



# Predicting Flight Departure Delays in an Air Traffic Network

Capstone Project 1

# The problem

## Potential Clients

- Companies distributing use of the Aircraft Situation Display to Industry (ASDI)
- Air Route Traffic Control Centers (ARTCC) looking to make better predictions of AAR's and plan efficient holding patterns.
- Implementers of Ground Delay Programs at Domestic Airports
- Traffic Management Personnel seeking to get the earliest Expect Departure Clearance Time (EDCT)

## Context - Cost of Delay

- \$40 billion in annual cost due to delays in 2015
- \$62.55 per minute average cost of aircraft block (taxi plus airborne) in 2016
- In 2016 nearly 40% of delays were due to the delayed arrival of the incoming aircraft, reflecting the high levels of interdependence in the delay dynamics

## Problem statement

Predict whether or not the departure delay on an Origin-Destination pair, with a 6-hour prediction horizon, will exceed a 15 min threshold; such that the future delay will fall into one of two classes, 'above threshold' (1) and 'below threshold' (0).

# Methodology



## Wrangling

1. Improve readability and address missing values
2. Aggregate to create mock network
3. Create target variable of departure delay state 6 hours in the future

## Storytelling

1. Flight Frequency
2. Delay Distribution
3. On-Time Performance
4. Network Visualization

## Inferential Statistics

1. Test for Normality and CLT in variables of interest
2. Regression Analysis
3. Hypothesis Tests

## Binary Classification

1. Baseline Logistic Regression Classifier
2. Resample data to address target class imbalance, both under- and over-sampling
3. Train Logistic Regression and Random Forest Classifiers under both resampling conditions

# Wrangling

Surface Cleaning and  
Pre-Processing

1. Import data and address missing values
  2. Compartmentalize columns and evaluate by category, then merge results into final Flights dataframe
  3. Create NetworkX Digraph from components
  4. Create mock air traffic network, a dataframe named Links\_d, containing target class for model construction
-

# Flights DataFrame

```
<class 'pandas.core.frame.DataFrame'>
```

```
MultiIndex: 1824403 entries, (ABE-ATL, 2016-01-01 07:00:00) to (YUM-PHX, 2016-12-31 19:15:00)
```

```
Data columns (total 39 columns):
```

quarter	int64	crs_dep_time	int64
month	int64	dep_time	float64
day_of_month	int64	dep_deviation	float64
day_of_week	int64	dep_delay	float64
fl_date	object	wheels_on	float64
Day_of_Week	object	taxi_in	float64
Month	object	crs_arr_time	int64
dt_index	datetime64[ns]	arr_time	float64
hour_of_day	object	arr_deviation	float64
unique_carrier	object	arr_delay	float64
fl_num	int64	crs_elapsed_time	float64
origin_airport_id	int64	actual_elapsed_time	float64
origin_city_market_id	int64	air_time	float64
origin	object	distance	float64
origin_city_name	object		
origin_state_abr	object		
origin_state_nm	object		
dest_airport_id	int64		
dest_city_market_id	int64		
dest	object		
dest_city_name	object		
dest_state_abr	object		
dest_state_nm	object		
link	object		
unique_carrier_nm	object		

```
dtypes: datetime64[ns](1), float64(12), int64(11), object(15)
```

```
memory usage: 554.5+ MB
```

# Links\_d DataFrame

```
<class 'pandas.core.frame.DataFrame'>
MultiIndex: 882815 entries, (ANC-SEA, 2016-01-01 00:00:00) to (TPA-ATL, 2016-12-31 19:00:00)
Columns: 561 entries, crs_dep_time to dd_binary
dtypes: float64(15), int32(1), uint8(545)
memory usage: 566.6+ MB
None
Index(['crs_dep_time', 'dep_time', 'dep_deviation', 'dep_delay', 'wheels_on',
       'taxi_in', 'crs_arr_time', 'arr_time', 'arr_deviation', 'arr_delay',
       ...,
       'dest_city_name_San Diego', 'dest_city_name_San Jose',
       'dest_city_name_San Juan', 'dest_city_name_Seattle',
       'dest_city_name_St. Louis', 'dest_city_name_Tampa',
       'dest_city_name_Washington',
       'dest_city_name_West Palm Beach/Palm Beach', 'dd_in_6hrs', 'dd_binary'],
      dtype='object', length=561)
```

# Storytelling

Exploratory Data Analysis

1. Flight Frequency
    - By Temporal Categories
    - By Commercial Airline Carrier
    - By Location Attributes
  2. Delay Distribution
    - Overall
    - By Commercial Airline Carrier
    - For Highest Traffic Airports
  3. On-Time Performance
  4. Network Visualization
-

# Air Traffic Statistics Explained

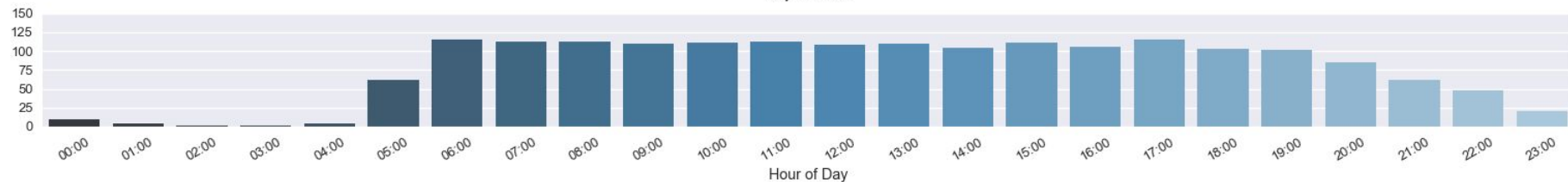
- Sample contains all commercial flights in January, August, November, and December of 2016, from the 12 Major Domestic Passenger Airlines.
- The U.S. Department of Transportation defines a major carrier or major airline carrier as a U.S.-based airline that posts more than \$1 billion in revenue during fiscal year, grouped accordingly as “Group III” (“Air Carrier Groupings 2016”, U.S. Bureau of Transportation)
- A flight is counted as “on time” if it operated less than 15 minutes after the scheduled time show in the carriers’ Computerized Reservation Systems (CRS), with departure and arrival times being calculated from gate to gate, not including taxi or airtime.



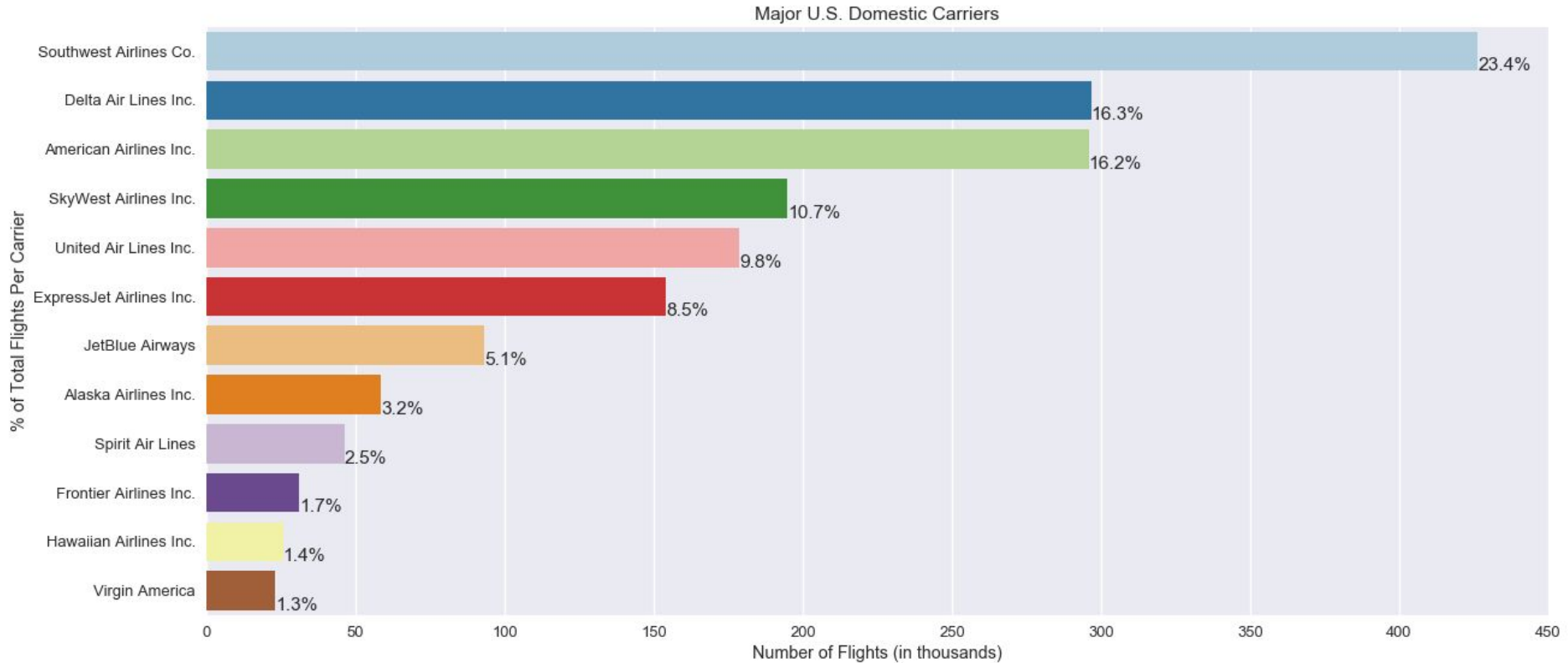


# 1.1 Frequency by Time

Flight Frequency



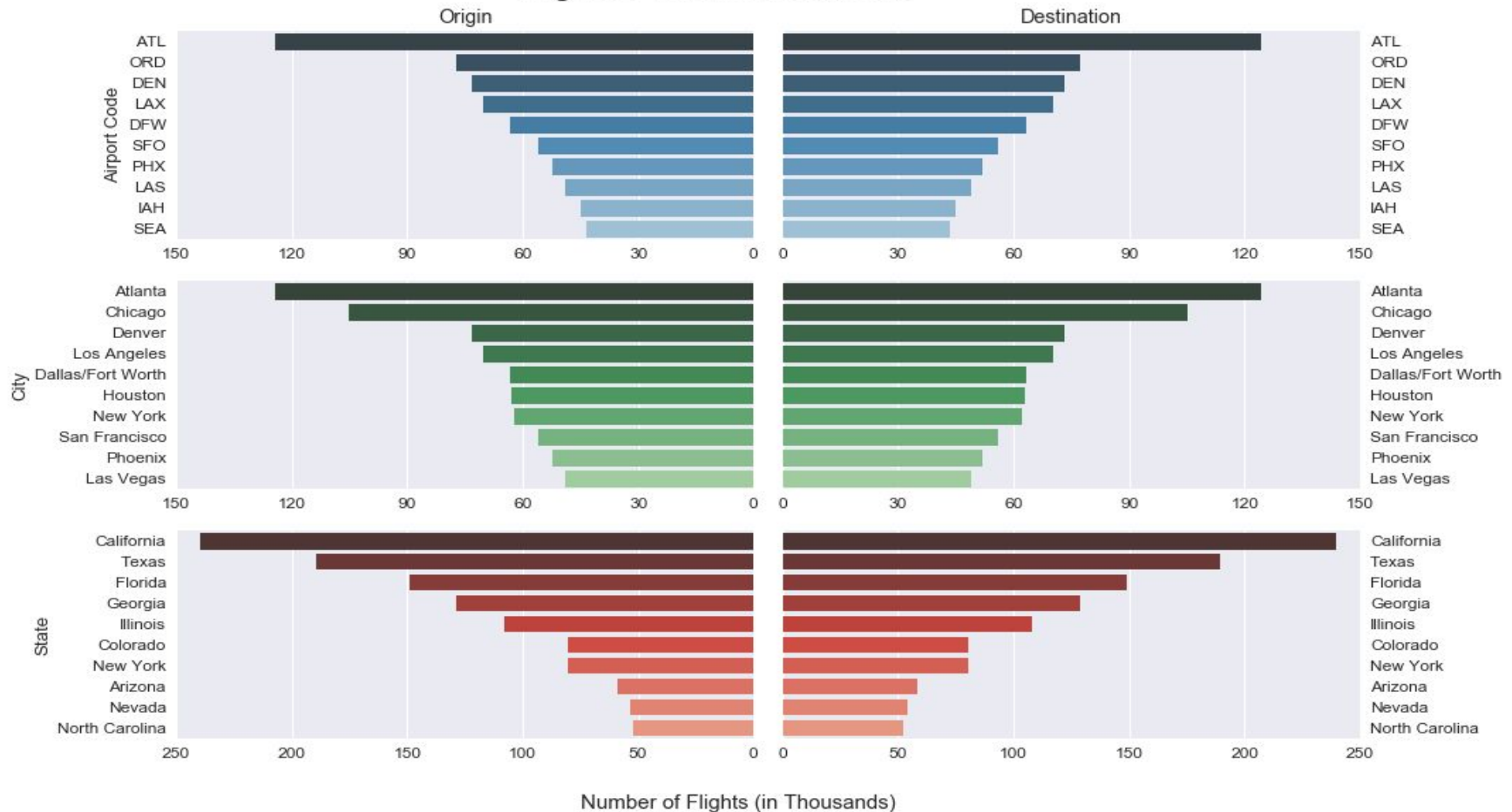
# 1.2 Frequency by Carrier



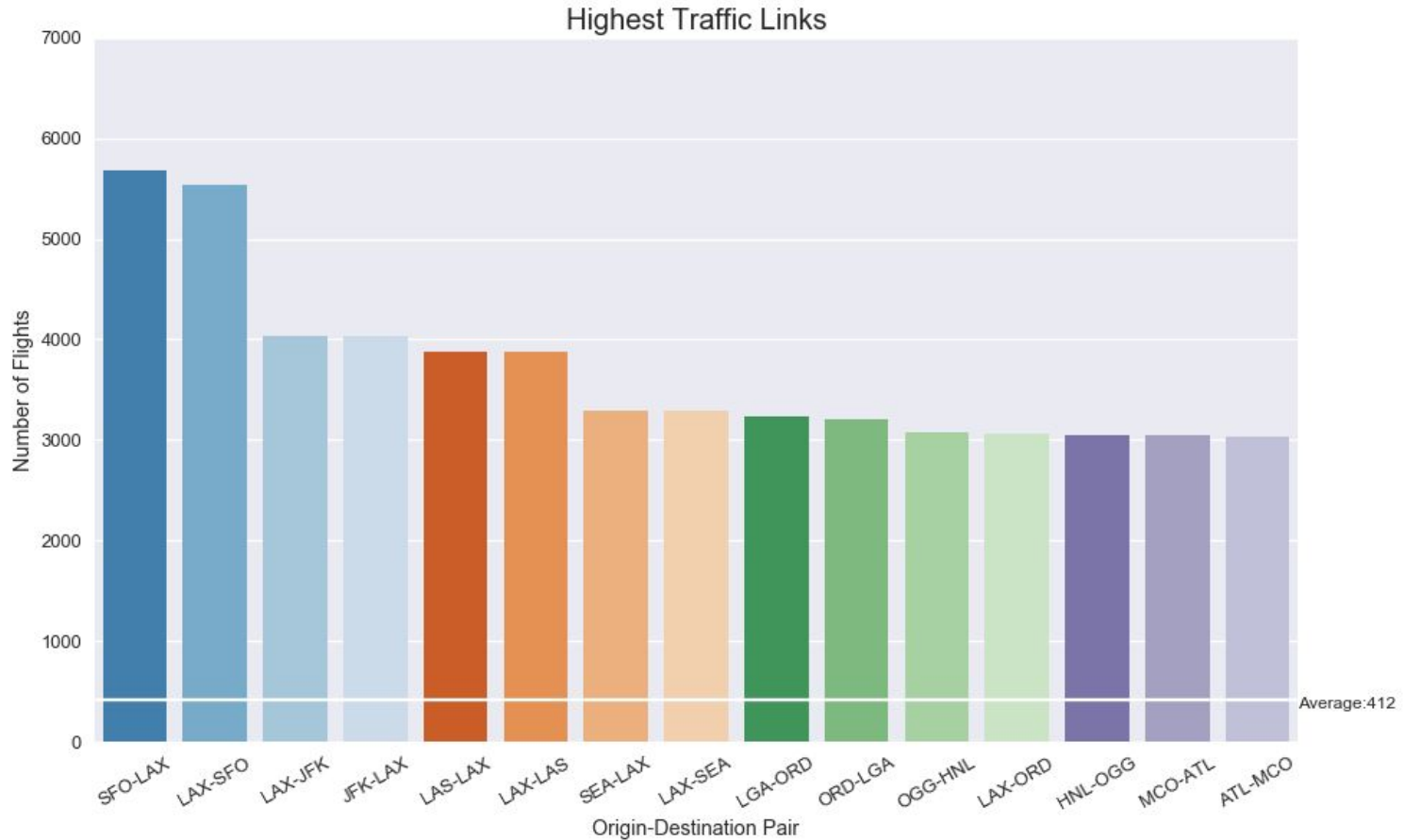
## 1.3 Frequency by Location

	Number of Locations Reported
Origin Airport	309
Origin City	297
Origin State	52
Destination Airport	308
Destination City	297
Destination State	52
Origin-Destination Pair (Link)	4431

# 1.3 Highest Traffic Locations

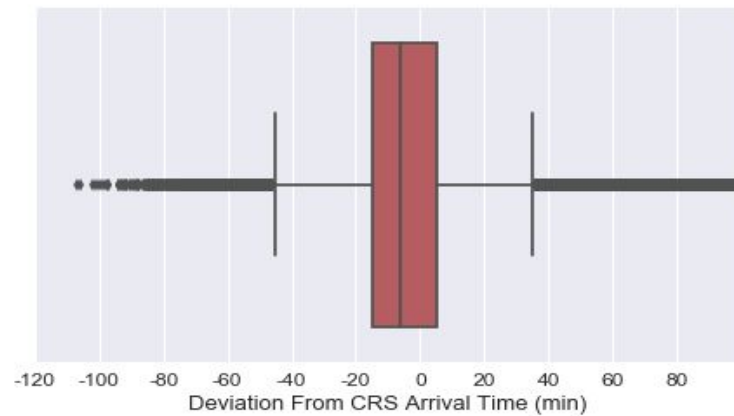
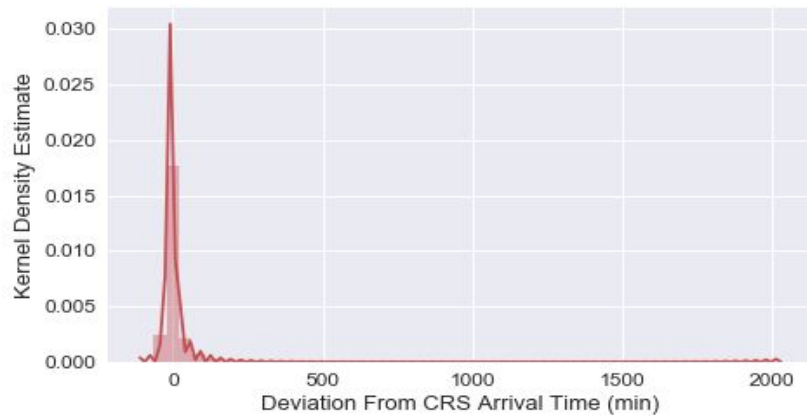
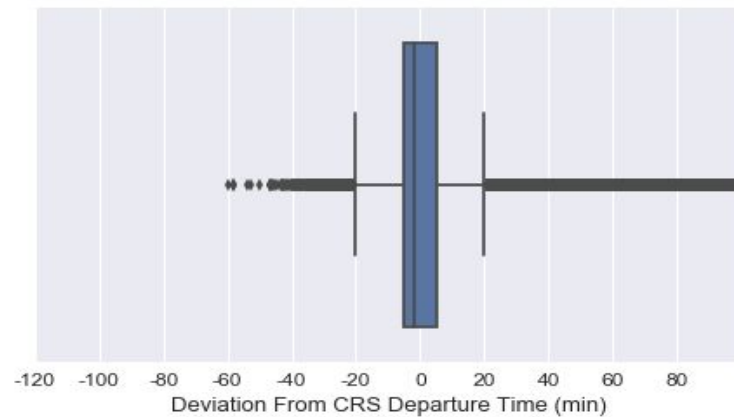
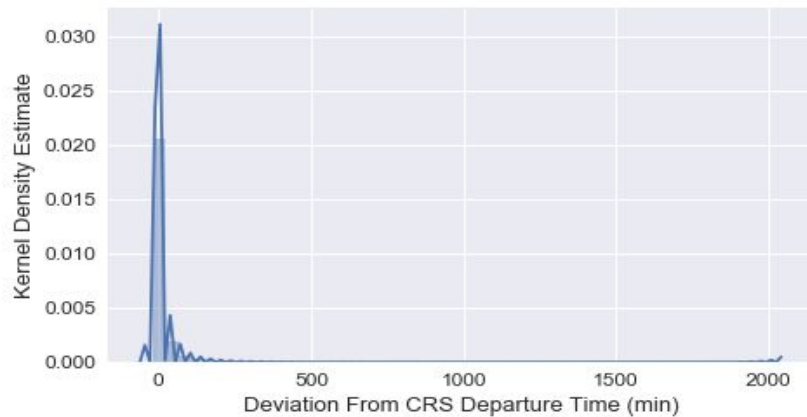


## 1.3 Highest Traffic Origin-Destination Pairs

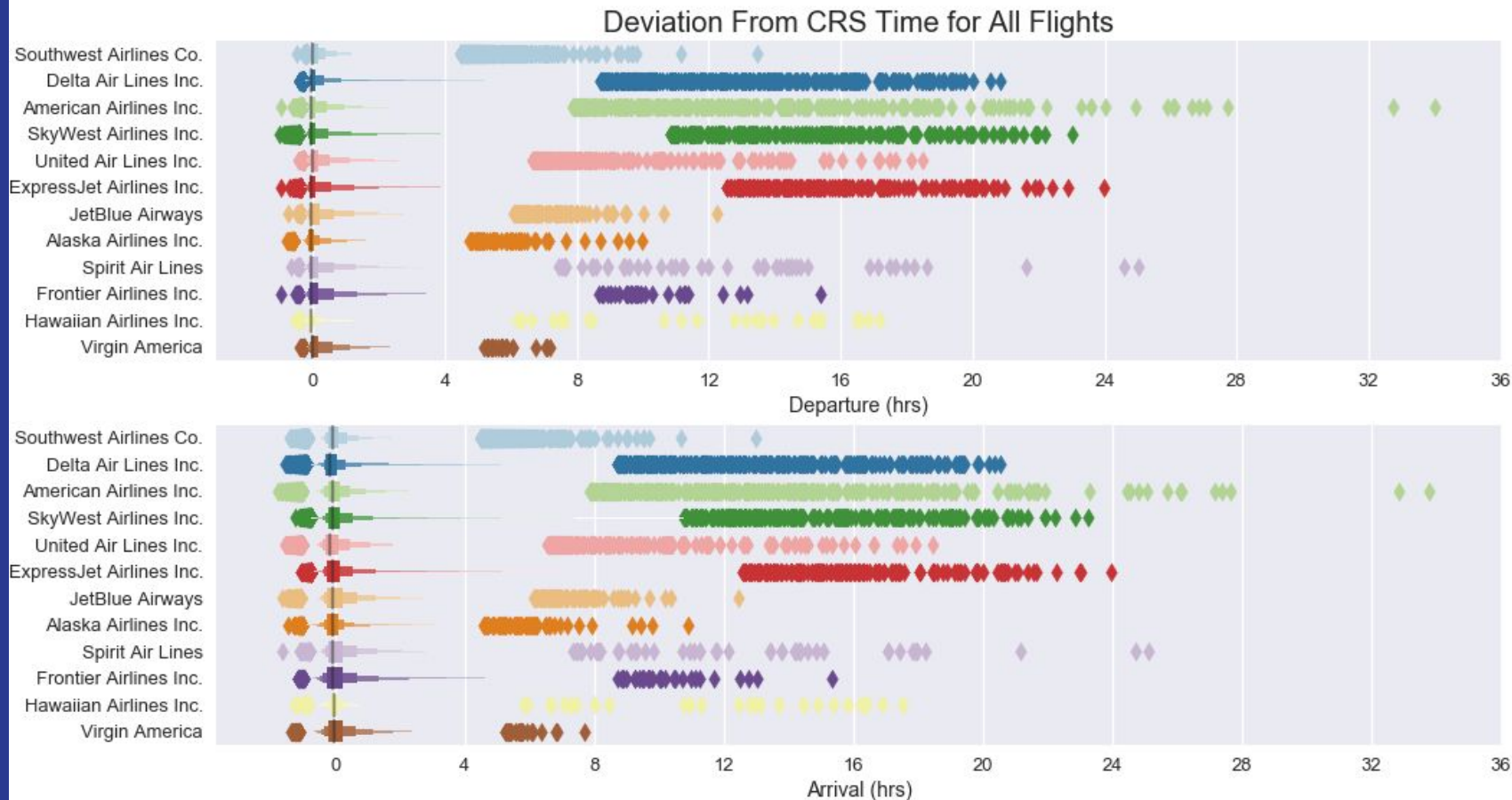


## 2.1 Overall Delay Distribution

Distribution of Delay Across All Flights

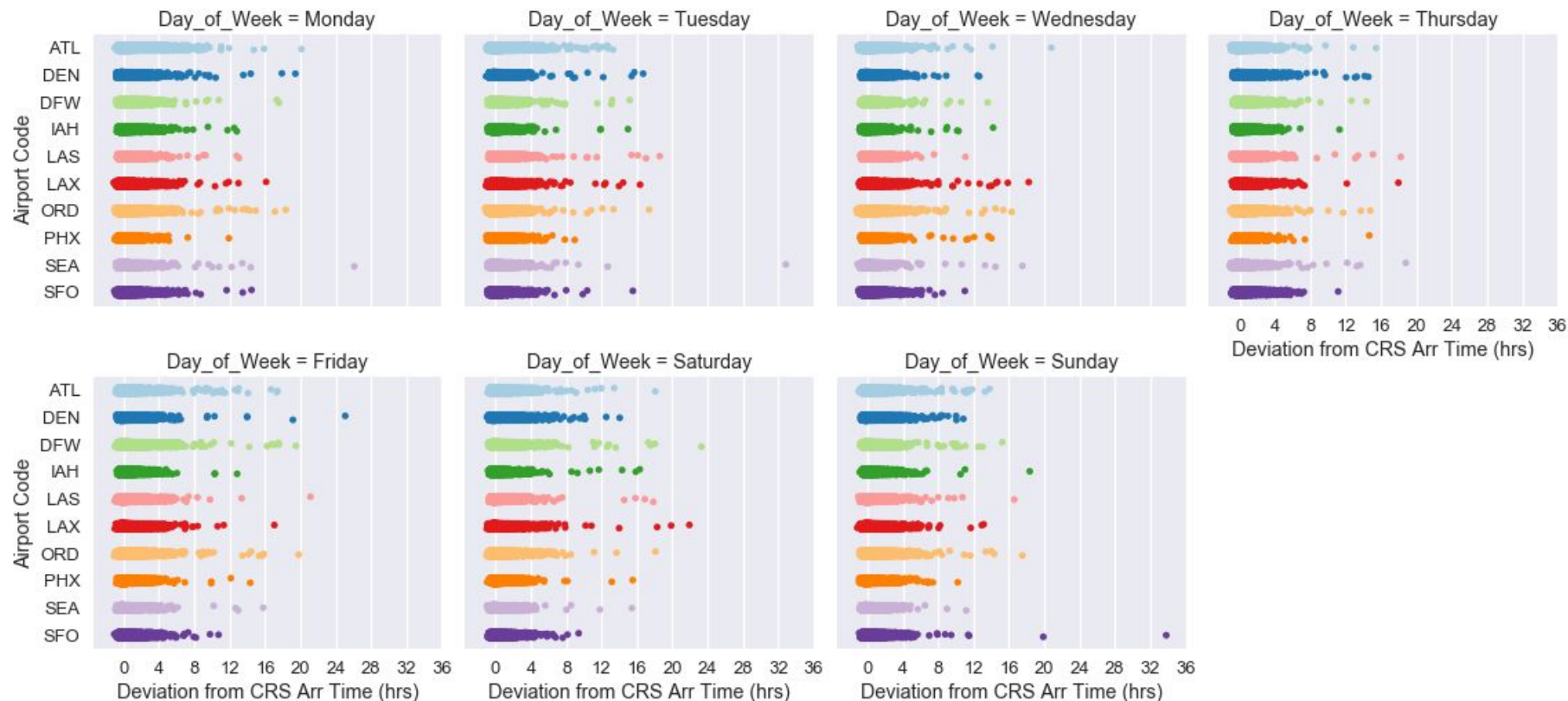


## 2.2 Distribution of Delay by Carrier



## 2.3 Distribution of Delay for Highest Traffic Airports

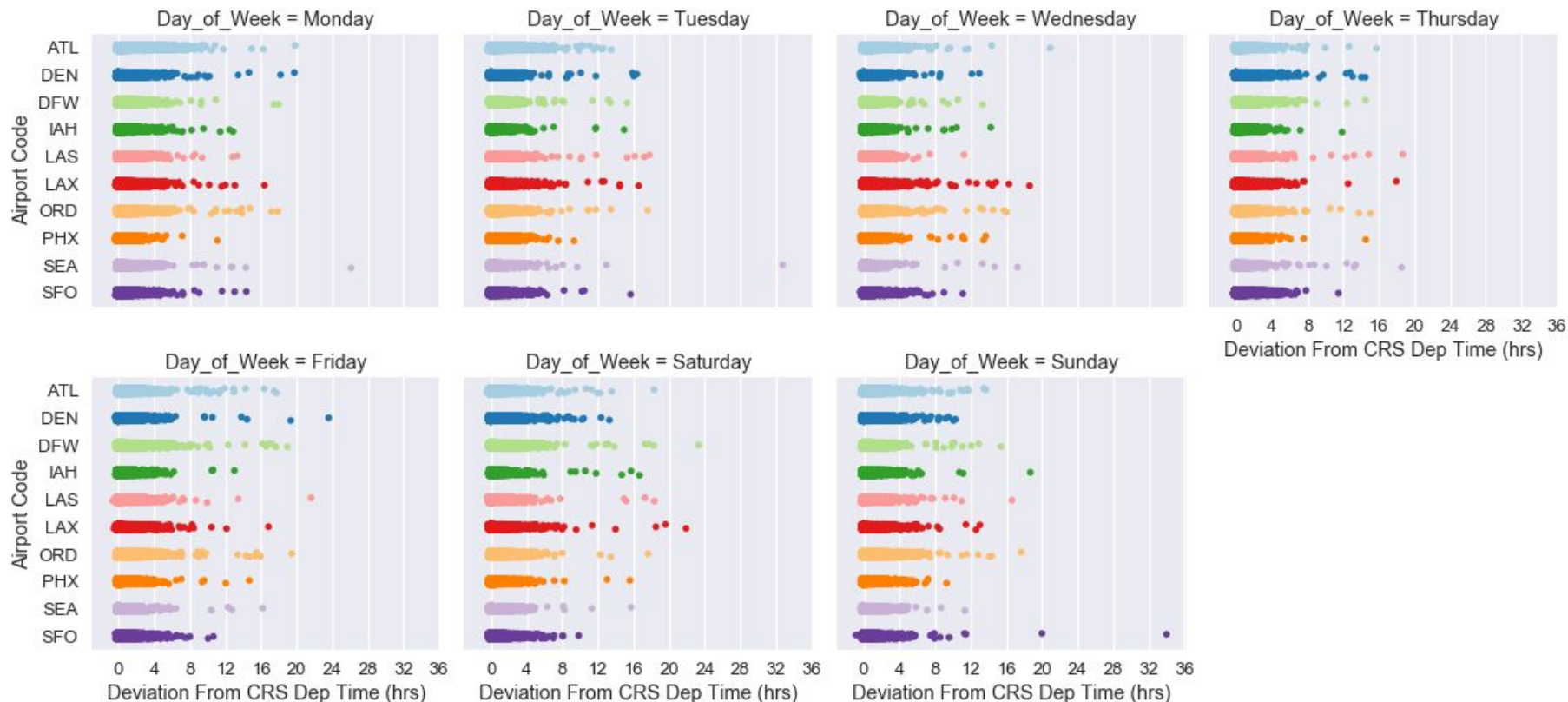
Arrival Distribution For Top 10 Airports



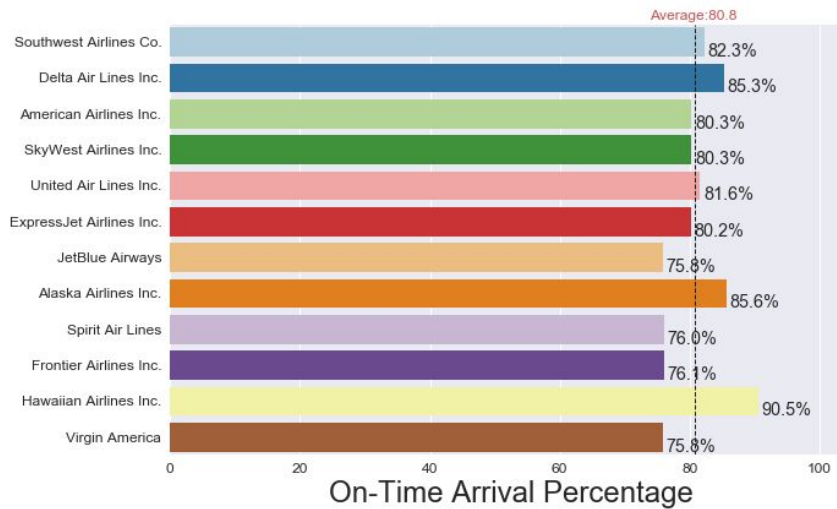
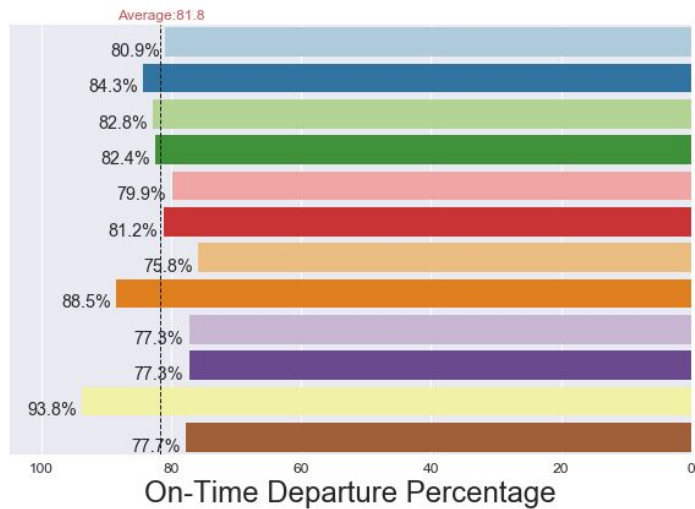
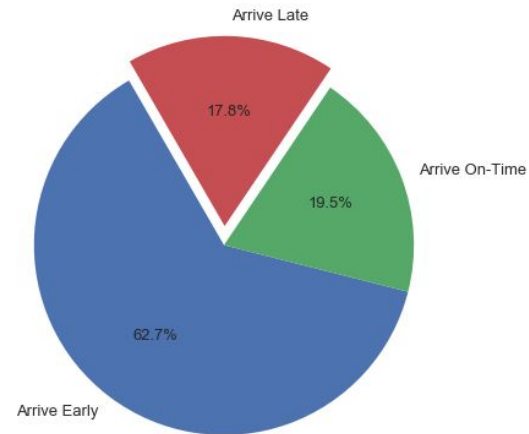
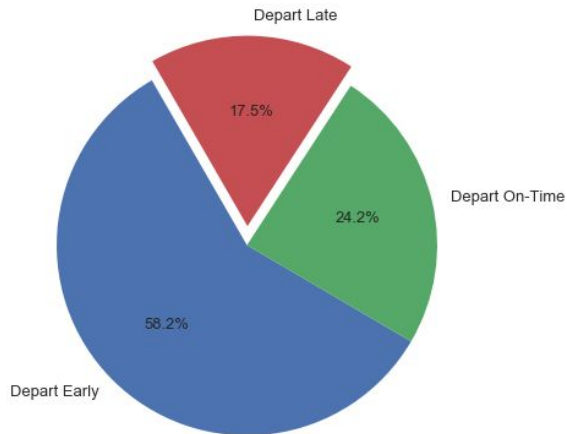


## 2.3 Distribution of Delay for Highest Traffic Airports

Departure Distribution For Top 10 Airports

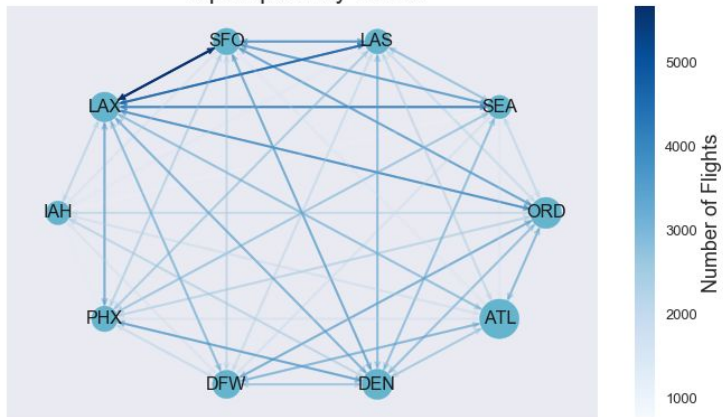


# On-Time Performance For All Flights

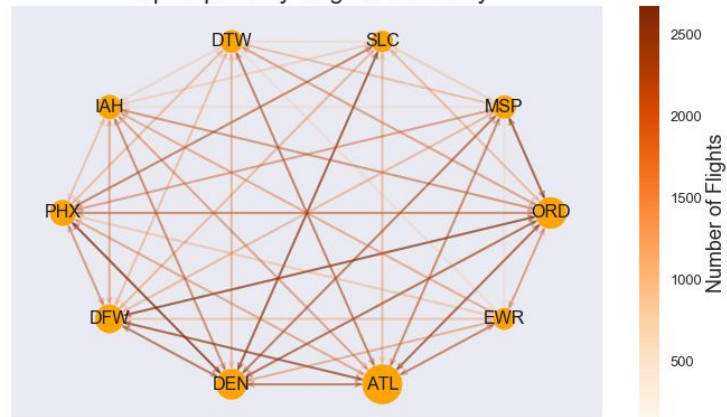


# 4. Network Visualization

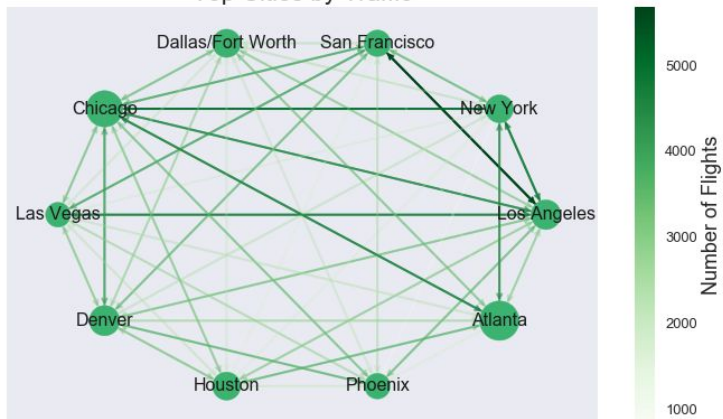
Top Airports by Traffic



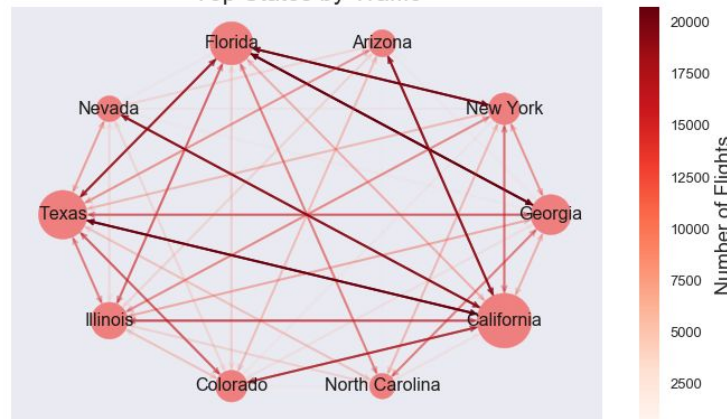
Top Airports by Degree Centrality



Top Cities by Traffic



Top States by Traffic



# Inferential Statistics

Initial Data Analysis


1. Test for Normality and CLT in variables of interest
  2. Regression Analysis
  3. Hypothesis Tests
    - T-Test of Arrival and Departure Deviation Means
    - Pearson's  $r$  Permutation of Median Arrival and Departure Delays
-

# Variables of Interest

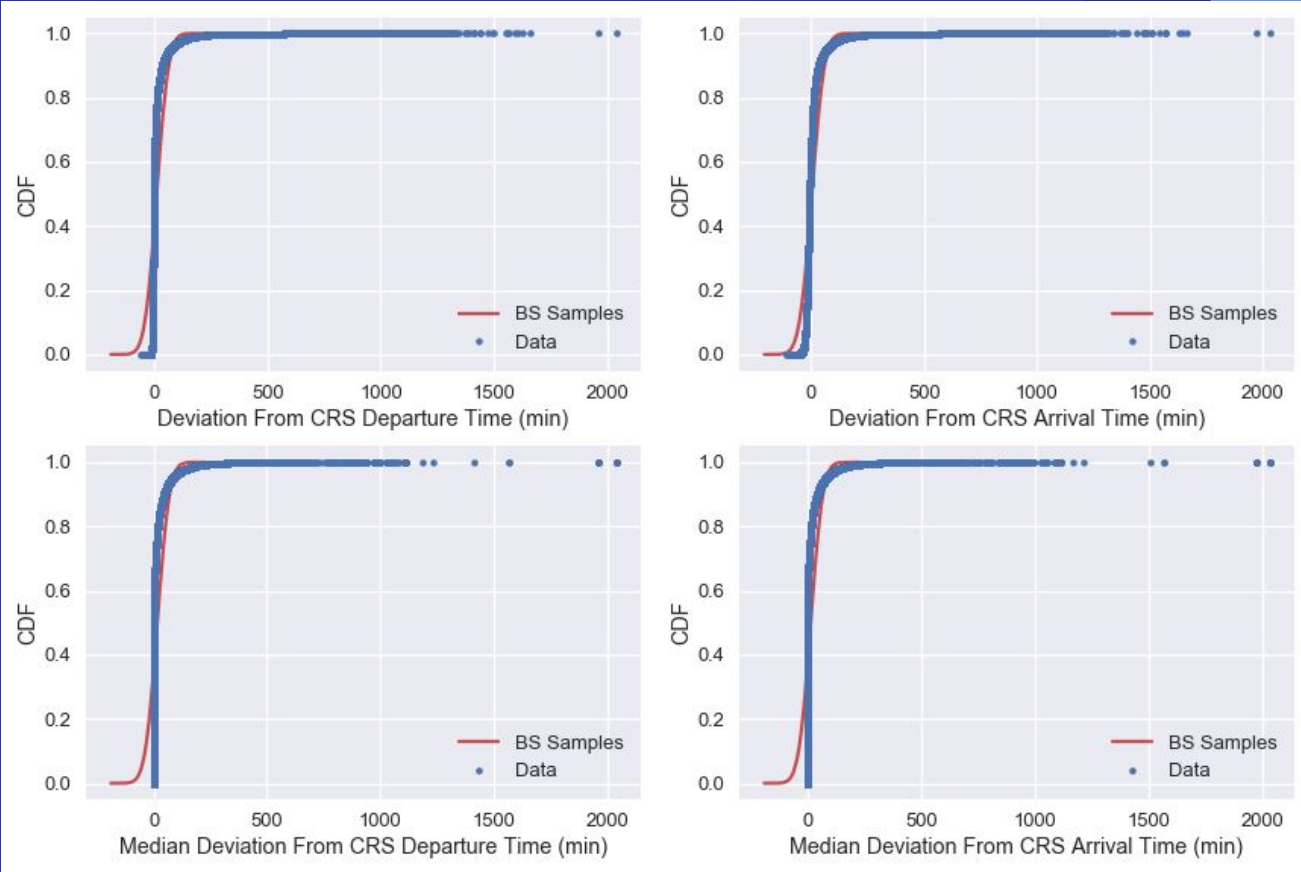
## **Flights DataFrame:**

- dep\_deviation: deviation from the actual departure time to the scheduled (CRS) departure time
- arr\_deviation: deviation from the actual arrival time to the scheduled (CRS) arrival time

## **Links\_d DataFrame:**

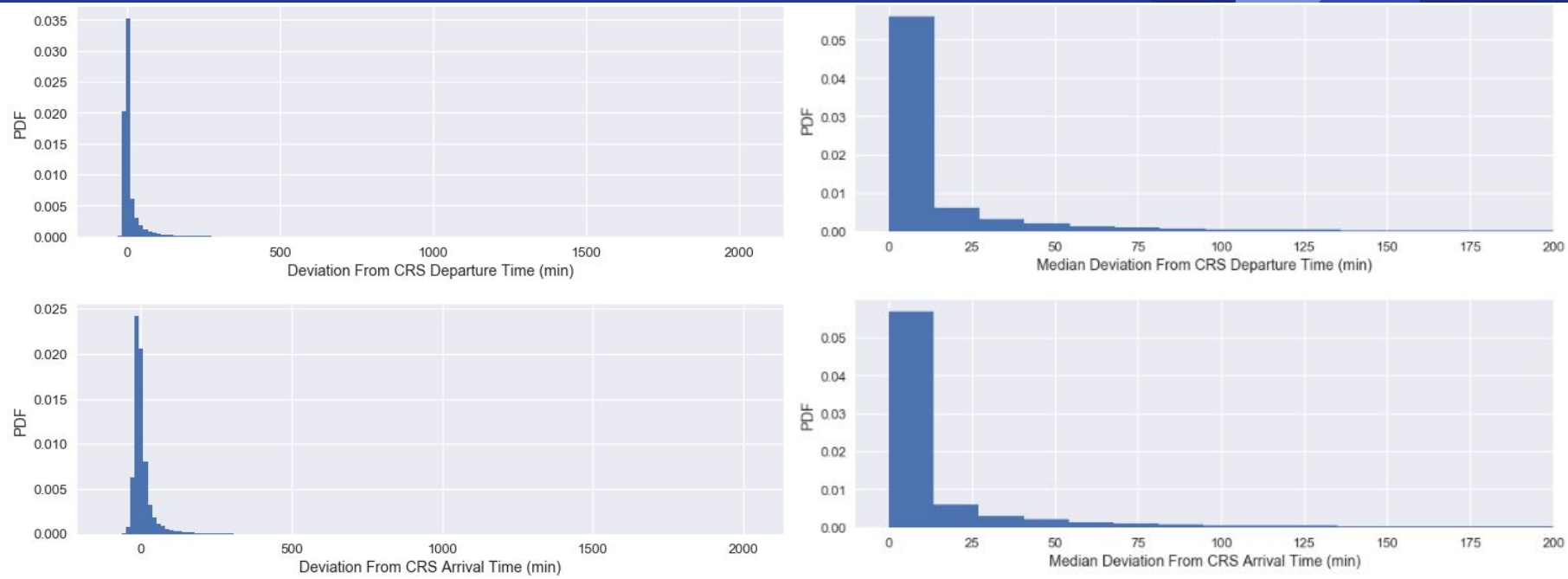
- dep\_delay: median deviation of the actual departure time from the scheduled (CRS) departure time, for an origin-destination pair, in an hour of day that had non-zero traffic
  - arr\_delay: median deviation of the actual arrival time from the scheduled (CRS) arrival time, for an origin-destination pair, in an hour of day that had non-zero traffic
- 

# Test for Normality in Variables of Interest



The Cumulative Density Function of each variable shows they all possess normal distributions

# Test for Central Limit Theorem

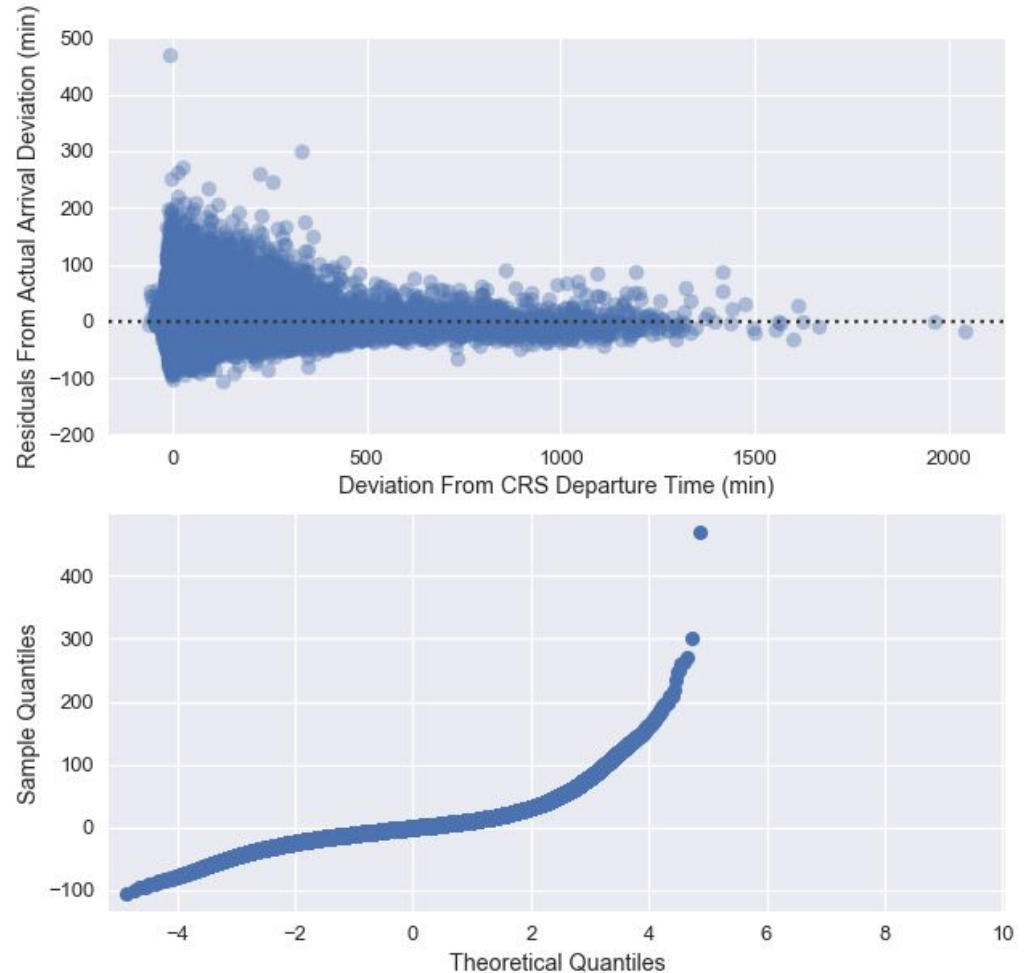


Central Limit Theorem applies as sample size is very large and the probability density function of each variable shows observations are independent



# Regression Analysis

- For the purpose of validating a directed network approach utilizing only departure delay as a predictor of overall (departure and arrival) delay
- The residuals between departure and arrival deviation are in a random pattern, supporting the use of a linear model. The quantile plot shows that the datasets are heavily skewed, and robust methods should be used in model construction to lessen the influence of extreme values.





# T-Test of Arrival and Departure Deviation Means

$$H_0 : \mu_d = \mu_a$$

$$H_a : \mu_d \neq \mu_a$$

Margin of Error: 0.045

Difference of Means: 5.549

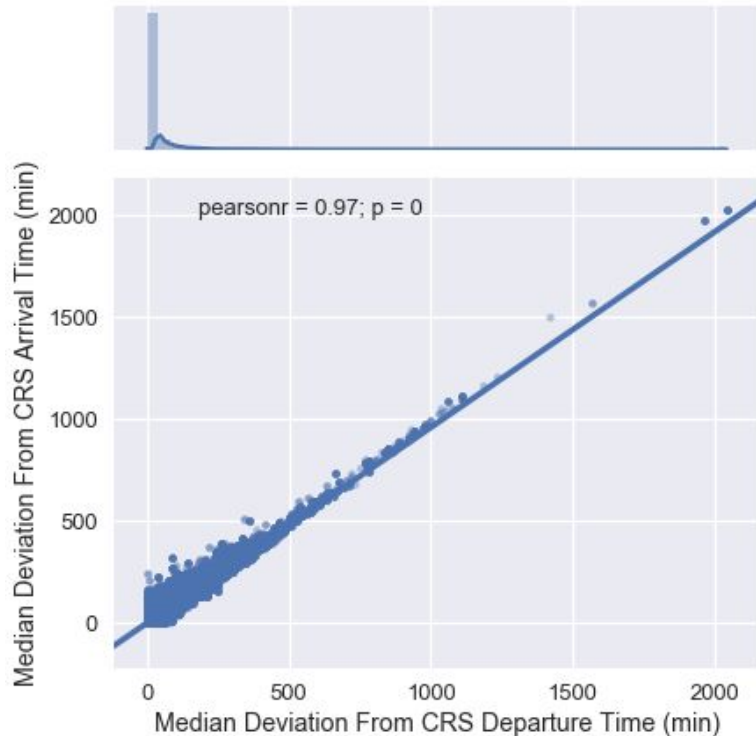
Departure Deviation 95% Confidence Interval: [-11. 112.] min

Arrival Deviation 95% Confidence Interval: [-31. 112.] min

T-Test: tstat = -124.015, P-value = 0.0000000000

With an alpha level of .01 ( $\alpha = .01$ ), the difference between the means of departure and arrival deviation from scheduled (CRS) time was statistically significant,  $p < .01$

# Pearson's r Permutation Test of Median Arrival and Departure Delays



$H_0$  : The correlation between the current median departure delay and median arrival delay for an Origin-Destination pair is not significant

$H_a$  : The correlation between the current median departure delay and median arrival delay for an Origin-Destination pair is significant

With an alpha level of .01 ( $\alpha = .01$ ), the correlation between the median departure delay and the median arrival delay for Origin-Destination pairs is statistically significant,  $p < .01$

# Model Construction

Binary Departure Delay  
Classification

1. Baseline Logistic Regression Classifier
  2. Resample data to address target class imbalance, both under- and over-sampling
  3. Train Logistic Regression and Random Forest Classifiers under both resampling conditions
-