# INFO 7250 Final Project Flight On-time Performance Analysis

YUNAN SHAO
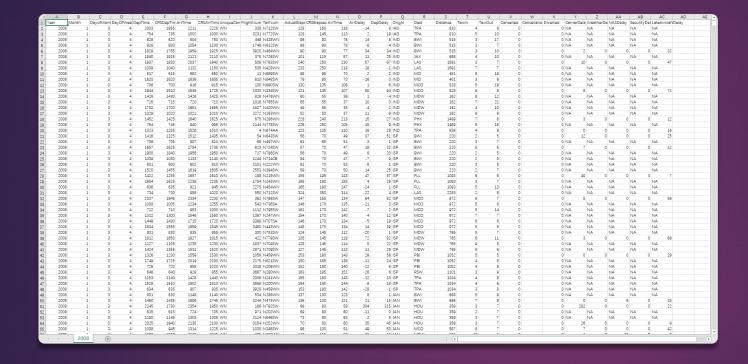
# Introduction

▶ Air transport is now very common for travelling. Although, it is the fastest method of transport, passengers waste both time and money when their flights are delayed.

▶ This project use the flight dataset to analyze the flight on-time performance and try to implement a machine learning model to predict if a flight will be delayed.

▶ Possible Questions:

   ▶ During which time period does flight delay the most?

   ▶ Relations between Airports' air traffic and flight delay.

   ▶ Airliner on-time performance.

# Objective & Techniques

**1. Analyze the flight on-time performance from different aspects (time, airport and carriers) and provide visualized results**

- MultipleOutput Binning
- MapReduce, Job Chaining
- MultipleInput Reducer side Join
- TopK Pattern
- Combiner, Memory-Conscious Implementation
- Custom Writable for processing selected attributes

**2. Try to implement machine learning model for prediction flight delay**

- Apache Mahout
- AWS Machine Learning

# Dataset Flight Data From 1987 to 2008

http://stat-computing.org/dataexpo/2009/the-data.html

# Field Attributes

**Variable descriptions**

| | Name | Description |
|---|---|---|
| 1 | Year | 1987-2008 |
| 2 | Month | 1-12 |
| 3 | DayofMonth | 1-31 |
| 4 | DayOfWeek | 1 (Monday) - 7 (Sunday) |
| 5 | DepTime | actual departure time (local, hhmm) |
| 6 | CRSDepTime | scheduled departure time (local, hhmm) |
| 7 | ArrTime | actual arrival time (local, hhmm) |
| 8 | CRSArrTime | scheduled arrival time (local, hhmm) |
| 9 | UniqueCarrier | unique carrier code |
| 10 | FlightNum | flight number |
| 11 | TailNum | plane tail number |
| 12 | ActualElapsedTime | in minutes |
| 13 | CRSElapsedTime | in minutes |
| 14 | AirTime | in minutes |
| 15 | ArrDelay | arrival delay, in minutes |
| 16 | DepDelay | departure delay, in minutes |
| 17 | Origin | origin IATA airport code |
| 18 | Dest | destination IATA airport code |
| 19 | Distance | in miles |
| 20 | TaxiIn | taxi in time, in minutes |
| 21 | TaxiOut | taxi out time in minutes |
| 22 | Cancelled | was the flight cancelled? |
| 23 | CancellationCode | reason for cancellation (A = carrier, B = weather, C = NAS, D = security) |
| 24 | Diverted | 1 = yes, 0 = no |
| 25 | CarrierDelay | in minutes |
| 26 | WeatherDelay | in minutes |
| 27 | NASDelay | in minutes |
| 28 | SecurityDelay | in minutes |
| 29 | LateAircraftDelay | in minutes |

Most of analyses are using files from 2003-2008 because these fields are 'NA' before 2003

# Type of Delay

▶ **Carrier Delay**

Carrier delay is within the control of the air carrier. Examples of occurrences that may determine carrier delay are: aircraft cleaning, aircraft damage, awaiting the arrival of connecting passengers or crew, baggage, bird strike, cargo loading, catering, computer, outage-carrier equipment, crew legality (pilot or attendant rest), damage by hazardous goods, engineering inspection, fueling, handling disabled passengers,   late crew, lavatory servicing, maintenance, oversales, potable water servicing, removal of unruly passenger, slow boarding or seating, stowing        carry-on baggage, weight and balance delays.

▶ **Late Arrival Delay**

Arrival delay at an airport due to the late arrival of the same aircraft at a previous airport. The ripple effect of an earlier delay at downstream airports is referred to as delay propagation.

▶ **NAS Delay**

Delay that is within the control of the National Airspace System (NAS) may include: non-extreme weather conditions, airport operations, heavy traffic volume, air traffic control, etc. Delays that occur after Actual Gate Out are usually attributed to the NAS and are also reported through OPSNET.

▶ **Security Delay**

Security delay is caused by evacuation of a terminal or concourse, re-boarding of aircraft because of security breach, inoperative screening equipment and/or long lines in excess of 29 minutes at screening areas.

▶ **Weather Delay**

Weather delay is caused by extreme or hazardous weather conditions that are forecasted or manifest themselves on point of departure,    enroute, or on point of arrival.

# Supplemental Data

▶ Airport and Carrier tables for JOIN in MapReduce Jobs

▶ * The Federal Aviation Administration (FAA) considers a flight to be delayed when it is 15 minutes later than its scheduled time.

**Airports**

airports.csv describes the locations of US airports, with the fields:

- iata: the international airport abbreviation code
- name of the airport
- city and country in which airport is located.
- lat and long: the latitude and longitude of the airport

This majority of this data comes from the FAA, but a few extra airports (mainly military bases and US protectorates) were collected from other web sources by Ryan Hafen and Hadley Wickham.

**Carrier codes**

Listing of carrier codes with full names: carriers.csv

# Inspection on the Dataset

▶ Because the dataset is too large, even one csv file takes too long to load. In order to inspect the data, I use MultipleOutput for binning. (mapper and reducer side)
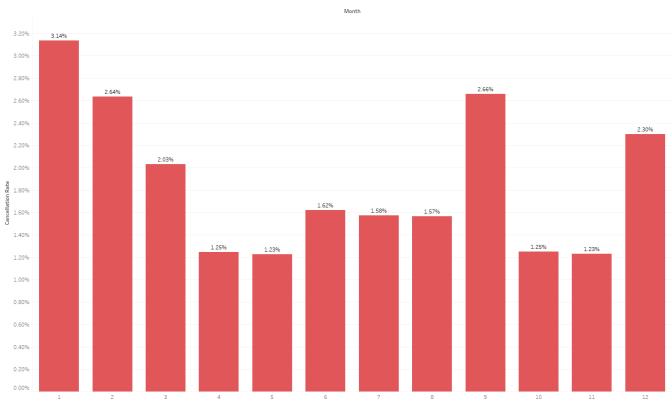
# Causes of Delay
# 2003-2008



**Main Causes**

1. Carrier Delay

2. NAS Delay (National Airspace System)

3. Late Aircraft Delay

*Detailed results for each year can be accessed
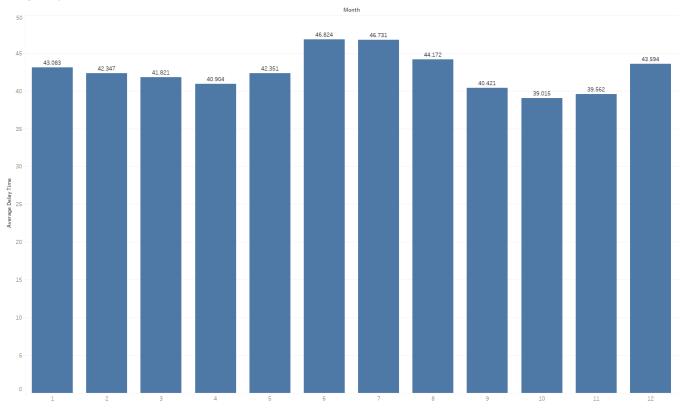
in Tableau file.

Cancellation Rate By Month

Delay Rate By Month

Average Delay Time By Time Range

# Time Analysis Summary & Conclusion

- ▶ January, September and December have the highest cancellation rates December, January, and June have the highest delay rates. Possible causes are high traffic volume during winter holidays and Christmas from Dec to Jan. Summer holidays starts in June and ends in Aug.

- ▶ Friday has the highest delay rate.

- ▶ Evening period (especially 6pm – 8pm) has the highest delay rate and relatively high cancellation rate.

- ▶ Not enough information to know why night especially midnight has higher cancellation rate. The cancellation rates are 'NA' even in most of the records in 03-08 files.

- ▶ *Detailed results for each year can be accessed in Tableau file.

# Top 20 Airports 1987-2008
## Number of Flights per Year

▶ Dynamic chart. See d3js folder in codes



**William B Hartsfield-Atlanta Intl**

| Airport | Value |
|---|---|
| William B Hartsfield-Atlanta Intl | 414513 |
| Chicago O'Hare International | 350380 |
| Dallas-Fort Worth International | 281281 |
| Denver Intl | 241443 |
| Los Angeles International | 215608 |
| Phoenix Sky Harbor International | 199408 |
| George Bush Intercontinental | 185172 |
| McCarran International | 172876 |
| Detroit Metropolitan-Wayne County | 161989 |
| San Francisco International | 140587 |
| Salt Lake City Intl | 139088 |
| Newark Intl | 138506 |
| Orlando International | 130872 |
| Minneapolis-St Paul Intl | 130289 |
| Charlotte/Douglas International | 126045 |
| LaGuardia | 119135 |
| John F Kennedy Intl | 118804 |
| Gen Edw L Logan Intl | 117915 |
| Seattle-Tacoma Intl | 109069 |
| more-Washington International | 104074 |

2008

Top 10 Origin Airports Delay Rate

# Top 10 Destination Airports Delay Rate

Delay Rate

Airport / Day of Week
ORD

Cancellation Rate

Airport / Day of Week
ORD

# Airport Time Example
# Chicago O'Hare

## Airport Analysis Conclusion

- ▶ Atlanta, Chicago and Dallas have top total number of flights through all time.

- ▶ Newark Intl is the worst airport to travel as both origin and destination, it has the highest delay rate and average delay time.

- ▶ Chicago, Atlanta, Detroit have relatively high delay rate as origin

- ▶ As a destination, San Francisco has almost the same delay rate as Newark

- ▶ If consider the airports with high number of delayed flights and delayed rate both as busy, most of them have also higher total number of flights

- ▶ *Detailed results for each year can be accessed in Tableau file.

# Total Number of Flights
By Carriers 1987-2008

▶ Dynamic chart. See d3js folder in codes

Cancellation Rate

Average Delay Rate by Carrier

Delay Rate

Delay Rate Line Chart 2003-2008

## Carrier Analysis Conclusion

- ▶ The carrier performance analysis only counts the records with carrier delay as the main reason of the delay

- ▶ Pinnacle and Mesa have highest average cancellation rates

- ▶ Atlantic Southeast and Alaska have highest average delay rates

- ▶ Airlines that delay rates keep raising: American Airlines, Continental, ExpressJet

- ▶ Alaska, Us Airways, Atlantic Southeast, US Airways, SkyWest and Northwest have big improvements in recent one year

- ▶ Aloha Airline is a small carrier and started service at 2006 but has lowest average delay rate and still improving

- ▶ *Detailed results for each year can be accessed in Tableau file.

# Travel Advice

- Passengers should avoid peak time period such as the evening. Avoid travelling in busy month is not really possible if they have a travel plan but they can still try to book flights in the morning, afternoon or early night.

- Another possible way is to rearrange the day you fly, for example flying on Saturday will be better than Friday.

- If your origin or destination airport is a big one, consider again to avoid peak time. (ORD example shows Saturday is much more better even for a busy airport)

- Booking flights operated by carriers have better on-time performance and avoid bad ones such as Atlantic SouthWest.

# Machine Learning (Attempt)
# Data Preparation

Use 2008.csv as sample data. Similar patterns for both implementations. MapReduce to get training data then train the model with selected fields.

## Fields Selected

- 1. Month
- 2. DayofWeek
- 3. CRSDepTime (Scheduled Departure Time) - Hour
- 4. Carrier
- 5. Origin
- 6. Dest
- *7. Delayed(0-not delayed, 1-delayed) arrdelay>=15 minutes

```java
            }
        return result;
    }

    public OnlineLogisticRegression train(List<Observation> trainData) {
        // System.out.println(trainData.size());
        System.out.println("Start Training");
        OnlineLogisticRegression olr = new OnlineLogisticRegression(2, 7, new L1());
        // Train the model
        for (int pass = 0; pass < 5; pass++) {
            for (Observation observation : trainData) {
                olr.train(observation.getActual(), observation.getVector());
            }

            if (pass % 1 == 0) {
                Auc eval = new Auc(0.5);
                for (Observation observation : trainData) {
                    eval.add(observation.getActual(), olr.classifyScalar(observation.getVector()));
                }

                System.out.format("Pass: %2d, Accuracy: %2.4f\n", pass + 1, eval.auc());
            }
        }
        return olr;
    }

    void testModel(OnlineLogisticRegression olr) {
        Observation newObservation = new Observation(new String[] { "12", "5", "19", "EV", "LAS", "PHX", "0" });
        Vector result = olr.classifyFull(newObservation.getVector());

        System.out.println("------------- Testing -------------");
        System.out.format("Probability of not Delay (0) = %.3f\n", result.get(0));
        System.out.format("Probability of Delay (1)     = %.3f\n", result.get(1));
    }

    class Observation {
        private DenseVector vector = new DenseVector(7);
        private int actual;
```

Console:

```
<terminated> LogisticRegression [Java Application] /usr/local/lib/jdk1.8.0_192/bin/java (Dec 12, 2018, 12:30:47 AM)
Start Importing
Start Training
Pass:  1, Accuracy: 0.5495
Pass:  2, Accuracy: 0.5481
Pass:  3, Accuracy: 0.5456
Pass:  4, Accuracy: 0.5477
Pass:  5, Accuracy: 0.5511
------------- Testing -------------
Probability of not Delay (0) = 0.717
Probability of Delay (1)      = 0.283
```

# Apache Mahout
# OnlineLogisticRegression

## Try real-time predictions

You submitted 6 out of 6 data values for this prediction.

Try generating real-time predictions for free using the web browser on this page. To request a real-time prediction, complete the following form or provide a single data record in CSV format. To provide a data record, choose the **Paste a record** button. [ Paste a record ]

**Prediction results**

| | | |
|---|---|---|
| Target name | y | |
| ML model type | BINARY | |
| Predicted label | 0 | |

| | Name | Type | Value |
|---|---|---|---|
| 1 | month | Categorical | 12 |
| 2 | dayofweek | Categorical | 5 |
| 3 | hour | Categorical | 19 |
| 4 | carrier | Categorical | EV |
| 5 | origin | Categorical | LAS |
| 6 | dest | Categorical | PHX |
| 7 | y | Binary | Target |

Items per page: 10 ▾    « ‹ 1 - 7 of 7 › »

« ‹ 1 - 7 of 7 › »

Clear data    [ Create prediction ]

```
{
  "Prediction": {
    "details": {
      "Algorithm": "SGD",
      "PredictiveModelType": "BINARY"
    },
    "predictedLabel": "0",
    "predictedScores": {
      "0": 0.43838316202163696
    }
  }
}
```

## ML model summary

| | |
|---|---|
| ID | ml-VUHidVYnpT1 |
| Name | ML model: flights ✏ |
| Type | Binary classification |
| Creation time | Dec 9, 2018 3:21:31 PM |
| Completion time | 11 mins. ⓘ |
| Compute Time (Approximate) | 9 mins. ⓘ |
| Status | Completed |
| Log | Download log |

## Datasource (training)

| | |
|---|---|
| Datasource ID | ds-kUXsd0OhMGW |
| Target | y |
| Input schema | View input schema |

## Evaluations

| | |
|---|---|
| Evaluations created | 1 |
| Latest evaluation result | 0.673 (AUC) |

[ Perform another Evaluation ]

# AWS Machine Learning

# Machine Learning Conclusion

▶ The accuracy of my prediction model is not good enough (best 60% and average 55%). One of the reason could be that Apache Mahout is normally used for recommendation and classification model. It's very hard to find a detailed guide showing how to implement the logistic regression and tune the model.

▶ The AWS Machine Learning Model has a better score (67%) but still not good enough. I have tried to used MapReduce to get training data with different parameters but this is the only one data that I have successfully imported to the model due to the permission problem of S3 bucket.

# Additional Ideas

▶ Selecting datasets with such big time range is good for showing results through years such as performance trending. However, making recommendations with average/overall result may not be a good idea since airports and carriers operate differently because of the increased air traffic and more advanced technologies they use.

▶ For prediction using machine learning, the idea is similar to the first one. My opinion is to use most recent data (2-3 years) to increase the performance of the prediction model. Another possible reason that my models don't get ideal scores is I don't have the weather data. There are weather cancellations and delays in the data and the weather data could be related to the month results in some seasons such as Summer and Winter (Storm, Snow).

# References

- ▶ Apache Mahout OnlineLogisticRegression example

  http://technobium.com/logistic-regression-using-apache-mahout/

- ▶ D3.js

  https://github.com/d3/d3

  https://github.com/Jannchie/Historical-ranking-data-visualization-based-on-d3.js