# HW4

## CS 5665

Tutorial referred : http://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/ , also discussed with Bhagyashree about the approaches to follow for questions 2-c and 2-d.

## Overview

In this homework, you will write Map and Reduce functions to perform following two tasks:

**Task 1: Word Count**

A) Given the provided file (Tolstoy's War and Peace), create a complete count of each word that appears in the text. Which word appears the most?

Ans: -  Word "the" appears the most in the document warandpeace.txt:

**Command used to copy the file into hdfs:**  hdfs dfs -copyFromLocal warandpeace.txt

**Command used to run the mapper and reducer on Hadoop:**

hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoop-streaming-2.6.0-mr1-cdh5.8.0.jar -Dmapred.reduce.task=1 -file /home/cloudera/Desktop/mapper.py /home/cloudera/Desktop/reducer.py -mapper "python mapper.py" -reducer "python reducer.py" -input /user/cloudera/warandpeace.txt -output /user/cloudera/output1.a

Converted the output file generated(part-00000) to output1.a.txt for better visualization and readability.Output1.a:

```
 1    the 34721
 2    and 22287
 3    to  16749
 4    of  15001
 5    a   10599
 6    he  10000
 7    in  8960
 8    that    8195
 9    his 7983
10    was 7351
11    with    5709
12    it  5590
13    had 5365
14    her 4704
15    not 4692
16    him 4565
17    at  4547
18    i   4523
19    s   4410
20    but 4055
21    as  4032
22    on  4002
23    you 3858
24    for 3542
25    she 3484
26    is  3340
27    all 2797
```

B) Create a count of all the palindromes that occur in the text. Which palindrome occurs most often?

Ans :- The most occurring palindrome in the file is : "a".

Converted the output file generated(part-00000) to output1.b.txt for better visualization and readability. Output1.b:

```
1   a    10502
2   i    4078
3   did 1476
4   anna    293
5   '    197
6   )    90
7   eye 54
8   1    48
9   o    46
10  e    26
11  v    22
12  deed    21
13  iii 21
14  sees    21
15  x    21
16  ii   17
17  -    14
18  eve 14
19  m    14
20  s    14
21  f    12
22  level    12
23  xix 12
24  n    11
25  3    10
26  madam    10
27  b    9
28  xx   9
```

## Task 2: Election Fraud

In this task your job is to investigate whether there was election fraud in 2008. You have 2006 and 2008

election data files: (i) 2006 data file; and (ii) 2008 data file. The files are of the format where each line is a vote

in the election.

The format of the text file is:

VoterID \t CountyID \t PartyID

A) Which party won the election in 2008?

Ans:- **Party 3** won the election in 2008.

Converted the output file generated(part-00000) to output2.a.txt for better visualization and readability. Output2.a:
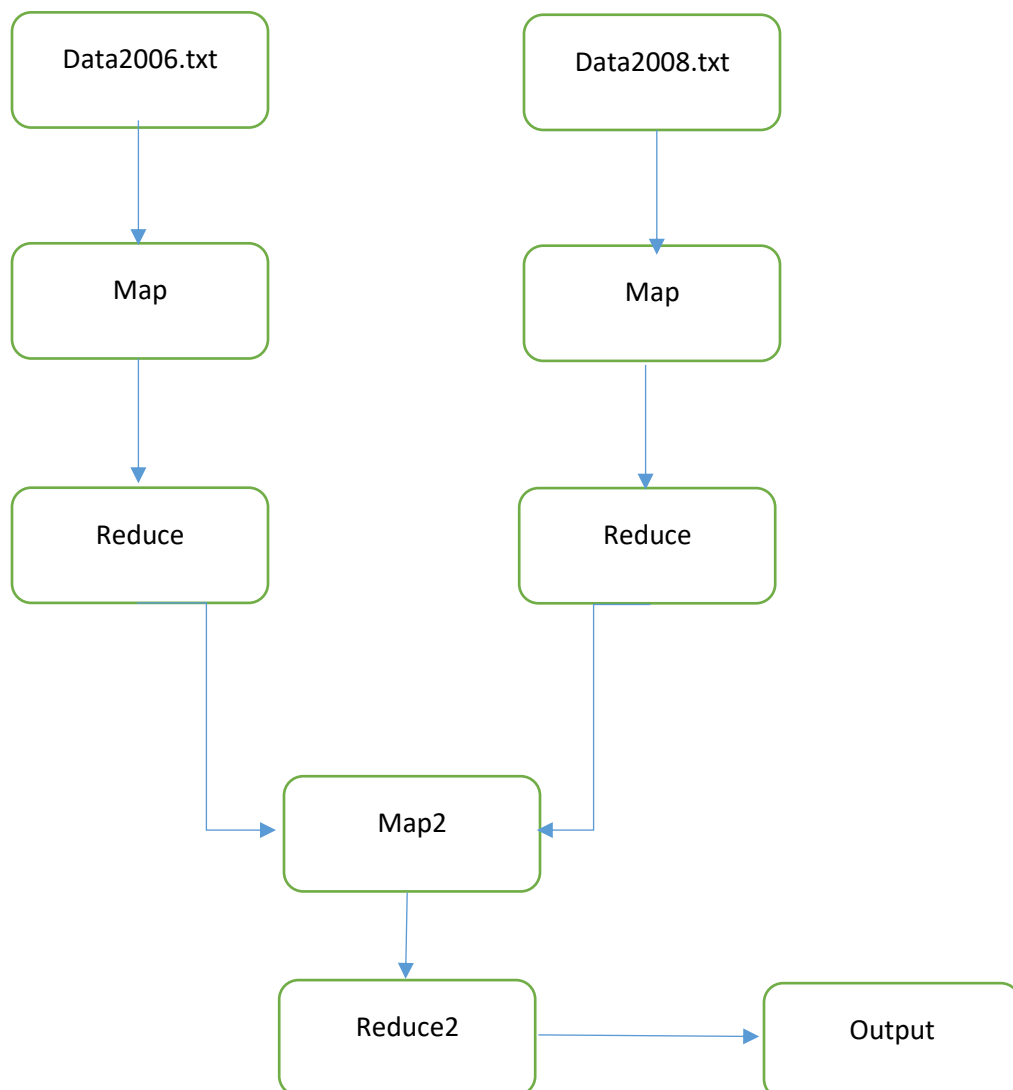
```
1   1    9408
2   2    10112
3   3    12071
```

B) In 2006, which county was the most monolithic in the manner in which they voted? (i.e. which county came

closest to voting 100% for a single party).

Ans :-  County # 277 was most monolithic in the manner they voted with party 3 getting 51.6% votes.

Hadoop architecture used to solve this problem:

```
  Data2006.txt              Data2008.txt
        |                         |
        v                         v
      Map                       Map
        |                         |
        v                         v
     Reduce                    Reduce
        |                         |
        +-----> Map2 <------------+
                  |
                  v
              Reduce2 ---------> Output
```

Converted the output file generated(part-00000) to output2.a.txt for better visualization and readability. Output2.a:

```
 1   100.000000   38.461538   24.615385   36.923077
 2   101.000000   24.242424   43.939394   31.818182
 3   102.000000   33.333333   33.333333   33.333333
 4   103.000000   32.926829   28.048780   39.024390
 5   104.000000   39.240506   29.113924   31.645570
 6   105.000000   37.837838   29.729730   32.432432
 7   106.000000   37.000000   28.000000   35.000000
 8   107.000000   35.616438   32.876712   31.506849
 9   108.000000   39.130435   28.985507   31.884058
10   109.000000   39.240506   24.050633   36.708861
11   110.000000   32.051282   30.769231   37.179487
12   111.000000   34.782609   27.173913   38.043478
13   112.000000   23.684211   34.210526   42.105263
14   113.000000   30.666667   37.333333   32.000000
15   114.000000   43.055556   22.222222   34.722222
16   115.000000   27.272727   36.363636   36.363636
17   116.000000   35.632184   37.931034   26.436782
```

.

.

.

.

```
178   277.000000   19.354839   29.032258   51.612903
```

C) Studies have shown if a political party gains more than 50% in voting percentage from one election cycle to

the next, then most likely fraud has occurred. (Example, if party A received 100 votes in 2006 in county B, then
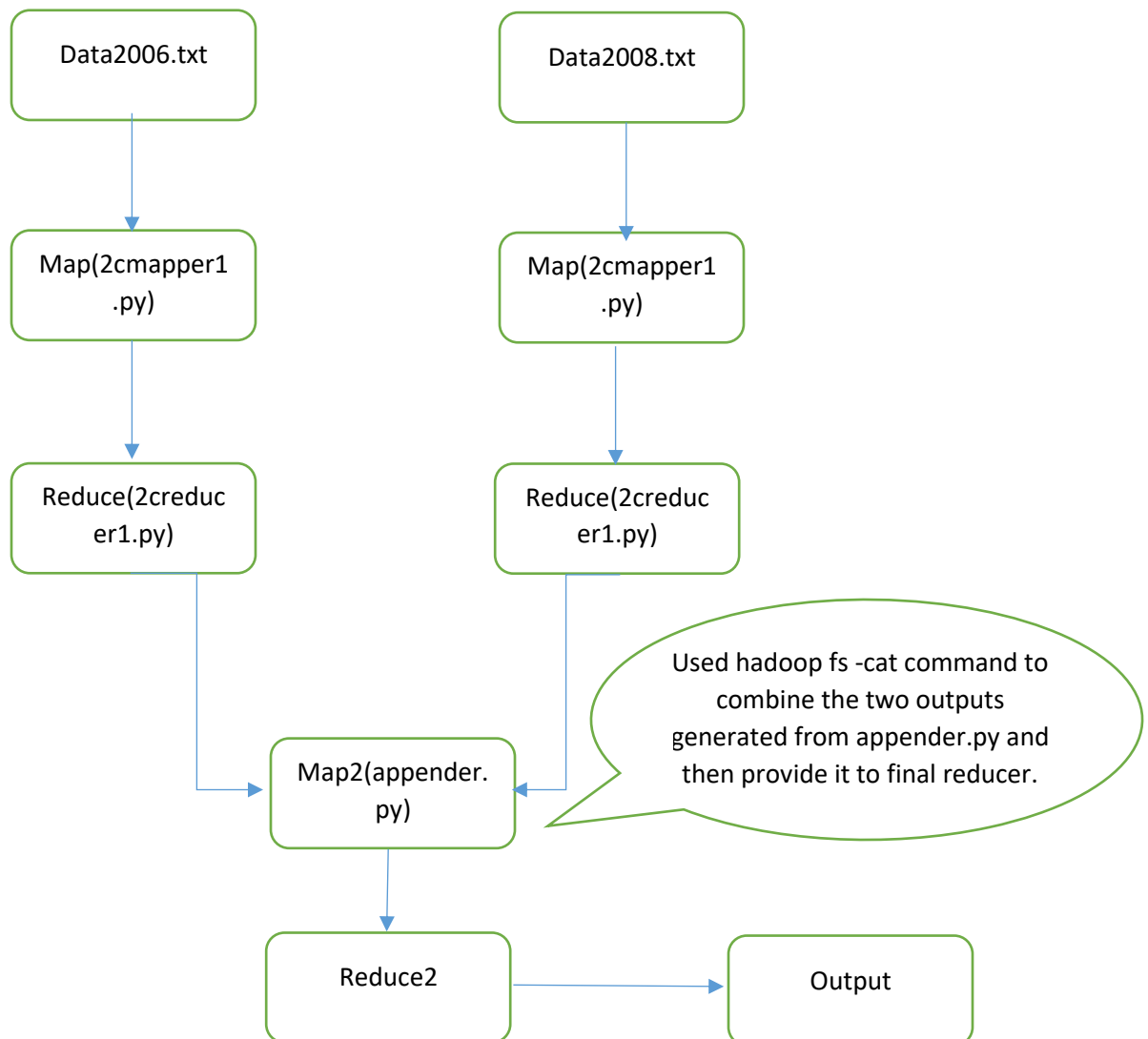
received 200 votes in 2008, fraud may have occurred). In which counties in 2008 did voter fraud likely occur?

Ans :- Counties in 2008 in which fraud occurred:

Converted the output file generated(part-00000) to output2.c.txt for better visualization and readability. Output2.c:

```
1    107 3    56.521739
2    178 3    72.222222
3    201 3    54.545455
4    220 3    60.000000
5    241 3    78.571429
6    244 3    60.000000
7    274 3    62.500000
8    332 3    65.217391
9    334 3    70.588235
10   359 3    52.173913
11   390 3    61.111111
12   424 3    52.631579
13   474 3    51.851852
```

Hadoop architecture used to solve this problem:

Data2006.txt

Data2008.txt

Map(2cmapper1.py)

Map(2cmapper1.py)

Reduce(2creducer1.py)

Reduce(2creducer1.py)

Map2(appender.py)

Used hadoop fs -cat command to combine the two outputs generated from appender.py and then provide it to final reducer.

Reduce2

Output

D) From 2006 to 2008 how many voters changed which party they voted for? What is the most common type of change?

Ans :- 6297 voters changed the party they voted for from 2006 to 2008, most common type of change was from party 1 to party 3. Following is the output:

Converted the output file generated(part-00000) to output2.d.txt for better visualization and readability. Output2.d:

```
1   party 1 to party 2  911
2   party 1 to party 3  1564
3   party 2 to party 1  668
4   party 2 to party 3  1563
5   party 3 to party 1  703
6   party 3 to party 2  888
```

Combined the two input files into one and then processed the file as a whole using single mapper and reducer.