

## Cloud Computing Capstone Task II

Baoshi Sun, February 17, 2016

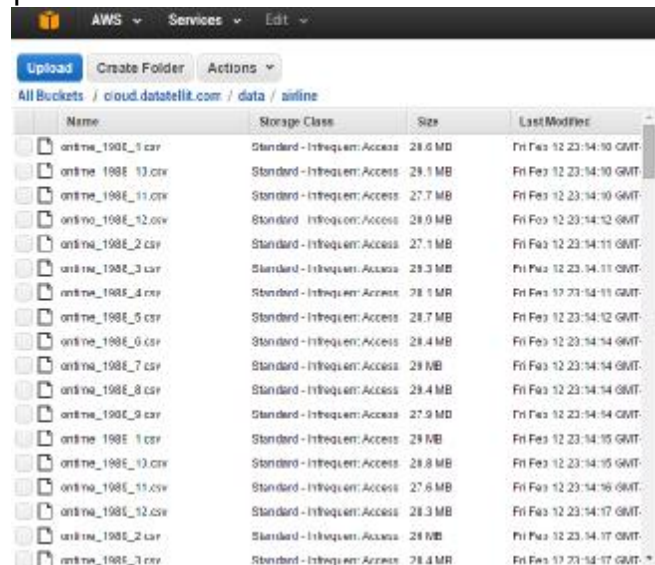
### Experiment Environment

Instead of using Hadoop & Cassandra on AWS EC2 IaaS like what I did in task I, for this task I turned to EMR and DynamoDB.

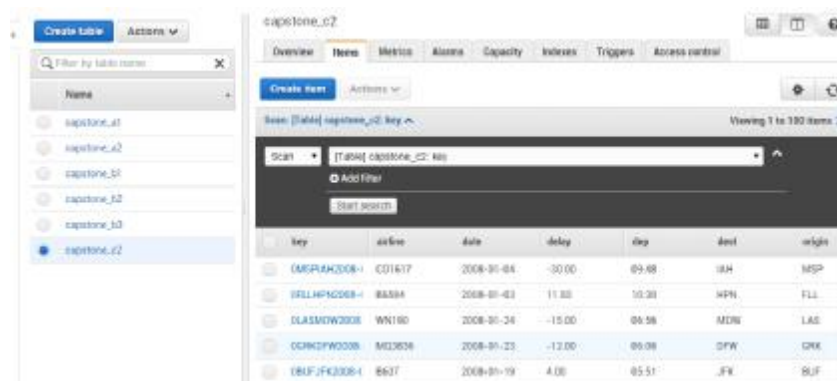
- AWS Resources
  - 3 m3.xlarge EC instances compose an EMR cluster
  - Volume: total 80G
  - S3 Storage
  - DynamoDB
- 3<sup>rd</sup> Party Components
  - Kafka 0.9.0.x
- Development Tools
  - Python 2.7 + pyspark
  - Boto 3

### Data Loading & Cleaning

The methods are almost the same as that in task I, expect that a few Hadoop file operation commands are replaced by AWS S3 CLI. The picture shows the cleaned csv files.



Name	Storage Class	Size	Last Modified
online_198E_1.csv	Standard - Infrequent Access	28.6 MB	Fri Feb 12 23:14:10 GMT
online_198E_13.csv	Standard - Infrequent Access	29.1 MB	Fri Feb 12 23:14:10 GMT
online_198E_11.csv	Standard - Infrequent Access	27.7 MB	Fri Feb 12 23:14:10 GMT
online_198E_12.csv	Standard - Infrequent Access	28.0 MB	Fri Feb 12 23:14:12 GMT
online_198E_2.csv	Standard - Infrequent Access	27.1 MB	Fri Feb 12 23:14:11 GMT
online_198E_3.csv	Standard - Infrequent Access	29.3 MB	Fri Feb 12 23:14:11 GMT
online_198E_4.csv	Standard - Infrequent Access	28.1 MB	Fri Feb 12 23:14:11 GMT
online_198E_5.csv	Standard - Infrequent Access	28.7 MB	Fri Feb 12 23:14:12 GMT
online_198E_6.csv	Standard - Infrequent Access	28.4 MB	Fri Feb 12 23:14:14 GMT
online_198E_7.csv	Standard - Infrequent Access	29.1 MB	Fri Feb 12 23:14:14 GMT
online_198E_8.csv	Standard - Infrequent Access	29.4 MB	Fri Feb 12 23:14:14 GMT
online_198E_9.csv	Standard - Infrequent Access	27.9 MB	Fri Feb 12 23:14:14 GMT
online_198E_1.csv	Standard - Infrequent Access	29.1 MB	Fri Feb 12 23:14:15 GMT
online_198E_10.csv	Standard - Infrequent Access	28.8 MB	Fri Feb 12 23:14:15 GMT
online_198E_11.csv	Standard - Infrequent Access	27.6 MB	Fri Feb 12 23:14:16 GMT
online_198E_12.csv	Standard - Infrequent Access	28.3 MB	Fri Feb 12 23:14:17 GMT
online_198E_2.csv	Standard - Infrequent Access	28.1 MB	Fri Feb 12 23:14:17 GMT
online_198E_3.csv	Standard - Infrequent Access	28.4 MB	Fri Feb 12 23:14:17 GMT

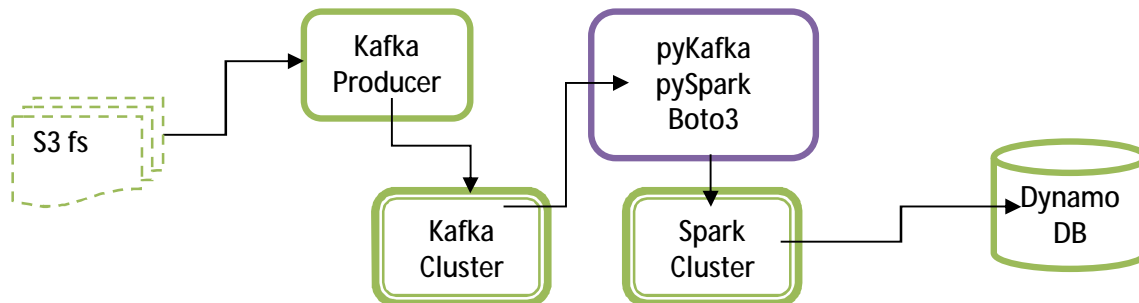


key	active	date	delay	log	dest	weight
CMGPAK008	CD1617	2008-01-04	-30.00	09:08	US	MSP
IRLLHPND08	BA584	2008-01-03	11.00	10:30	WPN	FLS
CLASHEW008	WN180	2008-01-24	-15.00	08:58	MEN	LAS
COHGPW008	MD3636	2008-01-23	-12.00	08:08	OPW	CRK
CRUF-FK2008-1	8637	2008-01-19	4.00	05:51	JFK	BUF

## Methodology

Although I tried S3 direct file streaming and Kafka consumer streaming approach, eventually I decided to use Kafka direct streaming, and it works very well.

Results of all calculations are stored into DynamoDB. As can be seen in above picture, one table corresponds to one question.



One single spark streaming program is deployed, which processes all 6 tasks (Q1.1, Q1.2, Q2.1, Q.2.2, Q2.3 and Q3.2) in a row, so that the data stream needs to be played only once.

I also wrote a small data feeding tool to generate stream to data consumers. The tool now supports reading data from S3, HDFS and local files.

As to the outputs, at first I wanted to write everything directly to DynamoDB via Boto interface. However, it turns out that the throughput of DynamoDB is a big problem. Even if I increased the write/read units, there was no effect. I guess the bottle neck could be the Boto interface. Elaborated program may help, but no time to do.

To balance the stream feeding speed and the program processing efficiency, a 10-second sleep is inserted between two feedings of data file. On the other hand, after test I set the streaming query interval to 6 seconds and check points is updated at every 60 seconds.

In addition, considering the program should execute for a long time but it can't stop even if there is no more data coming in, I added a simple self-check mechanism. If the number of total processed records stay unchanged for two minutes, the program will save necessary data and make a graceful exit.

Check the source code here:

- [Main spark streaming program](#)
- [Data feeding scripts](#)

## Results Report [\(Video link\)](#)

### Question 1.1 & Question 1.2

```

16/02/16 04:32:01 INFO DAGScheduler: ResultStage 12775 (runJob :
16/02/16 04:32:01 INFO DAGScheduler: Job 1746 finished: runJob :
Total: 233503980
ORD: 12449288
ATL: 11539676
DFW: 10799262
LAX: 7723452
PHX: 6585495
DEN: 6273780
DTW: 5636591
IAH: 5480672
MSP: 5199211
SFO: 5171014
16/02/16 04:32:01 INFO JobScheduler: Finished job streaming job
16/02/16 04:32:01 INFO JobScheduler: Total delay: 81.738 s for
16/02/16 13:08:51 INFO
HA: -1.01
AQ: 1.14
PS: 1.44
ML: 4.65
PA: 5.24
NW: 5.43
F9: 5.43
WN: 5.50
OO: 5.61
9E: 5.69
16/02/16 13:08:51 INFO

```

### Question 2.1

```

Result to question b1
CMI: (OH, 0.60), (US, 2.02), (TW, 4.12), (PI, 4.46), (DH, 6.02), (EV, 6.66), (MQ, 8.01)
Result to question b1
BWI: (F9, 0.76), (PA, 4.77), (CO, 5.18), (YV, 5.51), (NW, 5.71), (AA, 5.99), (9E, 7.24), (US, 7.49), (DL, 7.68), (UA, 7.74)
Result to question b1
MIA: (9E, -3.05), (EV, 1.20), (TZ, 1.78), (XE, 1.87), (PA, 4.20), (NW, 4.50), (US, 6.09), (UA, 6.87), (ML, 7.50), (FL, 8.56)
Result to question b1
LAX: (MQ, 2.41), (OO, 4.22), (FL, 4.73), (TZ, 4.76), (PS, 4.86), (NW, 5.12), (F9, 5.73), (HA, 5.81), (YV, 6.02), (US, 6.75)
Result to question b1
IAH: (NW, 3.56), (PA, 1), (3.98), (PI, 3.99), (US, 5.06), (F9, 5.55), (AA, 5.70), (TW, 6.05), (WN, 6.23), (OO, 6.59), (MQ, 6.71)
Result to question b1
SFO: (TZ, 3.96), (MQ, 4.83), (F9, 5.17), (PA, 5.28), (NW, 5.72), (PS, 6.31), (DL, 6.55), (CO, 7.06), (US, 7.52), (TW, 7.79)

```

### Question 2.2

```

Result to question b2
CMI: (ABI, -7.0), (PIT, 1.10), (CVG, 1.89), (DAY, 3.12), (STL, 3.98), (PIA, 4.59), (DFW, 5.94), (ATL, 6.67), (ORD, 8.19)
Result to question b2
BWI: (SAV, -6.98), (MLB, 1.16), (DAB, 1.47), (SRQ, 1.57), (IAD, 1.78), (UCA, 3.65), (CHO, 3.72), (GSP, 4.21), (SJU, 4.43), (OAJ, 4.47)
Result to question b2
MIA: (SHV, 0.00), (BUF, 1.00), (SAN, 1.71), (SLC, 2.51), (HOU, 2.91), (ISP, 3.65), (MEM, 3.75), (PSE, 3.98), (TLH, 4.26), (MCI, 4.61)
Result to question b2
LAX: (SDF, -16.01), (IDA, -6.98), (DRO, -5.99), (RSW, -3.01), (LAX, -2.02), (BZN, -0.73), (MAF, 0.00), (PIH, 0.00), (IYK, 1.27), (MFE, 1.38)
Result to question b2
IAH: (MSN, -2.0), (AGS, -0.62), (MLI, -0.5), (EFD, 1.89), (HOU, 2.17), (JAC, 2.57), (MTJ, 2.95), (RNO, 3.22), (BPT, 3.60), (VCT, 3.61)
Result to question b2
SFO: (SDF, -10.01), (MSO, -4.01), (PIH, -2.98), (LGA, -1.76), (PIE, -1.34), (OAK, -0.81), (FAR, 0.00), (BNA, 2.43), (MEM, 3.30), (SCK, 3.97)

```

### Question 2.3

```

[hadoop@ip-172-31-53-172 src]$ ./b3inp.sh
Result to question b3
CMI -> ORD: (MQ, 10.14)
Result to question b3
IND -> CMH: (CO, -2.55), (AA, 5.4), (HP, 5.72), (NW, 5.76), (US, 6.86), (DL, 10.69), (EA, 10.80)
Result to question b3
DFW -> IAH: (PA, -1.60), (EV, 5.09), (UA, 5.41), (CO, 6.49), (OO, 7.56), (XE, 8.09), (AA, 8.38), (DL, 8.60), (MQ, 9.10)
Result to question b3
LAX -> SFO: (TZ, -7.62), (PS, -2.15), (F9, -2.03), (EV, 6.96), (AA, 7.39), (MQ, 7.81), (US, 7.96), (WN, 8.79), (CO, 9.35), (NW, 9.85)
Result to question b3
JFK -> LAX: (UA, 3.31), (HP, 6.68), (AA, 6.90), (DL, 7.93), (PA, 11.02), (TW, 11.70)
Result to question b3
ATL -> PHX: (FL, 4.55), (US, 6.29), (HP, 8.48), (EA, 8.95), (DL, 9.81)
[hadoop@ip-172-31-53-172 src]$

```

### Question 3.2

```

Result to question c2
CMI -> ORD -> LAX, 04/03/2008: MQ4278,0710,-14,AA607,1952,-24
Result to question c2
JAX -> DFW -> CRP, 09/09/2008: AA845,0722,1,MQ3627,1648,-7
Result to question c2
SLC -> BFL -> LAX, 01/04/2008: 003755,1101,12,005429,1509,6
Result to question c2
LAX -> SFO -> PHX, 12/07/2008: WN3534,0650,-13,US412,1916,-19
Result to question c2
DFW -> ORD -> DFW, 10/06/2008: UA1104,0658,-21,AA2341,1650,-10
Result to question c2
LAX -> ORD -> JFK, 01/01/2008: UA944,0700,1,B6918,1853,-7

```

-- Query: CMI -> ORD -> LAX, 04/03/2008      -- Result: MQ4278,0710,-14,AA607,1952,-24

The result indicates the route is taking MQ4278 which departs at 07:10 on 04/03/2008 from CMI with 14 minutes earlier than the schedule, and taking AA607 at 19:52 on 06/03/2008 from ORD to LAX with 24 minutes earlier.

### Conclusions

1. By compared the results between task I and task II, I found there are slight differences. One possible reason is the algorithms are not identical. For example, in task II the 'cancel' flights are totally ignored, while they were counted as 'delay' in task I. However, I suspect some messages might be lost in streaming mode. If time allows, I'd like to investigate in detail.
2. What I learned from the project regarding the differences between MR and streaming include:
  - a. Streaming mode can perform many independent tasks upon one stream almost in one program space. But for MR, we need to launch different processes.
  - b. Spark requires much more server resources, especially memory, than MR. With the equivalent settings, Task I ran smoothly. However, during Task II I encountered numerous 'insufficient memory' errors and had to reboot the cluster again and again.
  - c. When dealing with streaming, it seems that more considerations should be put on optimization perspective.
  - d. Although EMR is easy to use, it is expensive and limited in many aspects. I would suggest my company to build its cloud platform from scratch (native Hadoop and Spark) on EC2, rather than use EMR.