

Airport Taxi Time Analysis

Group 6

0. Introduction - Data Creation

• What is Taxi Time?

Amount of time an aircraft spends on the surface of the airport while engine is running.

$$\text{TotalTaxi} = \text{Taxi Out} + \text{Taxi In}$$

• Dataset creation

Our dataset was handcrafted using different sources (4), each one containing different information we considered crucial for our "totalTaxi" analysis:

- 1 - Combined Flights 2019: 2019 USA domestic flights.
- 2 - Airport Weather 2019: USA airports weather conditions, with some corrections (ex. Phoenix).
- 3 - Aircraft dataset: aircrafts specs.
- 4 - Concurrent flights: meaning flights that arrived (arrivals) or took off (departures) from the same airport, the same day, at the same time block.

RESEARCH QUESTIONS:

- How do different airlines perform with respect to Taxi Time?
- What are the most impactful features that explain Taxi Time?
- Can we actually predict (accurately) Taxi Time?

• Airports selection

We kept 3 origins: JFK, BOS, and PHX.

Why?

First, they have a similar domestic traffic; Secondly, JFK has the highest TaxiTime; Lastly, we selected one in the North and one in the South.

• Final dataset

125k observations, in which each "route" is one individual, with 24 features, including our new "totalTaxi".

• Facts about consumption (B777)

APU	Ground	Flight
3.9kg Fuel	23.3kg Fuel	113.3kg Fuel
12.3kg CO ₂	73.4kg CO ₂	356.9kg CO ₂
per min	per min	per min

1. Data Preparation

• Downsampling

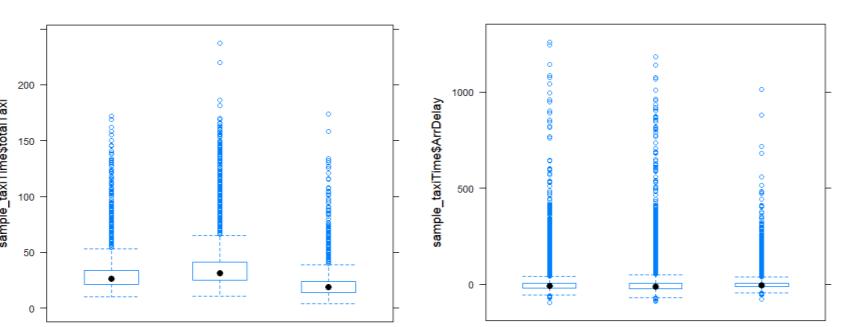
Keep only 20k samples per airport to make our analysis faster.

• Data conversion

Categorical variables conversion into factors and dates into datetime type.

• Outliers

taxiTime & ArrDelay outliers detected



We keep only rows s.t:

$$Q1-1.5*IQR < \text{taxiTime} < Q3+1.5*IQR$$

$$Q1-1.5*IQR < \text{ArrDelay} < Q3+1.5*IQR$$

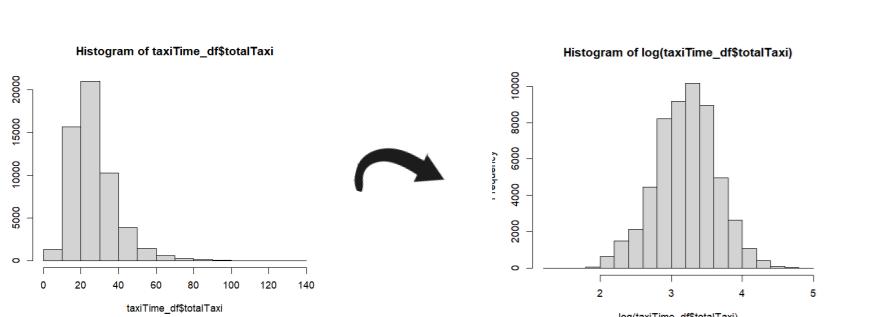
• Numerical feature scaling

To make the linear regression models more performant.

2. Feature Selection

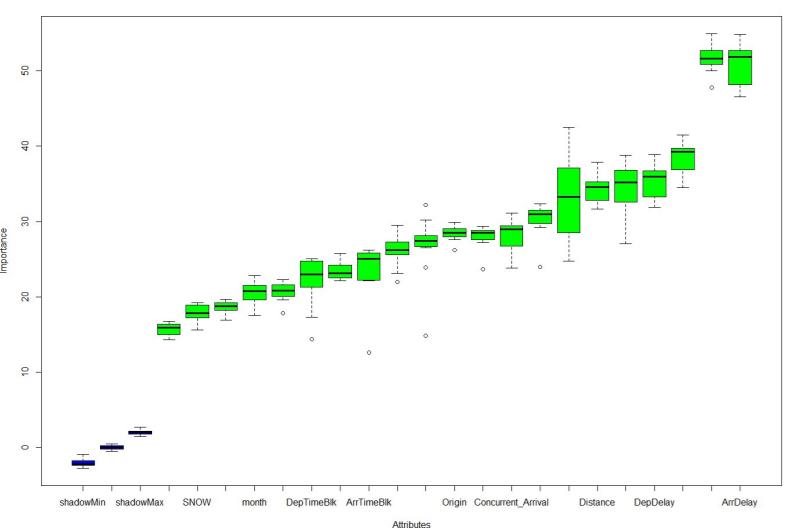
• Output log transformation

totalTaxi is left skewed -> normalization



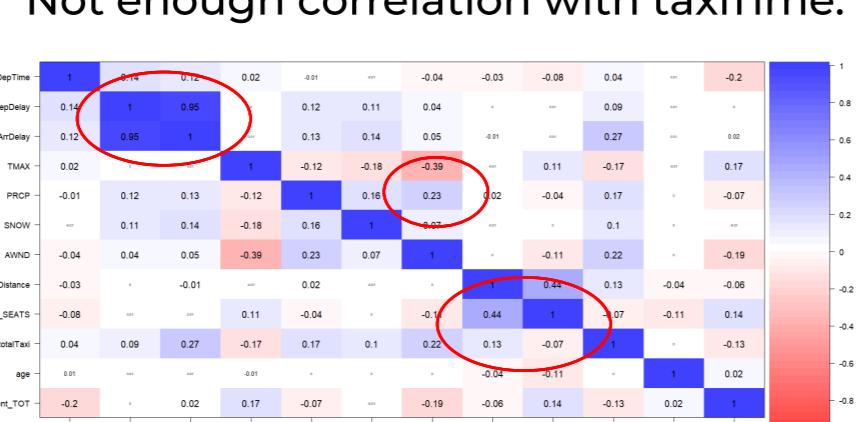
• Random forest selection

Get an overview of the most impactful features:



• Correlation plot

- Too much interdependant
- Not enough correlation with taxiTime:



• Backward linear selection

- First linear regression model by including all the variables
- Many of them are not significant enough -> backward selection to keep only the best predicates:

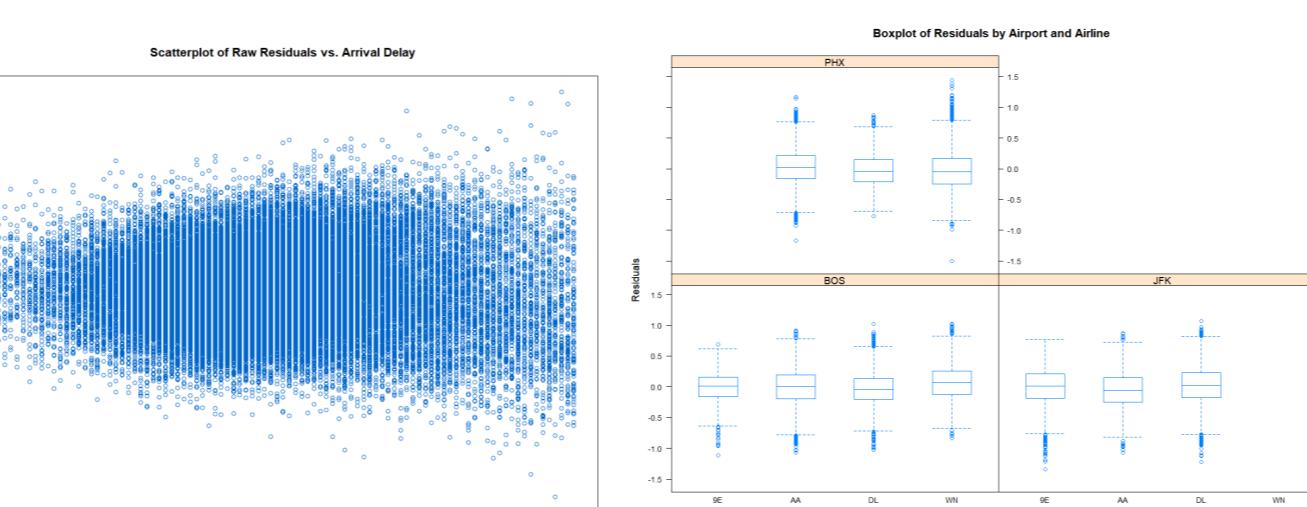
Coefficients:	Value	Std. Error	t-value	p-value
(Intercept)	3.423770	0.004855565	709.0922	0
Operating_AirlineAA	-0.049730	0.004932575	-10.0409	0
Operating_AirlineAL	-0.127620	0.004821799	-26.4682	0
Operating_AirlineAN	-0.142020	0.004821799	-30.4646	0
ArrDelay	0.173212	0.004361881	39.7104	0
Distance	0.052100	0.00423907	36.5948	0
OriginATL	0.143760	0.004344777	36.0422	0
DestATL	-0.131300	0.004344777	-36.0426	0
concurrent_TOT	0.049337	0.004204585	35.1257	0
Operating_AirlineAA:ArrDelay	-0.034312	0.004446464	-7.7324	0
Operating_AirlineAL:ArrDelay	-0.034312	0.004446464	-7.7324	0
Operating_AirlineAN:ArrDelay	-0.078890	0.005339964	-13.8803	0
ArrDelay:OriginATL	0.017751	0.00321152	5.4599	0
ArrDelay:DestATL	-0.022472	0.003631169	-6.1851	0

3. Variance Structure

Linear Model Formulation:

$$\log(\text{TaxiTime}) \sim \text{Airline} + \text{ArrDelay} + \text{Distance} + \text{OriginAirport} + \text{ConcurrentFlights} + \text{Airline:ArrDelay} + \text{Origin:ArrDelay} + \text{ArrDelay:Origin:ArrDelay}$$

Where do we see Heteroscedasticity?



- Clear Heteroscedasticity of residuals vs. Arrival Delay.
- Possible Heteroscedasticity with respect to Airline and Origin Airport.
- Residuals vs. Other Variables appear to be Homoscedastic.

Variance Structure Selection:

Step 1: Grouping Factor. Are Airlines and Airport significant grouping factors? Which one is better?

- M1 : weights = varIdent(form = ~1 | Airline)
M2 : weights = varIdent(form = ~1 | OriginAirport)

MODEL	AIC	BIC	LOGLIK
LM	19786.86	19920.43	-9878.432
M1	19613.04	19773.31	-9788.518
M2	19709.25	19860.62	-9837.624

Step 3: Combining Variance Covariate and Grouping Factor:

- M6 : weights = varConstPower(form = ~ ArrDelay | Airline)

MODEL	AIC	BIC	LOGLIK
LM	19786.86	19920.43	-9878.432
M4	18474.51	18625.89	-9220.257
M6	18188.05	18392.85	-9071.024

Step 4: Selecting the Best Model: M6.

Coefficients:	value	p-value	standardized residuals:
(Intercept)	-0.149688	0.0000	Q1 MAX
Operating_AirlineAA	-0.056968	0.0000	Q2 MAX
Operating_AirlineAL	-0.137988	0.0000	-0.709637
Operating_AirlineAN	-0.106000	0.0000	-0.690020
ArrDelay	0.176796	0.0000	Residual standard error: 0.0214337
Distance	0.145411	0.0000	Q3 MAX
OriginATL	-0.047490	0.0000	Parameter estimates:
concurrent_TOT	0.047490	0.0000	SE AA INN DL
Operating_AirlineAA:ArrDelay	-0.047490	0.0000	const 11.600414 12.13232 13.64476 12.126284
Operating_AirlineAL:ArrDelay	-0.047490	0.0000	2.369185 2.16988 2.01094 1.845901
Operating_AirlineAN:ArrDelay	-0.047490	0.0000	ArrDelay:OriginATL
ArrDelay:OriginATL	-0.051378	1e-04	-0.051378

4. Correlation Structure

Correlation Structure Selection:

Main Problem: Many of those variables led to problems of memory allocation.

Grouping Factor Selected:

- "route ID": 238 levels
- "FlightDate": 365 levels

Why CorCompSymm: Computational problem made not possible the use of CorSymm.

We introduced the grouping factor only, not having any position variable to use.

AIC test

Model	df	AIC	BIC	logLik
MRe.1	1 23	20213.50	20418.30	-10083.749
MRe.2	2 23	20187.39	20392.19	-10070.695
MRe.3	3 23	19837.29	20042.09	-9895.644
MRe.4	4 23	11292.06	11496.86	-5623.031
MRe.5	5 23	17813.65	18018.45	-8883.826
MRe.6	6 23	19651.94	19856.74	-9802.970
MRe.7	7 23	19464.32	19669.12	-9709.159

Formulation: $\log(\text{TaxiTime}) \sim \text{Airline} + \text{ArrDelay} + \text{Distance} + \text{Origin} + \text{ConcurrentFlights} + \text{Airline:ArrDelay} + \text{Origin:ArrDelay}$

weights: varConstPower (form = ~ ArrDelay | Airline)

No correlation structure

The best model with random intercept is MRe.4., which is the one associated to the lowest AIC (11.292) and uses as grouping factor routeID.

Plot of the Random Intercepts (MRe.4)

From the plot we can see that all the random intercept are very close to 0.

A smaller sample of the dataset was used to fit the models

<