

1. Give a brief overview of how you extracted and cleaned the data.

I launched an EC2 instance (t2.micro) with the EBS volume attached and then mounted:

<http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ebs-attaching-volume.html>
<http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ebs-using-volumes.html>

At the time data were available, I dig into the information and realized only files under “/aviation/airline_ontime” were required to complete all the question along the different groups. The fields imported were: (Year, Month, DayofMonth, Weekday, UniqueCarrier, FlightNum, Origin, Dest, CRSDepTime, DepDelay, CRSArtime, ArrDelay, Cancelled).

A python script was developed to extract, clean and import data to S3 bucket called “airlineontime”, different python modules was used to complete every step:

- The zipfile module was used to manipulate ZIP archive files.
- The os module was used to manipulate path and read all the lines in all the files.
- The csv module was used to manipulate all the files in csv format.
- The Boto module is the AWS SDK for Python, which allows write software that makes use of Amazon services like S3 and EC2.

2. Give a brief overview of how you integrated each system.

EC2 volume was loaded into S3 bucket called “airlineontime”. Public and secret key was provided to EC2 instance with a proper policy providing rights to load data into S3 bucket. Data from S3 bucket was available for HUE and by means of Hive SQL-like query language could be easily manipulated. S3 bucket was stored in HDFS, using Hive, in logical structures which can be easily queried by means HQL. Data from HDFS can be loaded into Dynamo DB tables by means of Hive queries:

```
CREATE EXTERNAL TABLE airline_ontime (year INT, month INT, day INT, weekday INT,  
carrier STRING, flight_num STRING, origin STRING, dest STRING, depdelay INT, depdelay  
INT, arrtime STRING, arrdelay INT, cancelled INT)  
PARTITIONED BY (date string)  
ROW FORMAT DELIMITED FIELDS TERMINATED BY ","  
LOCATION 's3n://airline-ontime/';  
MSCK REPAIR TABLE airline_ontime;
```

3. What approaches and algorithms did you use to answer each question?

Group 1, 2, 3: Hive SQL-like query language on top of Hadoop HDFS file system.

Group 1 ex 2 query example:

```
select carrier, sum(arrdelay)/count(arrdelay) as mean_delay from airline_ontime  
where cancelled = 0 group by carrier order by mean_delay asc limit 10;
```

Group 2, 3: Create Dynamo DB tables and insert data from Hive queries result.

Group 2 ex 4 query example:

```
insert overwrite table group2_ex4  
select origin, dest as destination, sum(arrdelay)/count(arrdelay) as mean_delay from  
airline_ontime where cancelled = 0 group by origin, dest;
```

4. What are the results of each question?

Group 1 (Answer any 2):

1. Rank the top 10 most popular airports by numbers of flights to/from the airport.

<i>airport</i>	<i>total</i>
ORD	12051796
ATL	11323515
DFW	10591818
LAX	7586304
PHX	6505078
DEN	6183518
DTW	5504120
IAH	5416653
MSP	5087036
SFO	5062339

2. Rank the top 10 airlines by on-time arrival performance.

<i>carrier</i>	<i>mean_delay</i>
HA	-1.01180434575
AQ	1.15692344248
PS	1.45063851278
ML (1)	4.74760919573
PA (1)	5.32243099993
F9	5.46588114882
NW	5.55778339267
WN	5.56077425988
OO	5.73631246366
9E	5.8671846617

3. Rank the days of the week by on-time arrival performance.

<i>weekday</i>	<i>mean_delay</i>
6	4.30166992608
2	5.99045884132
7	6.61328029244
1	6.71610280259
3	7.20365639467
4	9.09444100834
5	9.72103233759

Group 2 (Answer any 3):

1. For each airport X, rank the top-10 carriers in decreasing order of on-time departure performance from X.

1. CMI (University of Illinois Willard Airport)

<i>Airport</i>	<i>carrier</i>	<i>mean_delay</i>
CMI	OH	0
CMI	US	2
CMI	PI	4
CMI	TW	4

<i>CMI</i>	<i>DH</i>	<i>6</i>
<i>CMI</i>	<i>EV</i>	<i>6</i>
<i>CMI</i>	<i>MQ</i>	<i>8</i>

2. BWI (Baltimore-Washington International Airport)

<i>airport</i>	<i>carrier</i>	<i>mean_delay</i>
<i>BWI</i>	<i>F9</i>	<i>0</i>
<i>BWI</i>	<i>PA (1)</i>	<i>4</i>
<i>BWI</i>	<i>CO</i>	<i>5</i>
<i>BWI</i>	<i>NW</i>	<i>5</i>
<i>BWI</i>	<i>YV</i>	<i>5</i>
<i>BWI</i>	<i>AA</i>	<i>6</i>
<i>BWI</i>	<i>9E</i>	<i>7</i>
<i>BWI</i>	<i>DL</i>	<i>7</i>
<i>BWI</i>	<i>FL</i>	<i>7</i>
<i>BWI</i>	<i>UA</i>	<i>7</i>

3. MIA (Miami International Airport)

<i>airport</i>	<i>carrier</i>	<i>mean_delay</i>
<i>MIA</i>	<i>9E</i>	<i>-3</i>
<i>MIA</i>	<i>EV</i>	<i>1</i>
<i>MIA</i>	<i>TZ</i>	<i>1</i>
<i>MIA</i>	<i>XE</i>	<i>1</i>
<i>MIA</i>	<i>NW</i>	<i>4</i>
<i>MIA</i>	<i>PA (1)</i>	<i>4</i>
<i>MIA</i>	<i>UA</i>	<i>6</i>
<i>MIA</i>	<i>US</i>	<i>6</i>
<i>MIA</i>	<i>ML (1)</i>	<i>7</i>
<i>MIA</i>	<i>FL</i>	<i>8</i>

4. LAX (Los Angeles International Airport)

<i>airport</i>	<i>carrier</i>	<i>mean_delay</i>
<i>LAX</i>	<i>MQ</i>	<i>2</i>
<i>LAX</i>	<i>FL</i>	<i>4</i>
<i>LAX</i>	<i>OO</i>	<i>4</i>
<i>LAX</i>	<i>PS</i>	<i>4</i>
<i>LAX</i>	<i>TZ</i>	<i>4</i>
<i>LAX</i>	<i>F9</i>	<i>5</i>
<i>LAX</i>	<i>HA</i>	<i>5</i>
<i>LAX</i>	<i>NW</i>	<i>5</i>
<i>LAX</i>	<i>US</i>	<i>6</i>
<i>LAX</i>	<i>YV</i>	<i>6</i>

5. IAH (George Bush Intercontinental Airport)

<i>airport</i>	<i>carrier</i>	<i>mean_delay</i>
<i>IAH</i>	<i>NW</i>	<i>3</i>
<i>IAH</i>	<i>PA (1)</i>	<i>3</i>
<i>IAH</i>	<i>PI</i>	<i>3</i>

<i>IAH</i>	<i>AA</i>	<i>5</i>
<i>IAH</i>	<i>F9</i>	<i>5</i>
<i>IAH</i>	<i>US</i>	<i>5</i>
<i>IAH</i>	<i>HP</i>	<i>6</i>
<i>IAH</i>	<i>MQ</i>	<i>6</i>
<i>IAH</i>	<i>OO</i>	<i>6</i>
<i>IAH</i>	<i>TW</i>	<i>6</i>

6. SFO (San Francisco International Airport)

<i>airport</i>	<i>carrier</i>	<i>mean_delay</i>
<i>SFO</i>	<i>TZ</i>	<i>3</i>
<i>SFO</i>	<i>MQ</i>	<i>4</i>
<i>SFO</i>	<i>F9</i>	<i>5</i>
<i>SFO</i>	<i>NW</i>	<i>5</i>
<i>SFO</i>	<i>PA (1)</i>	<i>5</i>
<i>SFO</i>	<i>DL</i>	<i>6</i>
<i>SFO</i>	<i>PS</i>	<i>6</i>
<i>SFO</i>	<i>AA</i>	<i>7</i>
<i>SFO</i>	<i>CO</i>	<i>7</i>
<i>SFO</i>	<i>TW</i>	<i>7</i>

- For each airport X, rank the top-10 airports in decreasing order of on-time departure performance from X.

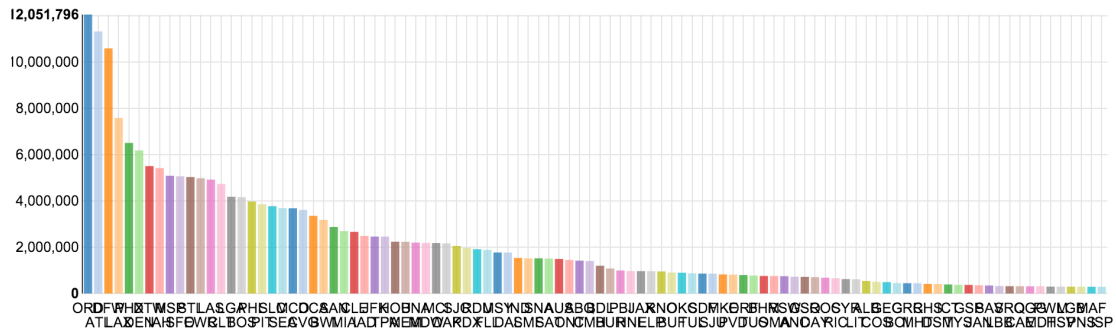
<i>airport</i>	<i>List of top-10 airports</i>
<i>CMI</i>	<i>ABI, PIT, CVG, DAY, STL, PIA, DFW, ATL, ORD</i>
<i>BWI</i>	<i>SAV, SRQ, DAB, IAD, MLB, UCA, CHO, DCA, IAH, OAJ</i>
<i>MIA</i>	<i>SHV, BUF, SAN, SLC, HOU, ISP, MEM, PSE, GNV, TLH</i>
<i>LAX</i>	<i>GRR, AZO, MSP, DTW, DAY, PIT, CVG, CLE, IAD, ATL</i>
<i>IAH</i>	<i>MSN, MLI, AGS, EFD, JAC, HOU, MTJ, VCT, RNO, BPT</i>
<i>SFO</i>	<i>SDF, MSO, PIH, LGA, PIE, FAR, OAK, BNA, MEM, SCK</i>

- For each source-destination pair X-Y, rank the top-10 carriers in decreasing order of on-time arrival performance at Y from X.
- For each source-destination pair X-Y, determine the mean arrival delay (in minutes) for a flight from X to Y.

<i>origin</i>	<i>destination</i>	<i>mean_delay</i>
<i>CMI</i>	<i>ORD</i>	<i>10</i>
<i>IND</i>	<i>CMH</i>	<i>2</i>
<i>DFW</i>	<i>IAH</i>	<i>7</i>
<i>LAX</i>	<i>SFO</i>	<i>9</i>
<i>JFK</i>	<i>LAX</i>	<i>6</i>
<i>ATL</i>	<i>PHX</i>	<i>9</i>

Group 3 (Answer both using only Hadoop and Spark):

- Does the popularity distribution of airports follow a Zipf distribution? If not, what distribution does it follow?



Do not follow a Zipf distribution, on the other hand it looks like a logarithmic distribution.

- Find, for each X-Y-Z and day/month combination in the year 2008, the two flights (X-Y and Y-Z) that satisfy constraints, if such flights exist.

<i>route</i>	<i>deptime</i>	<i>flight_xy</i>	<i>flight_yz</i>
CMI-ORD-LAX	2008-03-04	MQ4278	AA1345
CMI-ORD-LAX	2008-03-04	MQ4401	AA1345
JAX-DFW-CRP	2008-09-09	AA845	MQ3627
SLC-BFL-LAX	2008-04-01	OO3755	OO5429
LAX-SFO-PHX	2008-07-12	WN3534	US412
DFW-ORD-DFW	2008-06-10	UA1104	OO6119
LAX-ORD-JFK	2008-01-01	UA944	B6918

5. What system- or application-level optimizations (if any) did you employ?

DynamoDB sort keys were properly selected in order to obtain results faster. S3 buckets provide persistent storage which reduce wastes populating HDFS every time the Hadoop cluster is stopped.

6. Your opinion about whether the results make sense and are useful in any way.

Conclusion about size/popularity of different airport was quite predictable. The information obtained regarding different types of mean delay were not predictable at all and are quite useful from different points of views: carriers could take actions to reduce delays just in case these are not acceptable and passengers could adapt their preferences not only about carriers but also regarding flights' scheduling to avoid delays.

VIDEO DEMONSTRATION:

<https://sendvid.com/ngczdubk>