

# Data Mining: Final Project Write-up

*Misclassification-Terminator(Yuxi Chang, Mengran He, Alexander Lam, Xiaowen Yin)*

*May 10, 2017*

## Introduction

Every year, about 20% flights gets delayed or canceled across U.S., which could cause tremendous loss in money and time not only for travellers but also for airports or flight carriers. Pittsburgh as one of the main international airports in the United States, takes an important role in the transportation network. In this report, our goal is to predict whether a flight departing from Pittsburgh International Airport would be delayed at least 15 minutes using known flight information, along with weather information of both Pittsburgh International Airports and the destination airports.

## Data

There are three major data sources used for our prediction:

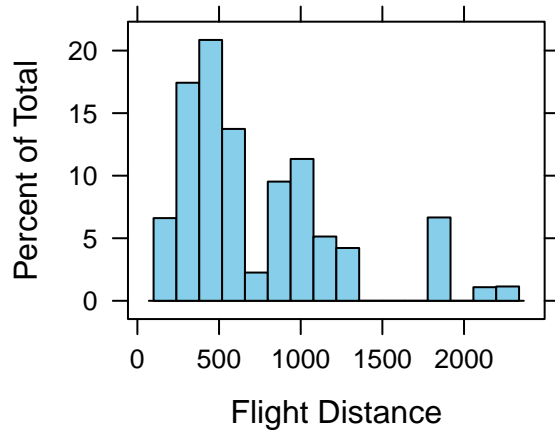
1. An extract from the Airline On-Time Performance Data made from the Bureau of Transportation Statistics of the U.S. Department of Transportation which reflects commercial flight activity to and from Pittsburgh International Airport.
2. Hourly weather data scraped from <http://www.wunderground.com>, which was merged into the Airline On-Time Performance Data based on the CRS departure time. The closest Pittsburgh hourly weather record pervious to the CRS departure time was used. All missing values for hourly data were imputed with the previous hour data.
3. Daily weather data for the arrival airports was achieved from <https://www7.ncdc.noaa.gov/>. This was also merged into the first dataset based on the flight date and arrival destination. All missing values for daily data were impyted with the previous day data.

## Exploratory Data Analysis

Our training dataset are flights and weather data from the year 2015 and there are 25870 records in our dataset in total, with 36 features. Among these features, 17 are features about flight activities, including departure time in different scales(hour, day, quarter, month, holiday), flight destination, flight carrier, CRS time, same hour flights in total, same day flights in total and whether the actual arrival time of the previous flight is already behind the CRS departure time of the next same flight. The features left are temperature features both at Pittsburgh International Airport and at the destination airport, including temperature, humidity, pressure, visibility, wind speed, snow, rain, fog, hail, thunder, or tornado conditions. Some of these features like snow, rain, fog, hail, thunder, and tornado were factorized into levels according to their occurance severities.

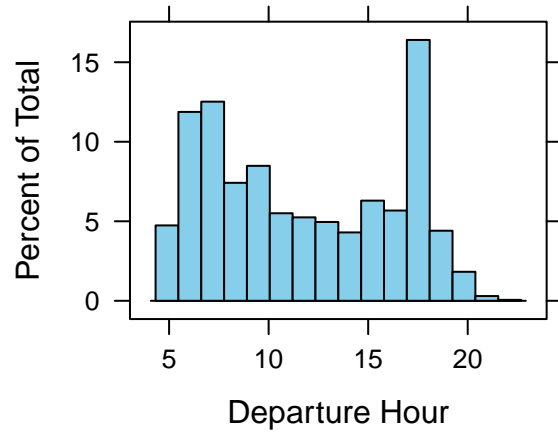
## Flight Activities

### Flight Distance Distribution



**Figure 1**

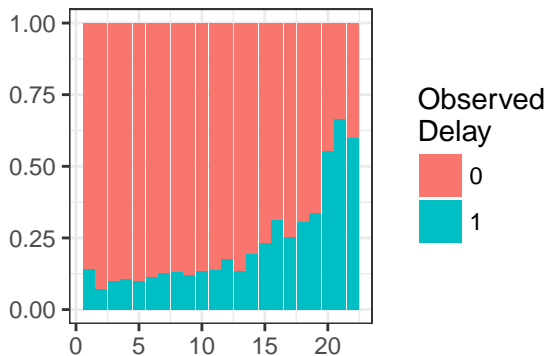
### Departure Time Distribution



**Figure 2**

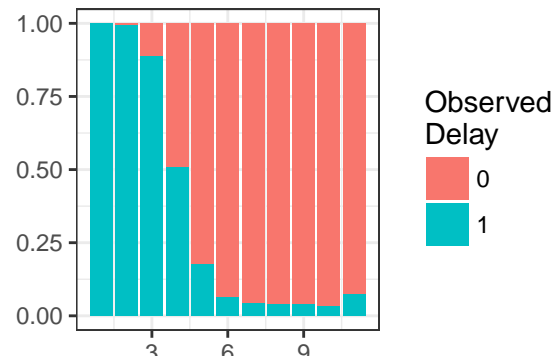
Figure 1 and Figure 2 above show the flight distance and departure time distribution across the year 2015. Most of the flight distance departing from Pittsburgh were in a range of 0-1000 miles. Morning (5-8am) is the busiest time for Pittsburgh International Airport during a day, and 4-7pm is the second busiest time. There is no departure flight after 11pm. In addition, weekdays are busier than weekends, and Saturdays have the fewest departure flights.

### Proportion of Delayed Departure Flights



**Figure 3**

### Proportion of Delayed Departure Flights



**Figure 4**

Figure 3 above shows the proportion of delayed departure flights versus the total number of departure and arrival flight at the same hour. From the plot, the proportion of PIT departure delayed flights increase with the total number of flights departing and arriving at the PIT.

Based on the plane registered tail number (TAIL\_NUM), we was able to identify if a planned departure plane was just arrived from another airport to PIT or not. We calculated time between the plane actual arrival plus taxiing time (ARR\_TIME + TAXI\_IN) and the plane planned departure time (CRS\_DEP\_TIME). If a scheduled departure plane wasn't from another city, time left for plane departure preparation was not applicable. In order to still use this variable, we discretized this time into 12 levels, with 1 meaning less than

5 minutes left, and 12 meaning more than 90 minutes left. Figure 4 shows the proportion of delayed departure flights versus the actual time left for departure preparation after the plane arrived at PIT.

### Feature of Departure Airport (Pittsburgh International Airport)

After some exploration in the weather information we scraped, we found that Pittsburgh is more humid from Dec to January and from June to August. We also found that the temperature follows a clear seasonal pattern (cold during winter months, hot during summer months). The visibility is very low during January through March. Wind speed is high during winter months. Heavy snows may occur from January to March. Visibility decreases along with the increase in rain and snow severity, so does pressure.

### Feature of Destination Airport

After some exploration in the weather information we scraped, we found that certain airports with more rains and snows have relatively low visibility and high wind speed, such airports include but limited to: Lambert St Louis International Airport, Boston Logan International Airport, Atlanta International Airport, Charlotte International Airport.

### Delay vs. Non-delay

The base rate of our dataset (i.e. proportion of delayed flights) is 13.77%. An interesting finding by comparing delay and non-delay flights is that there are more delay flights during evenings than in the mornings even though mornings are normally busier, as shown in figure 5 below. Moreover, temperatures at destination airports for non-delayed flights are much higher than those for delayed flights, as shown in figure 7.

### Delay vs. Not Delay: Departure Time vs. Not Delay: Arrival Temperature

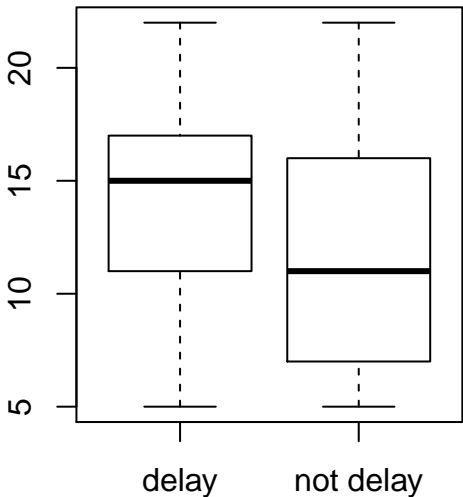


figure 5

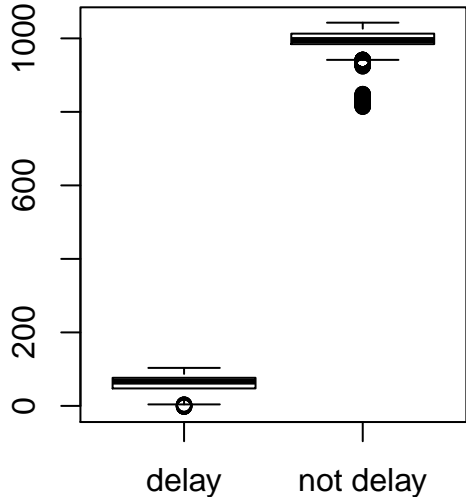


figure 6

## Unsupervised Learning: Clustering

Clustrings using k-means, k-medoids, complete linkage, and minimax linkage were examined. All methods gave 2 clusters, with K-means and K-medoids presenting similar results while complete linkage and minimax linkage giving slightly different results. The clustering results given by minimax linkage were the closest to the results we got from previous exploratory data analysis. As Figure 7 below shows, the average temperature of arrival destinations are higher for cluster 2 compared to that of cluster 1, which agrees with what we got from the previous section (figure 6).

## Supervised Learning

### Lasso logistic Regression

Lasso logistic regression as one of the most popular modern regularized regression models, could carry out the variable selection process automatically by shrinking the unnecessary variables to zero. We constructed model matrix including all features and fitted lasso logistic regression using the dataset of the year 2015 as our training dataset. To be more precise, we decided to use the minimum lambda instead of the 1SE lambda. If a simpler model is expected in the future, we could switch to the 1SE lambda. In this model, 88 of all 104 feature combinations were selected. For prediction, we did not think the threshold of 0.5 was good enough so that we designed a function which allows for best threshold selection by minimizing the test error with visible 2016 flights data as our testing dataset. As a result, a threshold of 0.45 was selected and the misclassification rate using the selected threshold was 6.7%.

```
bestThresh = function(y_hats, y_test) {  
  seqTry = seq(0, 1, 0.001)  
  
  current_misclass = 2  
  current_thresh = NA  
  for (ii in seqTry) {  
    y_hat01 = factor(ifelse(y_hats >= ii, 1, 0), levels = c(0, 1))  
    misclass = mean(y_hat01 != y_test)  
    thresh = ii  
    if (misclass < current_misclass) {  
      current_misclass = misclass  
      current_thresh = thresh  
    }  
  }  
  return(c(thresh = current_thresh, misclass = current_misclass))  
}
```

### Random Forests

Random forests are another robust method for classification, and is operated by constructing a group of decision trees. As an improvement of “bagging” strategy, random forests could make more independent trees by only allowing a small subset of variables to be considered at each split. We tried 500 trees to grow at first and then increased to 1000 trees for better classification. Again, we used the function designed for best threshold selection. As a result, a threshold of 0.304 was selected and the misclassification rate using this threshold was 5.6%, which was better than Lasso Logistic Regression. As you can see from the following variance importance plot (Figure 8) below:

(plot pending)

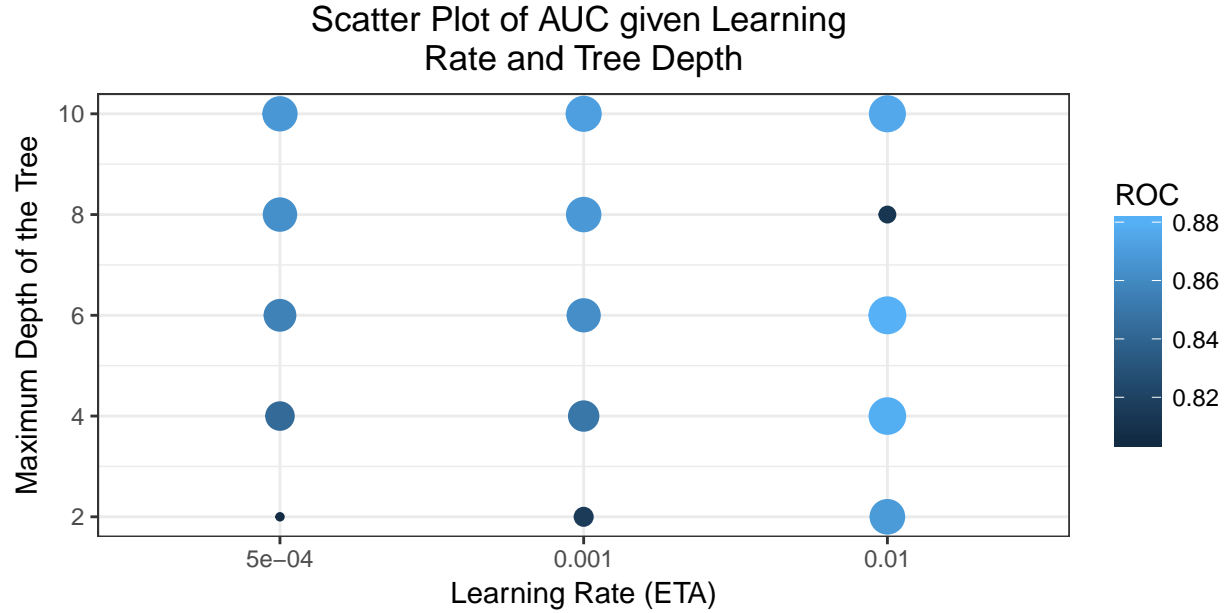


Figure 9

Figure 1: Scatter Plot of AUC given Learning Rate and Tree Depth on Training set

## K-Nearest-Neighbour (KNN)

K-Nearest-Neighbors (KNN) is a model-free classification method. It predicts an input according to its nearest-neighbors and then classifies according to a majority vote. Because KNN method heavily depends on distances between data points, we transferred all the categorical variables into numeric variables, except the flight carrier (CARRIER) and the flight destination (DEST). Also, to avoid variables with large variation determine the distances between points, we standardized all numeric variables before fitting the data in the KNN model. In order to choose the number of neighbors, parameter  $k$ , we ran KNN models with 1 to 40 nearest-neighbors, and  $k=13$  gives the lowest misclassification rates when comparing predictions to 16 test data.

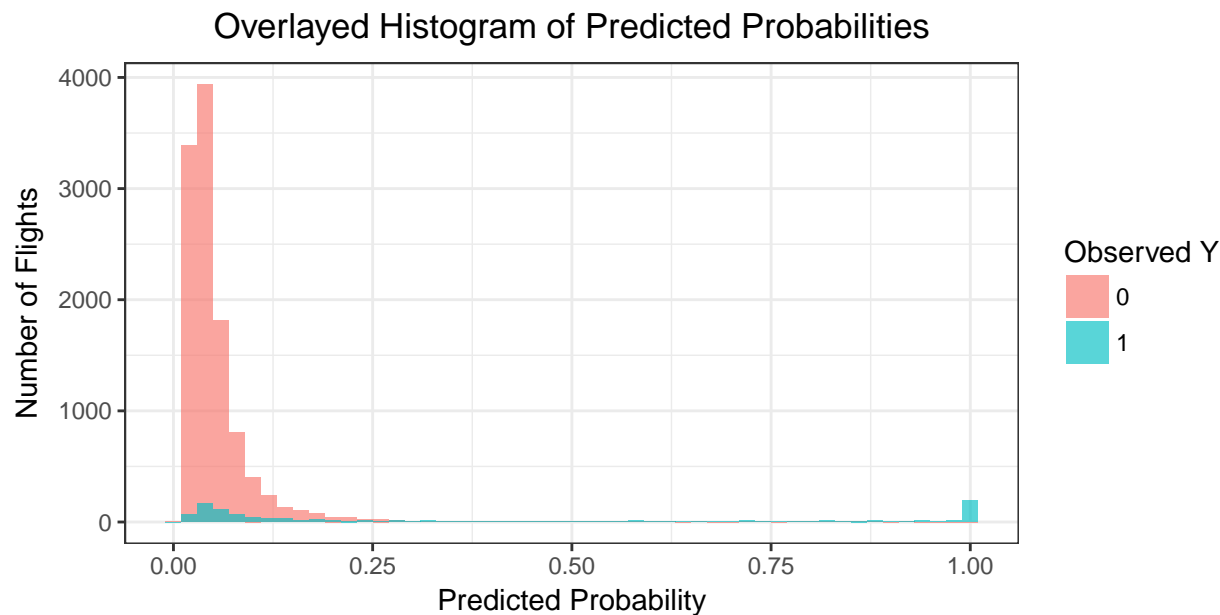
## Boosting

Extreme Gradient Boosting(xgboost) is an efficient algorithm that can be applied in supervised learning because it has high predictive power and fast computation. We decided to apply xgboost on training set (flight activities in 2015) and tried to find the best model by tuning booster parameters with 5 folds cross-validation. For the booster parameters, we tuned on learning rate(**eta**), **gamma**, maximum depth of the tree (**max\_depth**), and subsample of training instance (**subsample**) and built 1000 trees in the model (**nrounds**). The learning rate specifies step size shrinkage used in update to prevents overfitting, the smaller **eta** the more conservative the algorithm will be. **gamma** represents the minimum loss reduction required to make a further partition on a leaf node of the tree, so the larger the gamma more conservative the algorithm is.

```
xgb_grid1 = expand.grid(nrounds = 1000, eta = c(0.01, 0.001, 5e-04), max_depth = seq(2,
  10, 2), gamma = c(0.5, 1), colsample_bytree = 1, min_child_weight = 1, subsample = c(0.5,
  1))
```

After 5 folds of cross validation, we have the model with largest AUC (0.8806) value with **max\_dempth** = 6, **ets** = 0.01, **gamma** = 1, **subsample** = 0.5. Figure 9 shows the results of AUC given the values of learning rate and tree depth, the higher color and bigger size of the dot suggests higher AUC value.

Then we used the selected model to predict the flight delays on testing set (visible data in 2016). The figure 10 shows the histogram of the probability of departure delay given the actual flight delay information with two colors, where 0 means actual flight didn't delay, 1 means the flight delayed. In order to classify whether the flight delays or not, we tuned the cutoff value of probability of departure delay based on the missclassification rate, in such a way that we get the minimized missclassification rate 0.056 with cutoff value 31.2%. It means that if the probability of delay is greater than 31.2%, the flight is predicted to be delay.



Given the actual flight information in testing set, we generalized the following confusion matrix. Based on the following table 1, we predicted 11481 flights information correctly and 683 wrong.

Prediction/Observation	Not Delay	Delay
Not Delay	11061	607
Delay	76	420

Table 1: Confusion Matrix of Predicting Flight Delay and Actual Flight Delay

## Stacking

(Alex's paragraph pending) (ROC comparison plots and results comparison table across different methods discussed above pending)

## Discussion

(Go Alex!!!)