12:59:02 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1540842225502_0001/
18/10/29 12:59:02 INFO mapreduce.Job: Running job: job_1540842225502_0001
18/10/29 12:59:19 INFO mapreduce.Job: Job job_1540842225502_0001 running in uber mode : false
18/10/29 12:59:19 INFO mapreduce.Job: map 0% reduce 0%
18/10/29 12:59:39 INFO mapreduce.Job: map 100% reduce 0%
18/10/29 12:59:49 INFO mapreduce.Job: map 100% reduce 100%
18/10/29 12:59:51 INFO mapreduce.Job: Job job_1540842225502_0001 completed successfully
18/10/29 12:59:51 INFO mapreduce.Job: Counters: 49
```

```
Found 2 items
-rw-r--r-- 1 zhangxinrun supergroup 0 2018-10-29 12:59 /user/bigdata/handson-2/task3/output/_SUCCESS
-rw-r--r-- 1 zhangxinrun supergroup 50792 2018-10-29 12:59 /user/bigdata/handson-2/task3/output/part-00000
```

```
www.api.org 1
www.cfe.cornell.edu 6
www.hud.gov 2
wwwproxy.ac.il 6
wwwproxy.sanders.com 19
wwwproxy2.ca.sandia.gov 31
wyndmoor1-17.slip.netaxs.com 14
wzb145.wz-berlin.de 2
x121-8.dehs.umn.edu 16
x229-22.ppath.umn.edu 37
xs1.xs4all.nl 24
xslip48.csr.v.uidaho.edu 5
xyplex4-1-20.ucs.indiana.edu 22
y.nsf.gov 1
yaqut.ar.utexas.edu 9
yhqfm.yokogawa.co.jp 1
yukonho.cs.caltech.edu 8
yyj-ppp-13.cyberstore.ca 14
yyj-ppp-14.cyberstore.ca 4
yyj-ppp-8.cyberstore.ca 1
zeus.esy.com 86
zorba.wlctok1.epa.gov 3
zuni.r09.epa.gov 20
```

## 6. Code:

mapper.py:

```
mapper.py
~/handson-2/task3

#!/usr/bin/python3
import sys

x = 0
# input comes from STDIN (standard input)
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    attributes = line.split()
    #x = attributes[3].find('.gif')
    #y = attributes[3].find('.htm')
    #z = attributes[3].find('.xbm')
    #if x!= -1 or y!= -1 or z!= -1:
        #print('%s\t%s' % (attributes[0],attributes[3]))
    print('%s\t%s' % (attributes[0],attributes[3]))
```

reducer.py:

```
#!/usr/bin/python3
"""reducer.py"""

from operator import itemgetter
import sys

current_hostname = None
current_filename = None
hostname = None
filename = None
count = 1

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

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

    # convert count (currently a string) to int
    try:
        count = int(count)
    except ValueError:
        # count was not a number, so silently
        # ignore/discard this line
        continue

    # by key (here: word) before it is passed to the reducer
    if current_hostname == hostname and current_filename != filename:
        count += 1
    else:
        if current_hostname and current_hostname != hostname:
            # write result to ouput
            print('%s\t%s' % (current_hostname, count))
        count = 1
        current_hostname = hostname
        current_filename = filename

# put the last word into newDict{}
if current_hostname == hostname:
    print('%s\t%s' % (current_hostname, count))
```