



AMERICAN AIRLINES BAG CHECK PREDICTION

Alyssa Juarez

THE GOAL

The Goal of this project is to create an algorithm
that can accurately predict the number of bags that
will be on each flight from American Airlines.





PREPROCESSING

BASIC DATA CLEANING

Remove or replace any empty values & remove duplicate rows that occur in the set.

REMOVE NULLS & DUPLICATES

01

Trim Features down to only those that are relevant for our purposes.

TRIM FEATURES

02

Find the best way to group the data so that we can see what we can call normal and how each variable is effected

DATA GROUPING

03

Remove any anomalies.

REMOVE OUTLIERS

04



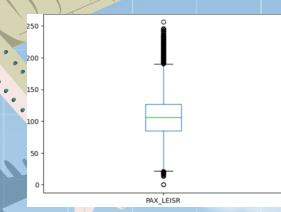
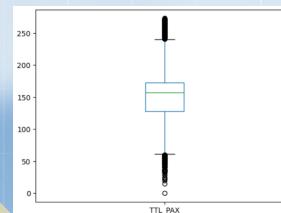
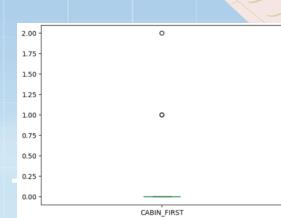
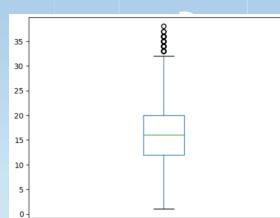
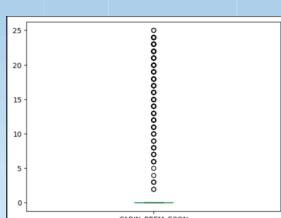
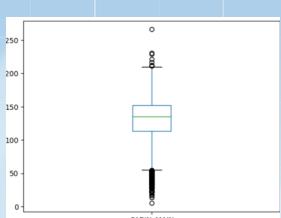
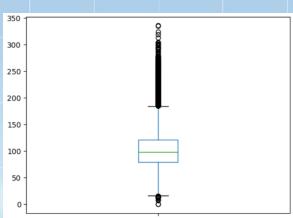
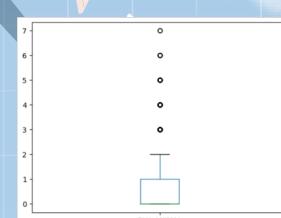
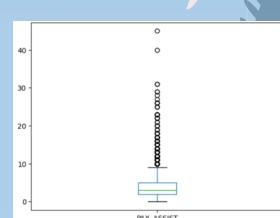
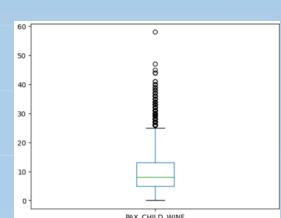
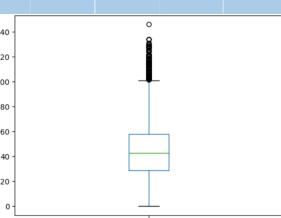
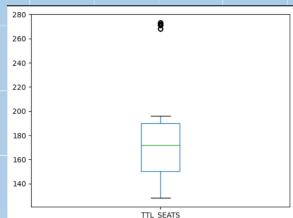
MY DATASET

The original dataset contained 74 variables, and 258,920 instances. Here are the features I found relevant

TTL_SEATS	Total number of seats available for the given aircraft type
TTL_PAX	Total number of passengers that traveled in the given flight
PAX_BUISNESS	Number of passengers travelling for business purposes booked through -third party travel agencies
PAX_LEISR	Number of passengers traveling for leisure - booked by third party, e.g., expedia
PAX_CHILD_WINF	Number of passengers travelling with child or infants
PAX_ASSIST	Number of passengers needing special assistance, e.g., wheel chair, etc.
PAX_ANIMAL	Number of passengers travelling with an animal
CABIN_FIRST	Number of passengers seating in First Class cabin
CABIN_BUISNESS	Number of passengers seating in First Class cabin
CABIN)PREM_ECON	Number of passengers seating in Premium Economy cabin
CABIN_MAIN	Number of passengers seating in Main or economy cabin
FLIGHT_ARVL_DT	Arrival date for flight departing from CLT
DEST_CITY_NM	Destination city name
TTL_BAGS	Total number of checked bags in the flight

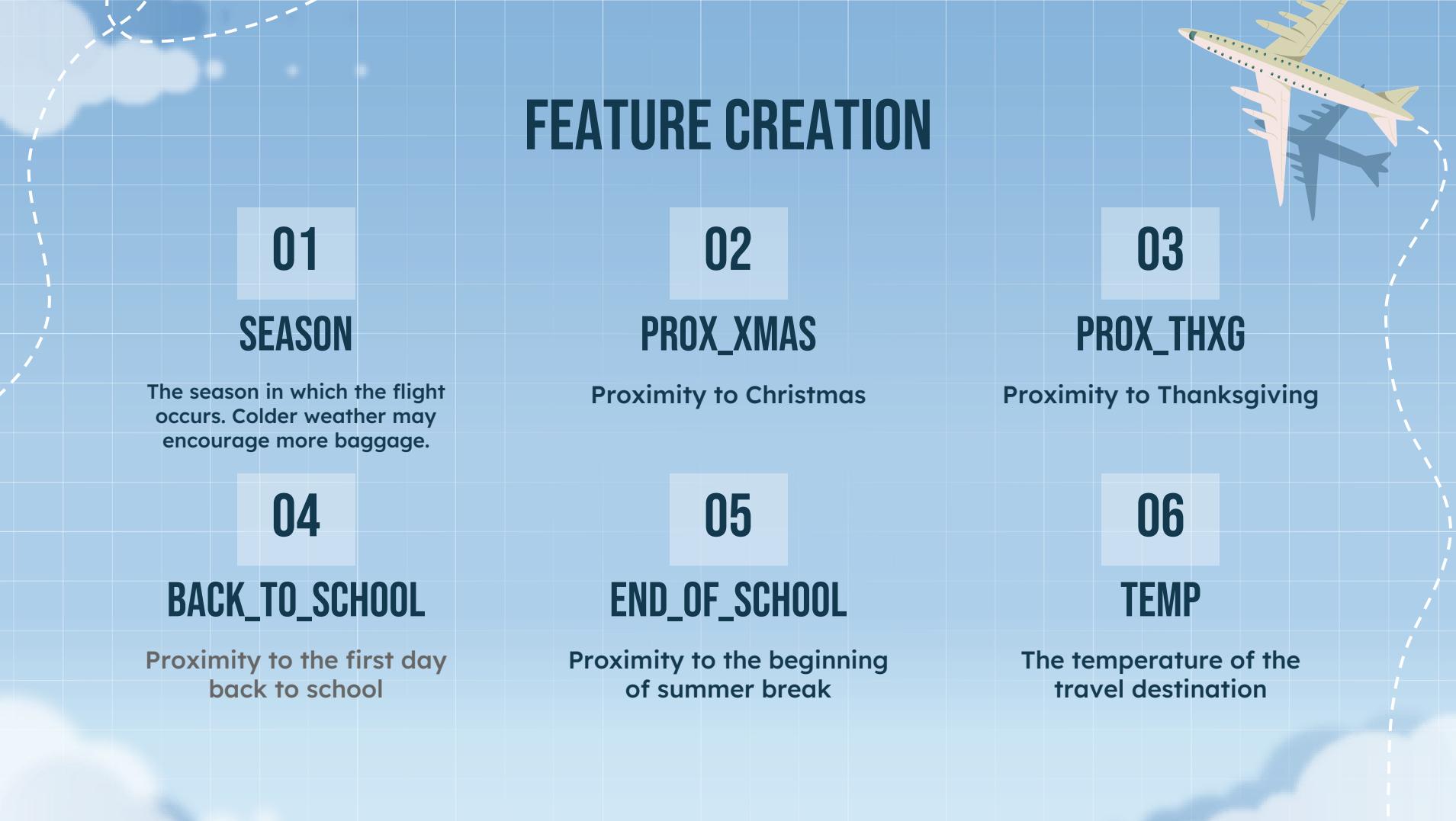
OUTLIERS

My dataset is very complicated with outliers, I will need to break apart the dataset into groups in order to get a more reliable read for what is normal. As removing outliers outright results in the removal of too many of my features.



DATA LOSS PER CLEANING STEP

Original size	258920 rows
After dropping null values	255982 rows
After dropping duplicates	33335 rows
After dropping outliers	11028 rows



FEATURE CREATION

01

SEASON

The season in which the flight occurs. Colder weather may encourage more baggage.

04

BACK_TO_SCHOOL

Proximity to the first day back to school

02

PROX_XMAS

Proximity to Christmas

05

END_OF SCHOOL

Proximity to the beginning of summer break

03

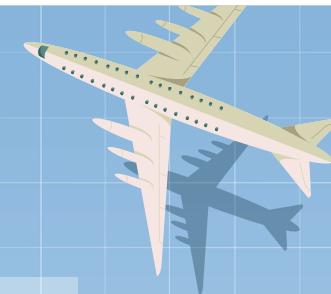
PROX_THXG

Proximity to Thanksgiving

06

TEMP

The temperature of the travel destination

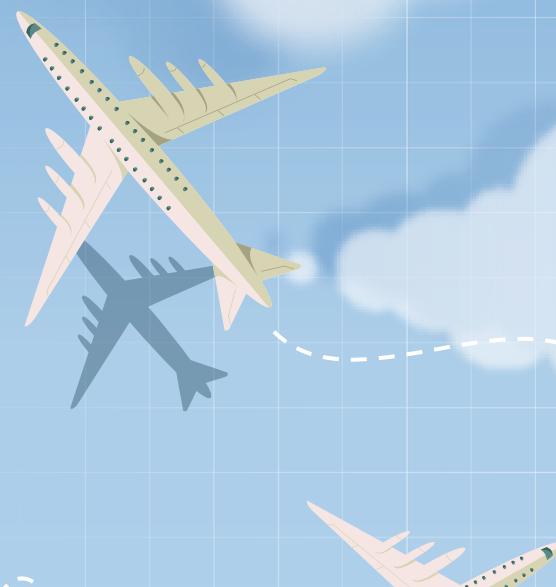


FEATURE CREATION

After some Analysis, I realized I was unable to do several of the features I would like to with this dataset given that the data only covers May to September.

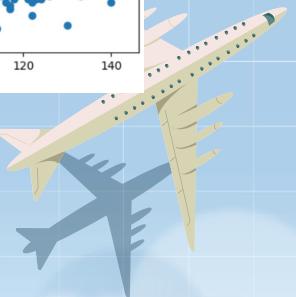
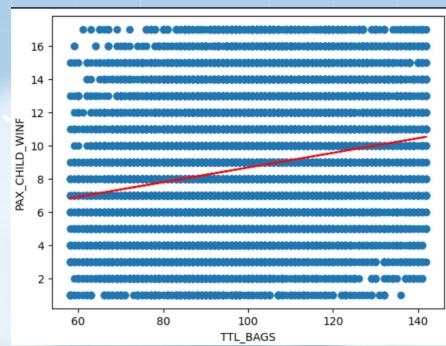
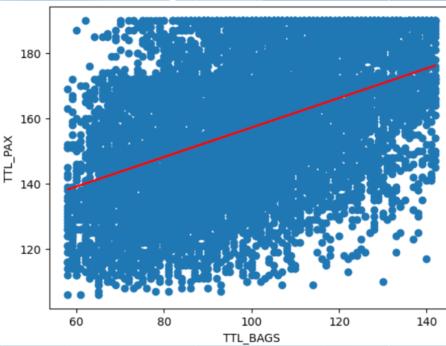
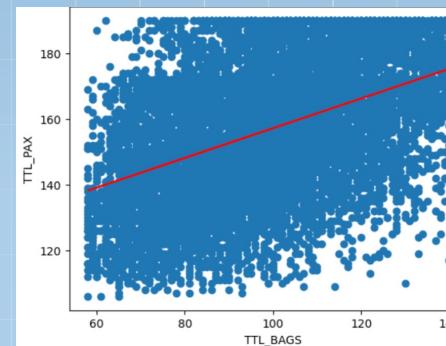
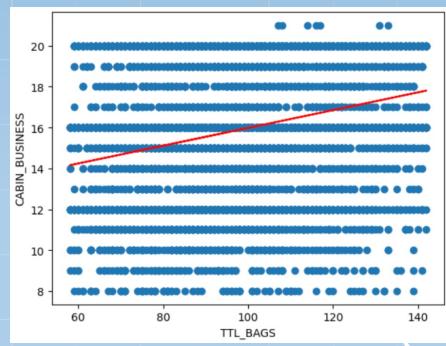
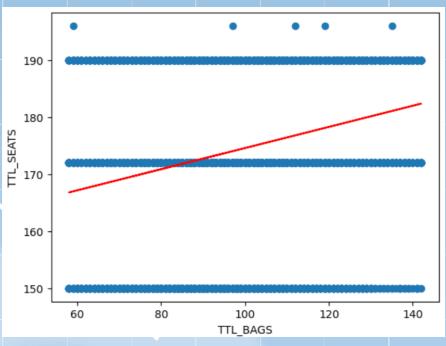
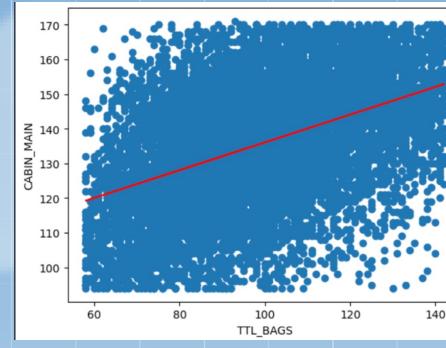
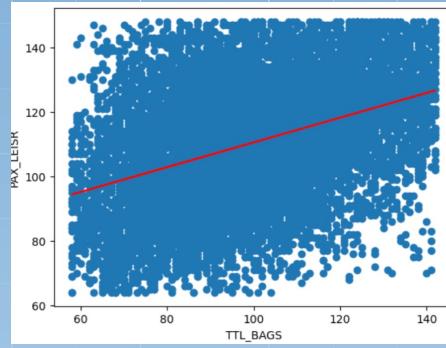
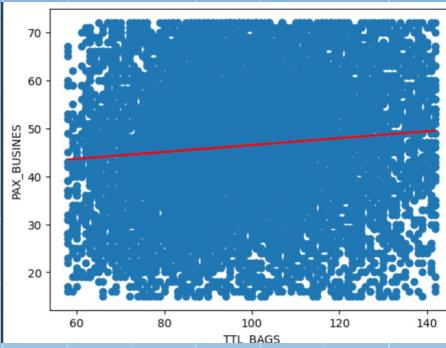
```
1 # END OF SCHOOL
2
3 df_cat['EOS1'] = '2022-05-08'
4 df_cat['EOS2'] = pd.to_datetime(df_cat['EOS1'])
5 df_cat['FAD'] = pd.to_datetime(df_cat['FLIGHT_ARVL_DT'])
6 df_cat.drop(['EOS1'], axis = 1)
7 df_cat['EOS'] = (df_cat['FAD'] - df_cat['EOS2']) / np.timedelta64(1, 'D')
8 dfcat = df_cat.drop(['EOS1', 'FAD', 'EOS2'], axis = 1)
9 dfcat
```

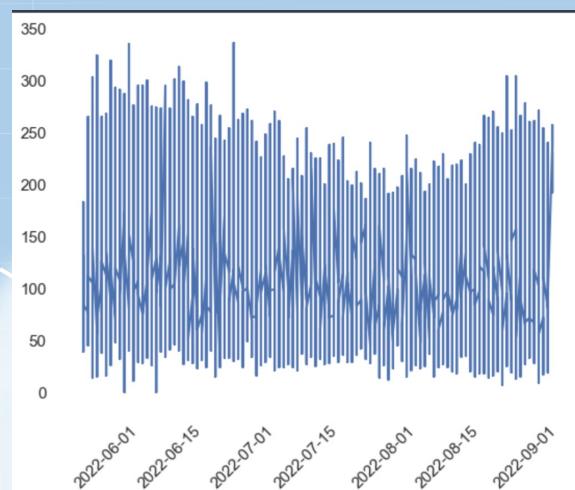
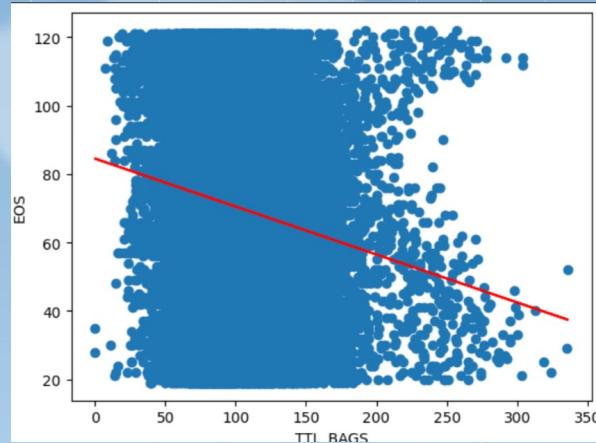
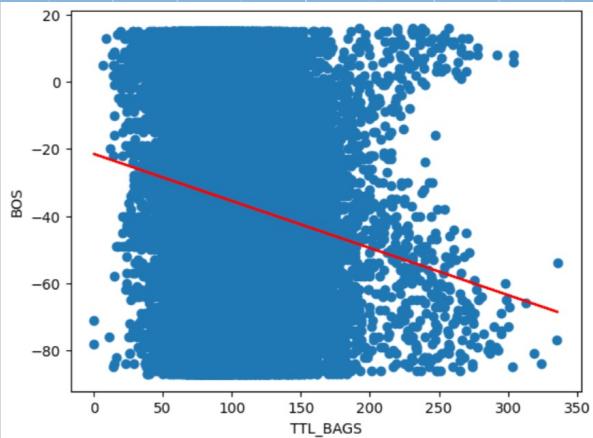
```
# BEGINNING OF SCHOOL
dfcat['BOS1'] = '2022-08-22'
dfcat['BOS2'] = pd.to_datetime(dfcat['BOS1'])
dfcat['FAD'] = pd.to_datetime(dfcat['FLIGHT_ARVL_DT'])
dfcat.drop(['BOS1'], axis = 1)
dfcat['BOS'] = (dfcat['FAD'] - dfcat['BOS2']) / np.timedelta64(1, 'D')
dfcat2 = dfcat.drop(['BOS1', 'FAD', 'BOS2'], axis = 1)
dfcat2
```

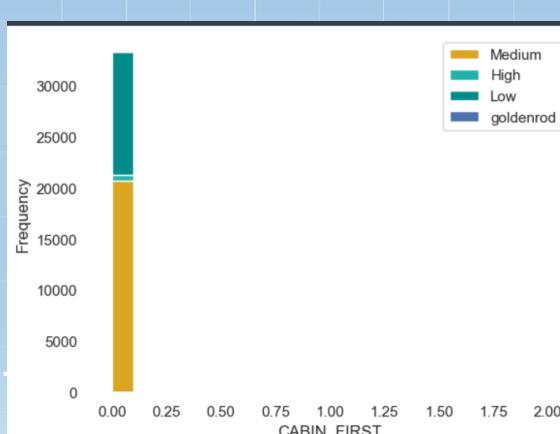
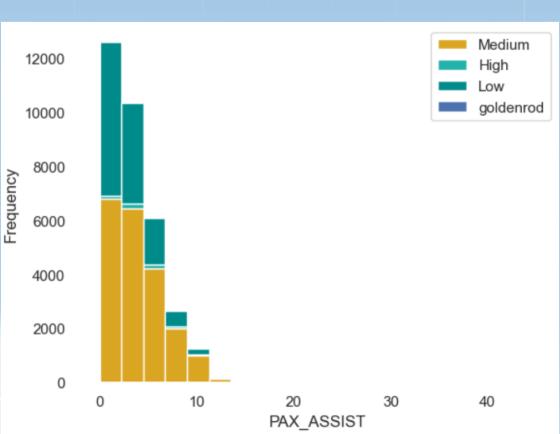
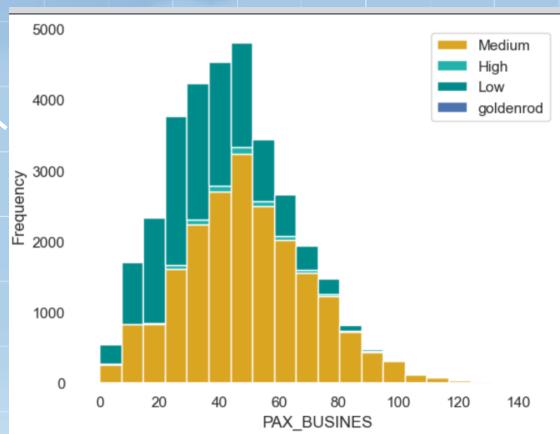
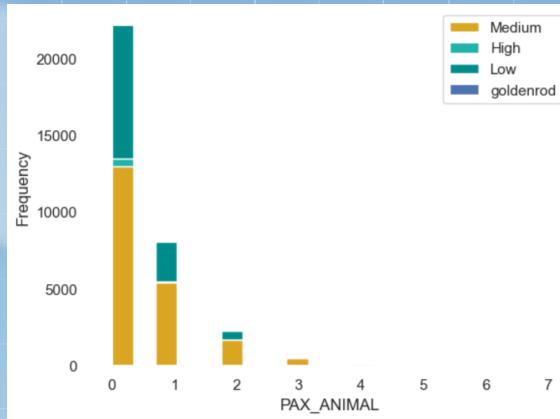
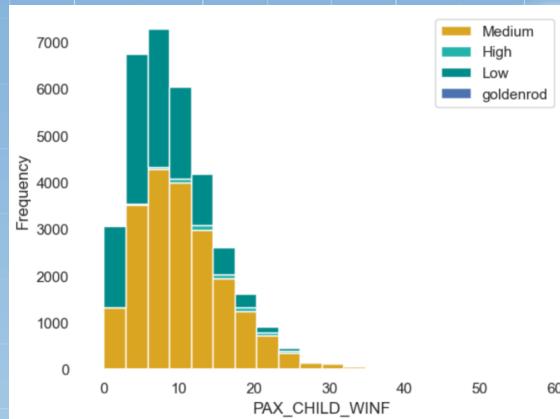
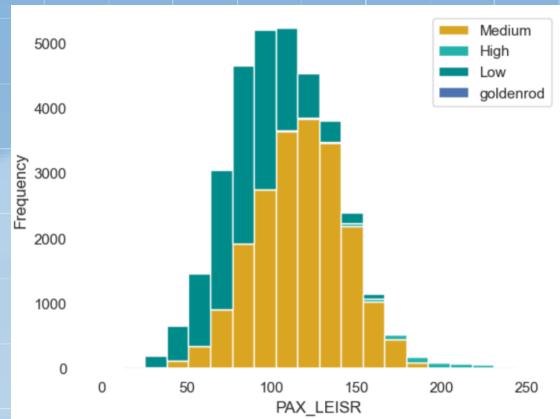


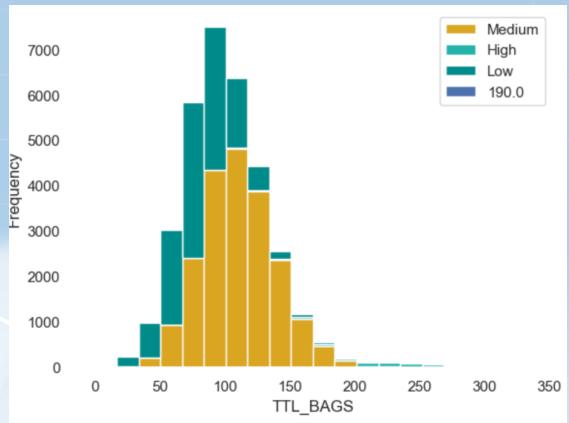
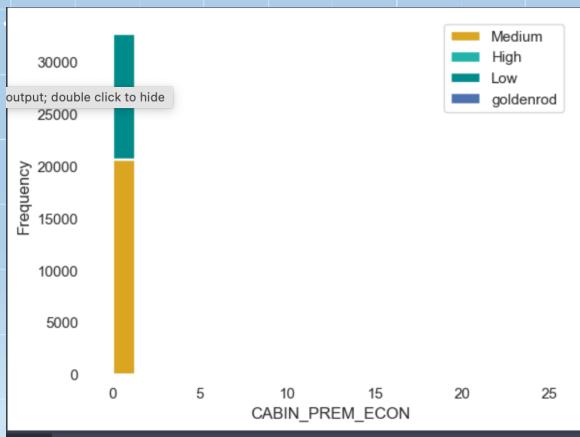
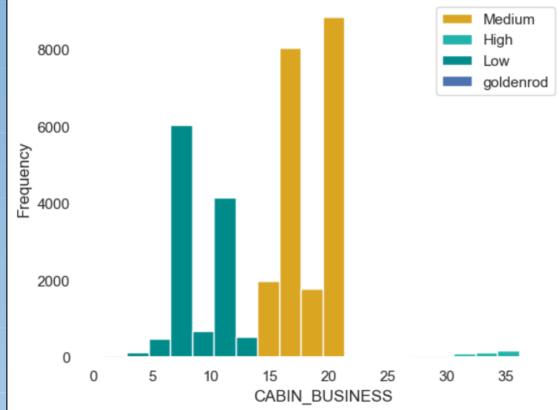
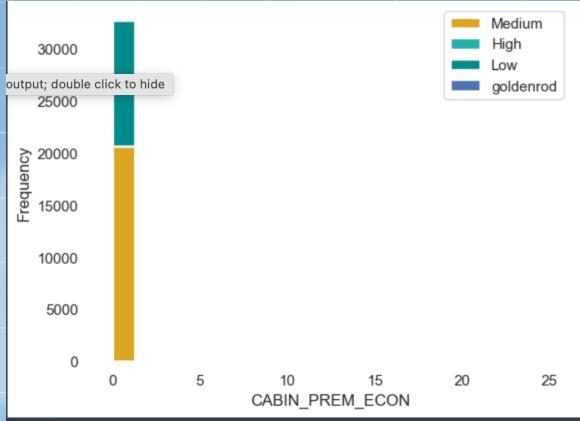


DATA ANALYSIS







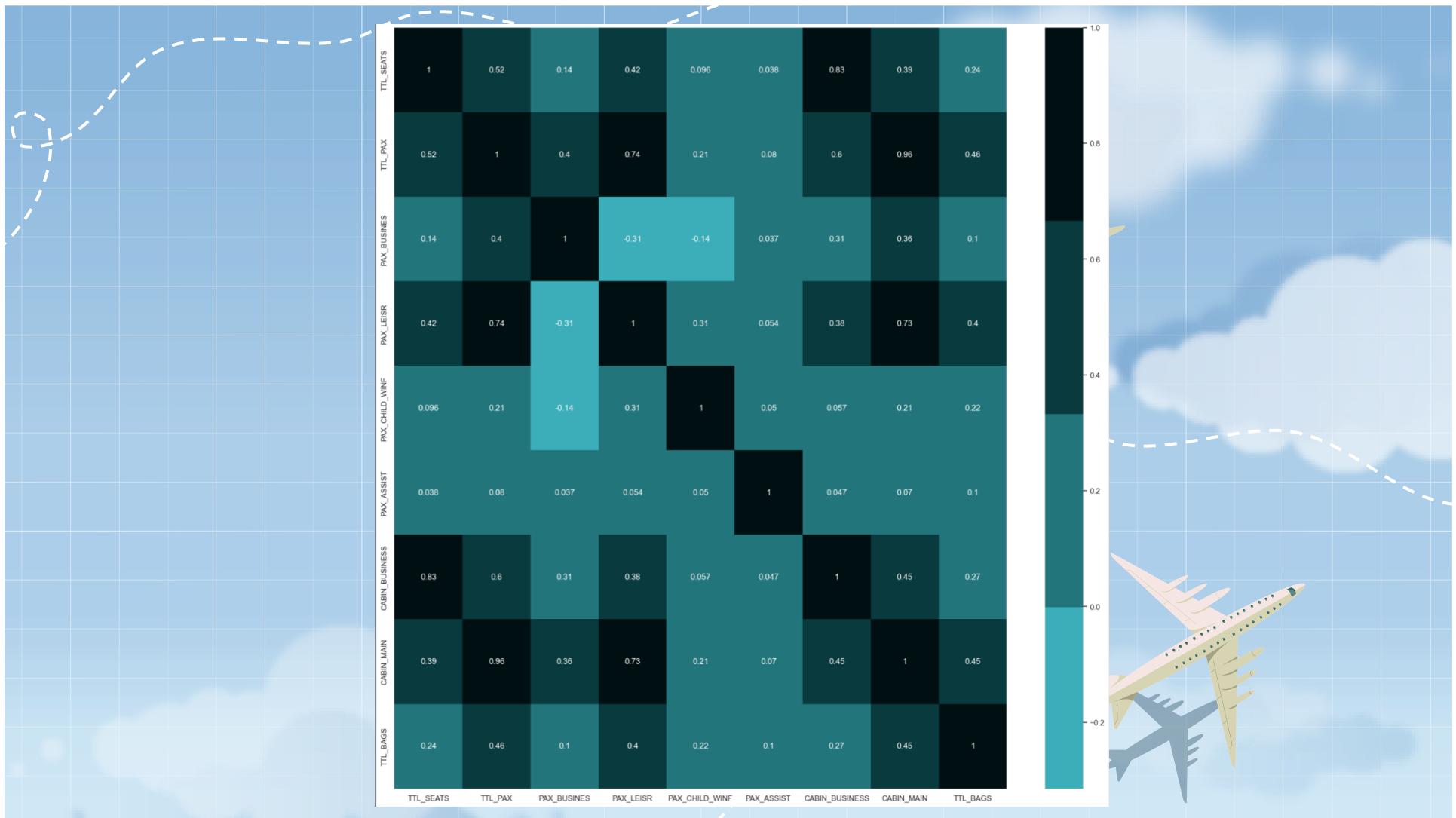


CHECKING FOR MULTICOLLINEARITY

Variance Inflation Factor (VIF) Results:

	Variable	VIF
0	TTL_SEATS	60.278522
1	TTL_PAX	5618.070894
2	PAX_BUSINES	454.811530
3	PAX_LEISR	2347.425826
4	PAX_CHILD_WINF	4.987747
5	PAX_ASSIST	3.519100
6	PAX_ANIMAL	1.524878
7	CABIN_FIRST	1.004760
8	CABIN_BUSINESS	44.406449
9	CABIN_PREM_ECON	1.803345
10	CABIN_MAIN	805.373017
11	TTL_BAGS	25.618489

	Correlation Coefficients				
	TTL_SEATS	TTL_PAX	PAX_BUSINES	PAX_LEISR	PAX_CHILD_WI
NF \					
TTL_SEATS	1.000000	0.517799	0.144910	0.417251	0.0957
25					
TTL_PAX	0.517799	1.000000	0.396037	0.743222	0.2051
18					
PAX_BUSINES	0.144910	0.396037	1.000000	-0.314204	-0.1385
41					
PAX_LEISR	0.417251	0.743222	-0.314204	1.000000	0.3099
29					
PAX_CHILD_WINF	0.095725	0.205118	-0.138541	0.309929	1.0000
00					
PAX_ASSIST	0.037917	0.079747	0.037475	0.053932	0.0502
98					
CABIN_BUSINESS	0.833595	0.601406	0.314189	0.384609	0.0568
23					
CABIN_MAIN	0.386758	0.956760	0.355939	0.730735	0.2086
56					
TTL_BAGS	0.244905	0.459196	0.101788	0.398418	0.2151
29					
	PAX_ASSIST	CABIN_BUSINESS	CABIN_MAIN	TTL_BAGS	
TTL_SEATS	0.037917	0.833595	0.386758	0.244905	
TTL_PAX	0.079747	0.601406	0.956760	0.459196	
PAX_BUSINES	0.037475	0.314189	0.355939	0.101788	
PAX_LEISR	0.053932	0.384609	0.730735	0.398418	
PAX_CHILD_WINF	0.050298	0.056823	0.208656	0.215129	
PAX_ASSIST	1.000000	0.047192	0.070195	0.103151	
CABIN_BUSINESS	0.047192	1.000000	0.453829	0.267997	
CABIN_MAIN	0.070195	0.453829	1.000000	0.446952	
TTL_BAGS	0.103151	0.267997	0.446952	1.000000	





THE NEXT STEPS

THE PLAN

Over the break I plan to finish my data analysis, I suspect I will need to create groups within the data set in order to get a more reliable view into how each variable effects the total bag count.

Next semester I would like to begin with building my model and testing its accuracy. I also would like to reach out to American Airlines to test the model on multiple Airports data.

