# Report

I declare that the assignment submitted on Elearning system is original except for source material explicitly acknowledged, and that the same or related material has not been previously submitted for another course. I also acknowledge that I am aware of University policy and regulations on honesty in academic work, and of the disciplinary guidelines and procedures applicable to breaches of such policy and regulations, as contained in the website http://www.cuhk.edu.hk/policy/academichonesty/.


Signed (Student ___*Cao Yuhang*___ )         Date: 3/17/2018
Name: Cao Yuhang                             SID:   1155092180

|  | part (a) | part (b) | part (d) |
|---|---|---|---|
| Job1 time |  | 31s | 25 |
| Job2 time |  | 27s | 28 |
| Total time | 2m50s | 58s | 53s |

* For part (a) and part (d), #mappers is 20, #reducers is 4
* I run part (a) in my own mac book, my aws has no credit
* All time are fast running time for that experiment

Theoretically, running time should be d < b < a, since in part (d) we use PCY algorithm to further filter the false positive compared with part (b); and by using map reduce, parallelization should beat local machine.

According to the form, we can find running time d < b < a, it prove our guess.

All command can be found in "run.sh" in corresponding part directory. Here is the screenshot of "run.sh" for part (b):

```
10 hdfs dfs -rm -r output2
11
12 input='input'
13 s=0.005
14 m1=20
15 r1=4
16
17 output="output"
18
19 hadoop jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar \
20     -D mapred.map.tasks=$m1 \
21     -D mapred.reduce.tasks=$r1 \
22     -D mapred.compress.map.output=ture \
23     -file mapper1.py -mapper mapper1.py \
24     -file reducer1.py -reducer reducer1.py \
25     -cmdenv "s=${s}" \
26     -input $input \
27     -output output \
28
29 hdfs dfs -get output ./
30
31 hadoop jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar \
32     -D mapred.map.tasks=$m1 \
33     -D mapred.reduce.tasks=$r1 \
34     -D mapred.compress.map.output=ture \
35     -file mapper2.py -mapper mapper2.py \
36     -file reducer2.py -reducer reducer2.py \
37     -file output \
38     -cmdenv "s=${s}" \
39     -input $input \
40     -output output2 \
41
42 hdfs dfs -get output2 ./
43
44 cat output2/* > part_b_tmp
45 sort -k3nr part_b_tmp | head -n 40 > part_b_res
46 rm -rf part_b_tmp
47 rm -rf output
48 rm -rf output2
```

All the results can be found in corresponding part directory. For example, the result of part (b) can be found in "part_b/part_b_res.txt"

Q2

(a)

$$\begin{cases} 1-\left(1-s^r\right)^B \ge P1 & s \ge T1 \\ 1-\left(1-s^r\right)^B \ge P2 & s \le T2 \end{cases}$$

(b)

$$\begin{cases} 1-\left(1-s^r\right)^{20} \ge 0.99 & s \ge 0.9 \\ 1-\left(1-s^r\right)^{20} \le 0.05 & s \le 0.7 \end{cases}$$
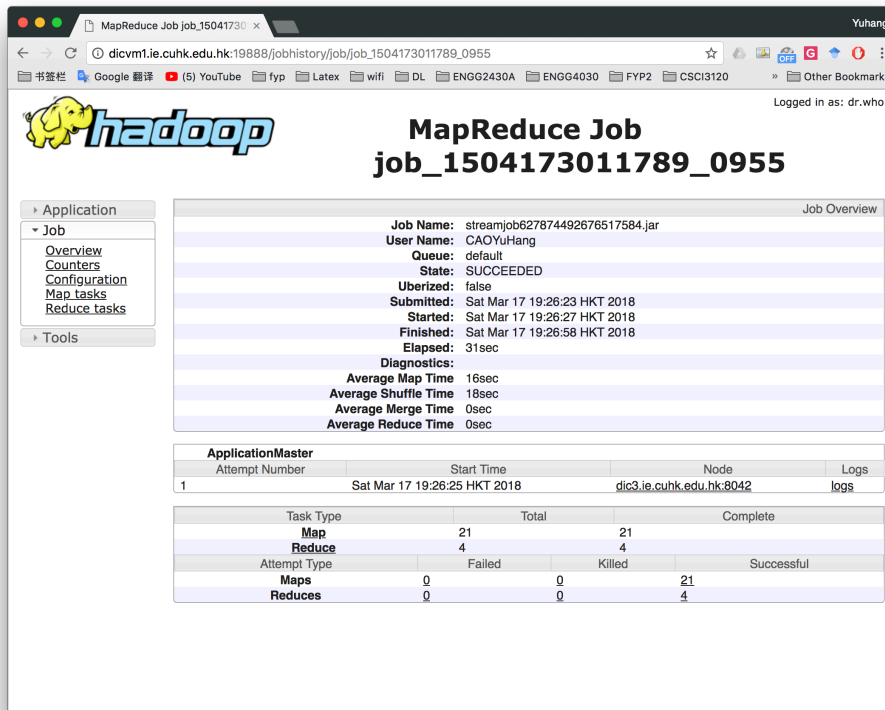
$$\begin{cases} r \ge \log_{0.9}\left(1-\left(1-0.99\right)^{\frac{1}{20}}\right)=15.01 \\ r \le \log_{0.7}\left(1-\left(1-0.05\right)^{\frac{1}{20}}\right)=16.73 \end{cases}$$

$$\begin{cases} B=20 \\ r=16 \end{cases}$$

Screen shot of running time of part (b) and part (d):

Part (b):

Job 1:



Job 2:

## Part (d)

### Job 1:



### Job2: