# JFK TAXI-OUT Linear Regression

Nov 2019 – Jan 2020

Deirdre Boland
07 Feb 2025

# Objective and dataset

**Objective**
- Create Linear Regression Model of JFK Taxi-Out, using data scraped from an Academic Paper under Review by IEEE transportation covering Nov 2019- Jan 2020 (D Kansal, Kaggle dataset)
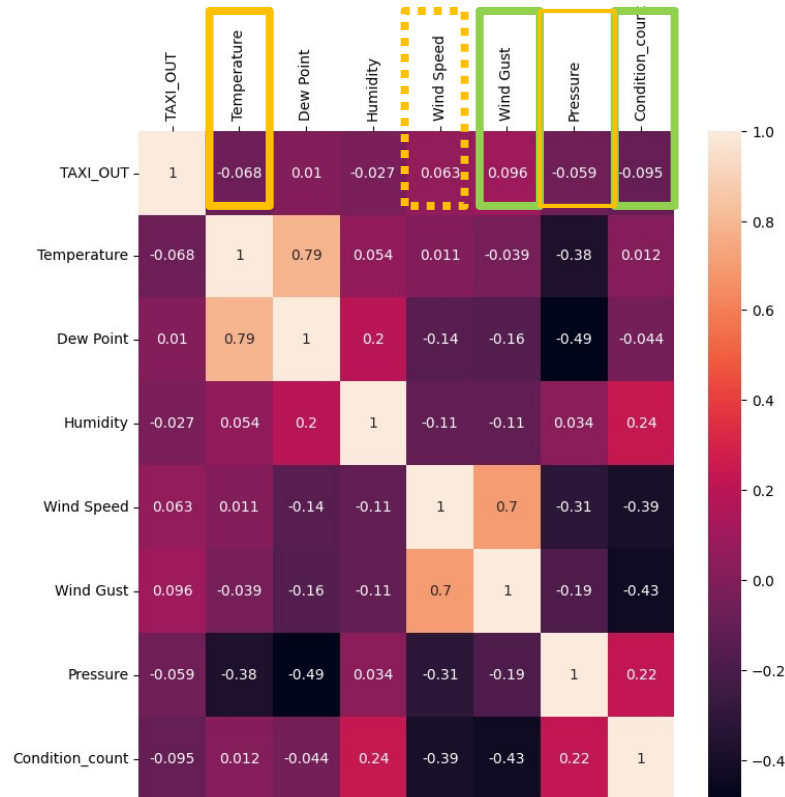
**Value**
- At JFK airport Taxi-Out prediction is an important concept for calculating runway time and directly impacts the cost of flights.
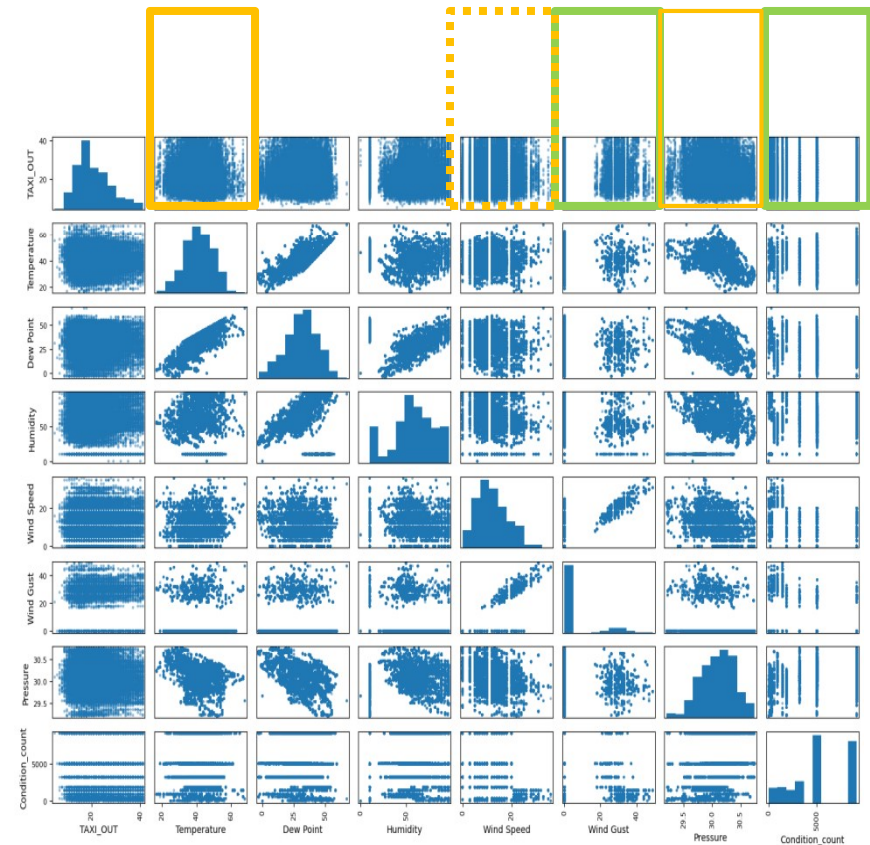
**Dataset**
- 5 text based features/variables
  - Airline and flight number indicators (TAIL_NUM, OP_UNIQUE_CARRIER)
  - Destination (DEST) covered by numeric features of distance and scheduled flight time
  - Wind direction (Wind e.g. NW, E etc) – have wind speed and gust in numeric
  - Climate/weather (Condition) – converted to numeric by using frequency counts (more severe weather less frequent)
- 18 numeric features, including target feature(TAXI_OUT)
  - Covering weather/conditions at flight time and flight details e.g. Time, duration
  - Scheduled Arrival Time (CRS_ARR_M) excluded – arrival at another airport, includes time difference offsets. This feature is covered by scheduled duration of flight feature
- Final 18 numeric features (17 original & 1 generated from Condition text feature)

- Reference – Kaggle page Flight Take Off Data - JFK Airport

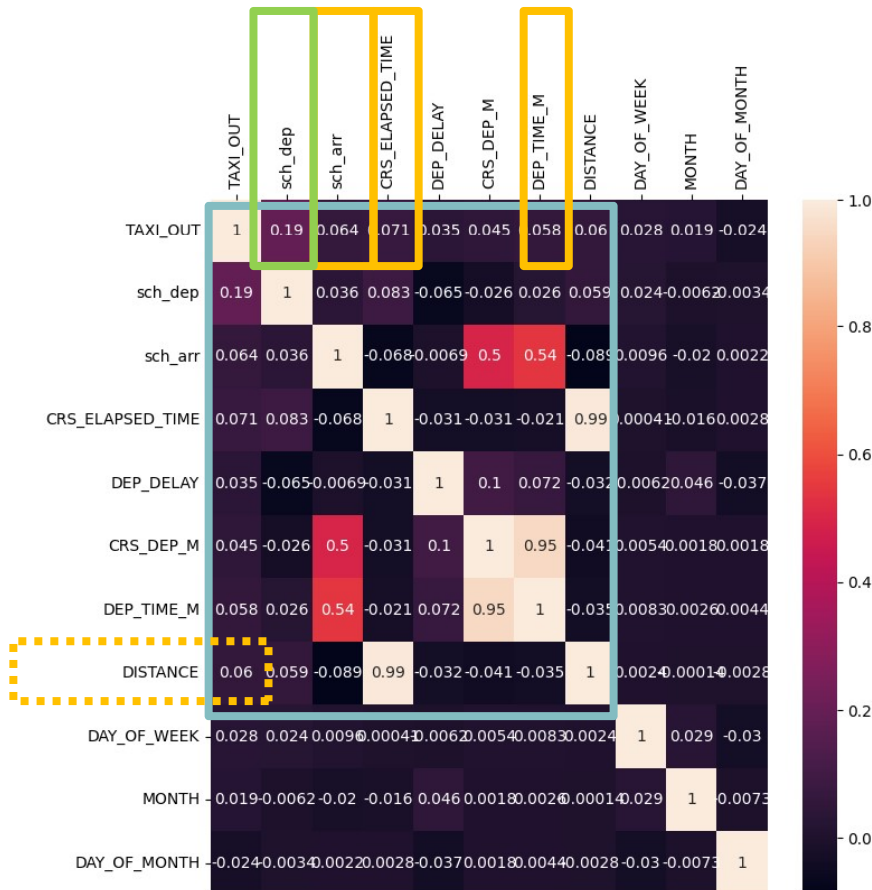# Correlation Matrix – Weather/Climate

# Correlation Matrix – Flight Details

## Correlation Matrix



## Scatter Matrix

- Number of flights scheduled for arrival./departure (sch_arr, sch_dep)
- Departure delay of the flight (DEP_DELAY) - is calculation of *Actual Departure Time (DEP_TIME_M ) subtract Scheduled Departure Time (CRS_DEP_M)
  * Gate checkout of the flight not the take off time



- Scheduled journey time of the flight (CRS_ELAPSED_TIME)

# Linear regression Model

**Feature Histograms**
(x- axis, horiz = units bin 40, y- axis, vert = frequency)



- *All models scaled using standard scaler. 80% train, 20% test*
- $R^2$ 0-1, higher is better fit

3 features (highest corrl, green)
- Multi Linear Regression (MLR) or Ordinary Least Squares (OLS) – $r^2$ 0.059
- Lasso – $r^2$ 0.057
- Ridge – $r^2$ 0.059

8 features (>0.5 corrl, ambr non-dash & green)
- MLR/OLS – $r^2$ 0.077
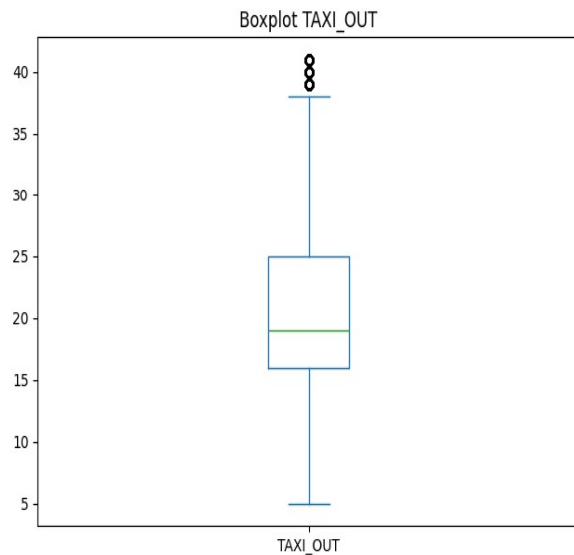- Lasso – $r^2$ 0.074
- Ridge – $r^2$ 0.077

14 features (all except red)
- MLR/OLS – $r^2$ 0.083
- Lasso – $r^2$ 0.080
- Ridge – $r^2$ 0.083

*Alpha tuning (set at 0.1 above)*
- *Regularisation penalty in ridge and lasso modelling.*
- *Increasing value on Lasso made model dramatically worse but minimal impact Ridge*
- *Expected as increasing alpha will get rid of features on Lasso*

# Conclusion and next steps

Boxplot TAXI_OUT



- Linear regression model not viable to model JFK Taxi-out with this data set
  - low correlation of features leading to underfitting (erroneous outcomes on new data)
- Limited dataset Nov 2019 – Jan 2020, includes holidays and wintry conditions
- Majority Taxi-out time under 25 mins with a min of 5 and max of 40mins
  - Need more information on Taxi-out metric and calculation
    - What is target? Is consistent 5-15 min realistic?
    - What other data not collected, might be impacting? E.g. impact of different terminals (5) and runways (4), staffing levels

Q&A