

# Cloud Computing HW2

## Use Hadoop with Docker Container

R06921105 葉津源

### Part 1: Set up Hadoop Docker environment

#### Steps:

1. create instance  
use mi.xlarge (need more than 8 GB RAM)
2. install docker  
use the following commands to install:

```
1  sudo apt-get remove docker docker-engine docker.io containerd runc
2  sudo apt-get update
3  sudo apt-get install \
4      apt-transport-https \
5      ca-certificates \
6      curl \
7      gnupg-agent \
8      software-properties-common
9  curl -fsSL https://download.docker.com/linux/ubuntu/gpg | sudo apt-key add -
10 sudo apt-key fingerprint 0EBFCD88
11 sudo add-apt-repository \
12     "deb [arch=amd64] https://download.docker.com/linux/ubuntu \
13     $(lsb_release -cs) \
14     stable"
15 sudo apt-get update
16 sudo apt-get install docker-ce docker-ce-cli containerd.io
17 sudo docker run hello-world
18
```

3. prepare files:

```
git clone https://github.com/sdwangntu/hadoop-cluster.git
cd hadoop-cluster

wget https://archive.apache.org/dist/hadoop/core/hadoop-3.1.2/hadoop-3.1.2.tar.gz
wget https://archive.apache.org/dist/hbase/1.4.9/hbase-1.4.9-bin.tar.gz
## scp to transfer hue-4.3.0.tgz
```

4. set up docker container with hadoop by following commands:

```
docker swarm init
docker network create --driver=overlay --attachable myattachable-network
docker build -t hadoop3cluster .

docker run --hostname=hadoop-master -p 8088:8088 -p 9870:9870 -p 9864:9864 \
    -p 19888:19888 -p 8042:8042 -p 8888:8888 --name hadoop-master \
    --network myattachable-network -d hadoop3cluster

docker run --hostname=hadoop-worker --name hadoop-worker --network myattachable-network -d hadoop3cluster
```

5. now docker container is set up and we can get into the docker containers to run hadoop programs by:

```
docker exec -it hadoop-master bash
docker exec -it hadoop-worker bash
```

## Part 2: Run Map Reduce in python for log analysis of apache2 web server access log

Steps: (collaborated with 陳冠甫)

1. write mapper.py

```
1  #!/usr/bin/python
2  import sys
3  from datetime import datetime
4  for line in sys.stdin:
5      line = line.strip()
6      if line == '':
7          continue
8      begin = line.find '['
9      end = line.find ']'
10     record = line[begin+1:end-6]
11     record = str(datetime.strptime(record,"%d/%b/%Y:%H:%M:%S"))
12     record = record[:len(record)-6]+' :00:00'
13     print(record+'\t1')
14
```

note that on line 6 and 7 we need to handle the case where there is an empty line, otherwise map reduce will fail.

2. write reducer.py

```
1  #!/usr/bin/python
2
3  from operator import itemgetter
4  import sys
5  current_word = None
6  current_count = 0
7  word = None
8  for line in sys.stdin:
9      line = line.strip()
10     word,count = line.split('\t',1)
11     try:
12         count = int(count)
13     except ValueError:
14         continue
15     if current_word == word:
16         current_count += count
17     else:
18         if current_word:
19             print("%s\t%s" % (current_word,current_count))
20             current_count = count
21             current_word = word
22         if current_word == word:
23             print("%s\t%s" % (current_word,current_count))
24
```

3. copy access.log into hadoop

do

hdfs dfs -mkdir /input

hdfs dfs -copyFromLocal access.log /input

4. run Map Reduce

do

hdfs dfs -rm -r /output

hadoop jar /opt/hadoop/share/hadoop/tools/lib/hadoop-streaming-3.1.2.jar -mapper "python \$PWD/mapper.py" -reducer "python \$PWD/reducer.py" -input "/input" -output "/output"

## Results:

2004-03-07 16:00:00	27
2004-03-07 17:00:00	25
2004-03-07 18:00:00	24
2004-03-07 19:00:00	26
2004-03-07 20:00:00	20
2004-03-07 21:00:00	23
2004-03-07 22:00:00	29
2004-03-07 23:00:00	22
2004-03-08 00:00:00	21
2004-03-08 01:00:00	21
2004-03-08 02:00:00	27
2004-03-08 03:00:00	22
2004-03-08 04:00:00	26
2004-03-08 05:00:00	37
2004-03-08 06:00:00	17
2004-03-08 07:00:00	31
2004-03-08 08:00:00	44
2004-03-08 09:00:00	63
2004-03-08 10:00:00	39
2004-03-08 11:00:00	34
2004-03-08 12:00:00	45
2004-03-08 13:00:00	37
2004-03-08 14:00:00	23
2004-03-08 15:00:00	9
2004-03-08 16:00:00	2
2004-03-08 17:00:00	2
2004-03-08 18:00:00	9
2004-03-08 19:00:00	6
2004-03-08 20:00:00	23
2004-03-08 22:00:00	20
2004-03-08 23:00:00	1
2004-03-09 01:00:00	12
2004-03-09 02:00:00	15
2004-03-09 03:00:00	1
2004-03-09 05:00:00	24
2004-03-09 06:00:00	29
2004-03-09 07:00:00	8
2004-03-09 08:00:00	27
2004-03-09 09:00:00	2

2004-03-09 10:00:00	11
2004-03-09 11:00:00	6
2004-03-09 12:00:00	9
2004-03-09 13:00:00	8
2004-03-09 14:00:00	14
2004-03-09 15:00:00	28
2004-03-09 16:00:00	2
2004-03-09 17:00:00	8
2004-03-09 18:00:00	12
2004-03-09 20:00:00	8
2004-03-09 21:00:00	3
2004-03-09 22:00:00	5
2004-03-09 23:00:00	1
2004-03-10 00:00:00	6
2004-03-10 02:00:00	5
2004-03-10 03:00:00	17
2004-03-10 05:00:00	1
2004-03-10 07:00:00	1
2004-03-10 08:00:00	38
2004-03-10 09:00:00	3
2004-03-10 10:00:00	6
2004-03-10 11:00:00	29
2004-03-10 12:00:00	102
2004-03-10 13:00:00	13
2004-03-10 14:00:00	3
2004-03-10 15:00:00	13
2004-03-10 16:00:00	6
2004-03-10 18:00:00	2
2004-03-10 19:00:00	2
2004-03-10 20:00:00	5
2004-03-10 21:00:00	12
2004-03-10 22:00:00	13
2004-03-10 23:00:00	9
2004-03-11 00:00:00	6
2004-03-11 01:00:00	1
2004-03-11 02:00:00	1
2004-03-11 03:00:00	6
2004-03-11 06:00:00	19
2004-03-11 07:00:00	6
2004-03-11 08:00:00	8
2004-03-11 10:00:00	3
2004-03-11 11:00:00	19

2004-03-11 12:00:00	17
2004-03-11 13:00:00	46
2004-03-11 14:00:00	15
2004-03-11 15:00:00	27
2004-03-11 16:00:00	6
2004-03-11 18:00:00	4
2004-03-11 20:00:00	14
2004-03-11 23:00:00	1
2004-03-12 01:00:00	1
2004-03-12 02:00:00	2
2004-03-12 03:00:00	1
2004-03-12 04:00:00	5
2004-03-12 05:00:00	11
2004-03-12 08:00:00	1
2004-03-12 09:00:00	2
2004-03-12 11:00:00	25
2004-03-12 12:00:00	22
2004-03-12 13:00:00	3