# Estimating flight prices from booking data

Group 7: Amy Merkelz, Jordan Stump, Laila Sultan, Nayoung Lim, Vani Kalra
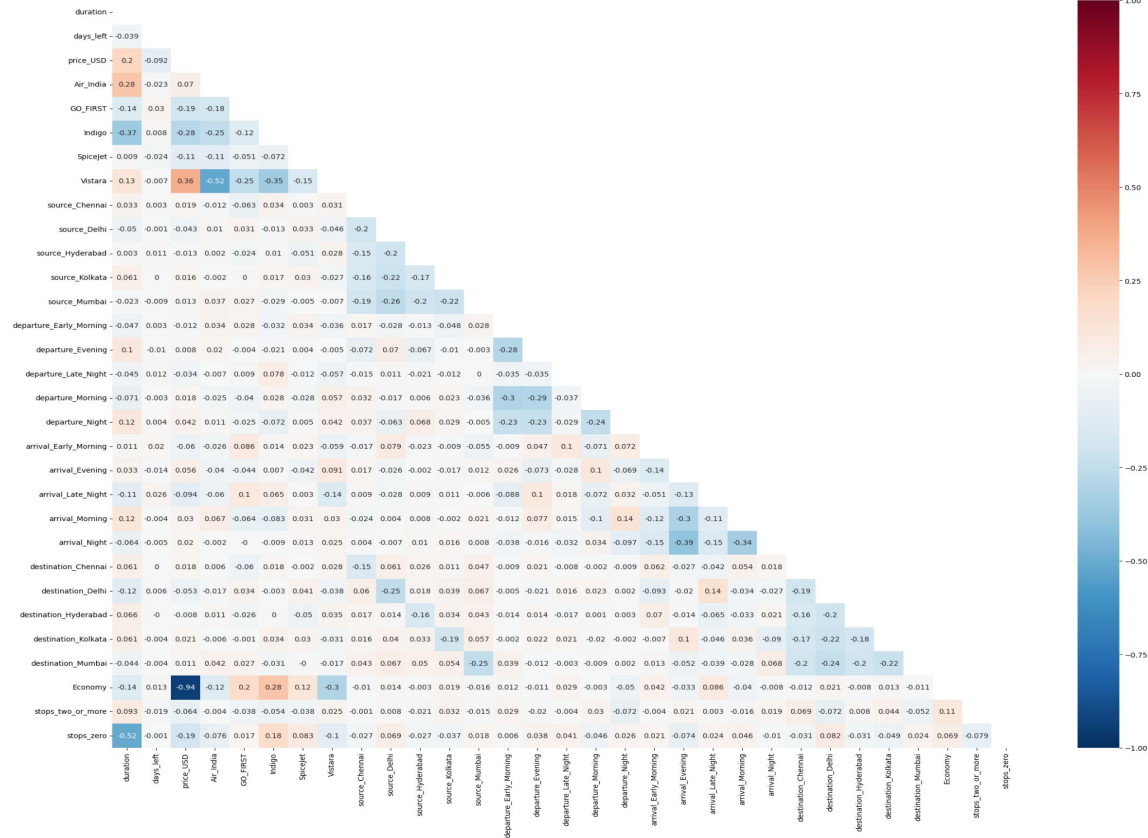
# Background

- Dataset is booking data from EaseMyTrip for 6 major metropolitan areas of India
- Data was collected from February 11th, 2022 to March 31, 2022
- Metadata:
  - 11 features, including airline, departure/arrival times and cities, days before departure that the flight was booked
  - Many categorical features that were one-hot encoded
  - 300,261 rows
- Our goal was to build a regression model to predict flight prices



publicdomainvectors.org
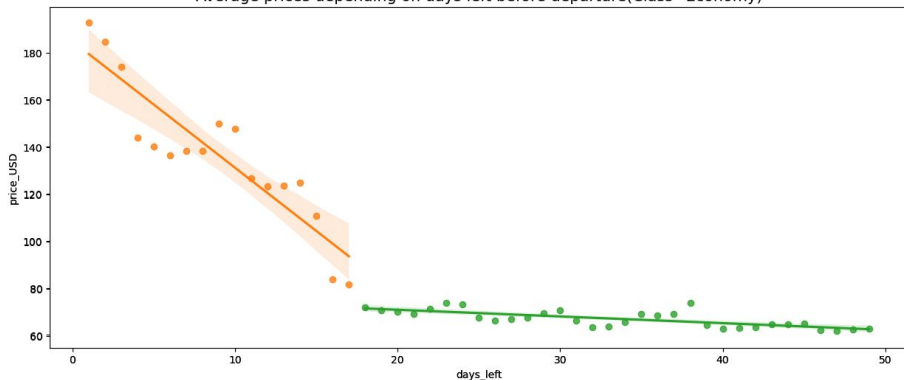
# Correlation Between Price and Other Features

Features :

- Class(Economy, Business)
- Airline
- Number of stops between source city and destination city
- Days left before departure
- Flight duration
- Departure time & Arrival time
- Duration

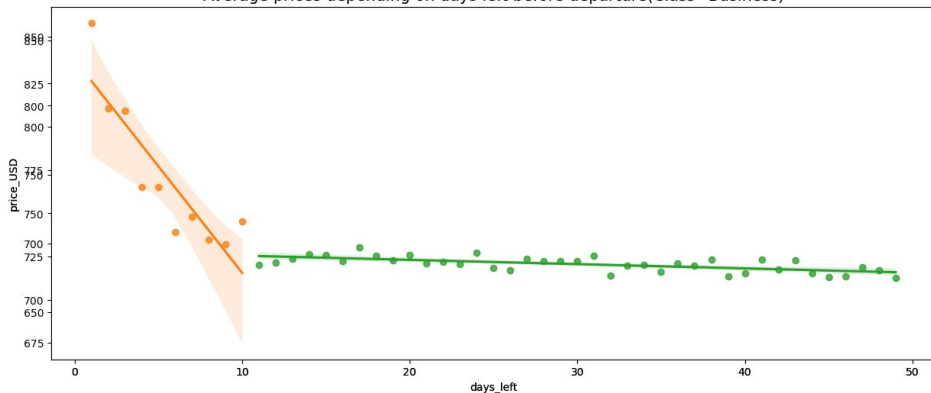# Average price depending on the number of days left before departure



Average prices depending on days left before departure(Class=Economy)

Class = **Economy**

- 18 to 50 days left before departure (green): Prices remain stable during this period.
- 1 to 18 days left before departure (orange): Prices rise starting from 18 days before departure and continue to rise.



Average prices depending on days left before departure(Class=Business)
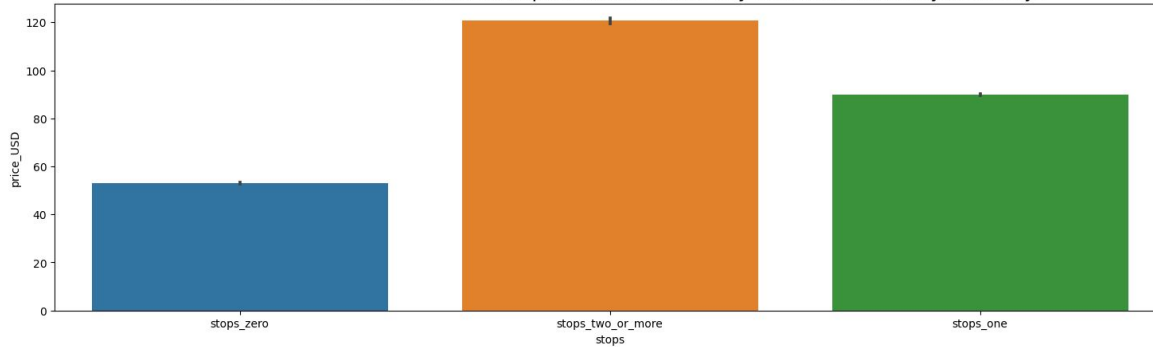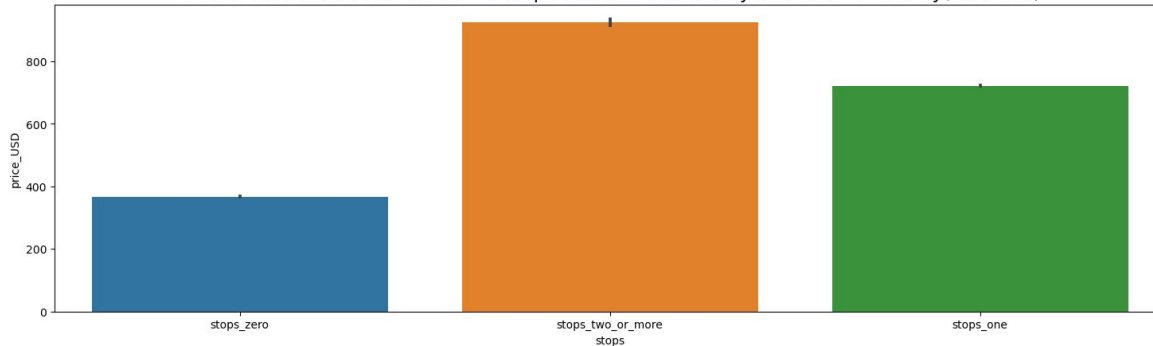
Class = **Business**

- 10 to 50 days left before departure (green): Prices remain stable during this period.
- 1 to 10 days left before departure (orange): Prices rise starting from 10 days before departure and continue to rise.

# Average price depending on number of stops between origin city and destination city


Ticket Prices based the number of stops between source city and destination city(Economy)


Ticket Prices based the number of stops between source city and destination city(Business)

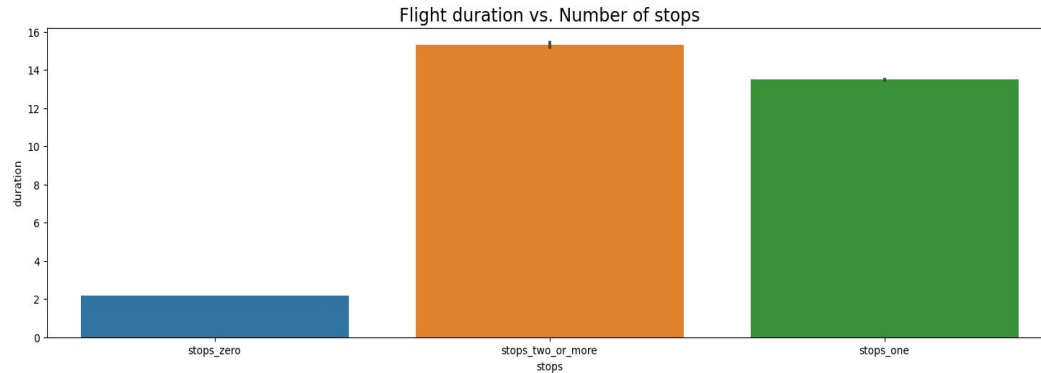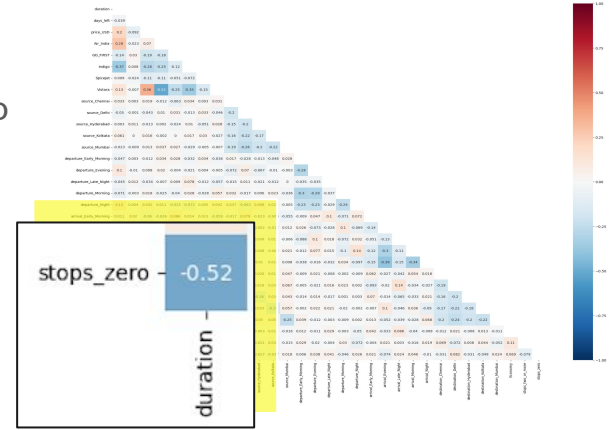"For both Economy and Business class, as the number of stops increases, the prices also rise."

# Correlation between Covariates

High correlation between covariates (abs>=0.5) : Flight duration vs. stops_zero

| | Correlation Value |
|---|---|
| Flight duration vs. stops_zero | -0.52 |





Flight duration vs. Number of stops

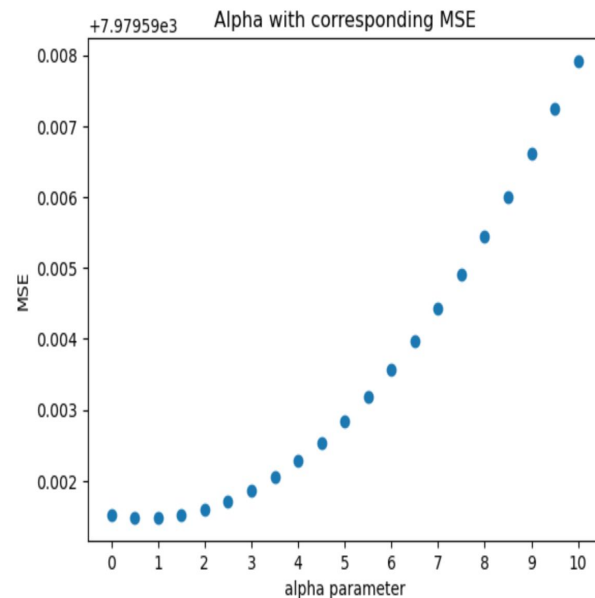"As the number of stops increases, the flight duration(time) increases."

# Ridge Regression

**Grid search** (hyperparameter tuning for alpha parameter)

- Tested values from 0 to 10 incremented by 0.5.
- Alpha of .5 provided the lowest MSE.
- Tested large alpha values, 950-1050 incremented by 5.
- Provided a slightly higher MSE and slightly lower R^squared value.
- Overall the tuning seemed to very slightly increase accuracy.

**Standard Scaler** ( for "distance" and "days_left" features)

- Decreased runtime of fit for grid search by ~40%
- Decrease in runtime of score by ~70%
- No significant impact observed on MSE or R^squared values



Alpha with corresponding MSE

|  | MSE | R^squared | Adj R^squared |
|---|---|---|---|
| **Train** | 7.98e+03 | .911 | .911 |
| **Validation** | 7.72e+03 | .913 | .913 |
| **Test** | 7.94e+03 | .911 | .911 |

# LASSO Regression

- Hyperparameter Tuning - generates 50 alpha values from {0.0001, 10} to tune model
- Best alpha - 0.000828642772854684
- Use grid search with 5-fold cross validation to find best alpha value, - using negative mean squared as the scoring metric
- Fit model on training data and uses validation and test sets for prediction
- Arrival time in the evening are the cheapest - largest coefficient of -549.7239506107213

|  | MSE | $R^2$ |
|---|---|---|
| Validation | 7721.137484452601 | 0.9129648213744022 |
| Test | 7943.1293464585 | 0.9112205609361346 |

# Decision Tree Regression: all data & random samples

| Sample Size | $R^2$ - Test | MSE - Test Data | $R^2$ - Validation | MSE - Validation |
|---|---|---|---|---|
| 100 | 0.862 | $12.3 \times 10^3$ | 0.869 | $11.7 \times 10^3$ |
| 1,000 | 0.896 | $9.30 \times 10^3$ | 0.898 | $9.06 \times 10^3$ |
| 10,000 | 0.934 | $5.87 \times 10^3$ | 0.937 | $5.60 \times 10^3$ |
| 100,000 | 0.979 | $1.91 \times 10^3$ | 0.978 | $1.95 \times 10^3$ |
| All data | 0.977 | $2.06 \times 10^3$ | 0.976 | $2.10 \times 10^3$ |

# Kernel Regression

- Hyperparameter Tuning:
  - Performed grid search for hyperparameters `alpha` and `gamma`.
  - Explored `alpha` values of 0.01 and 0.001, with `gamma` values spanning `np.logspace(-2, 0, 3)`.
  - Identified best parameters as `{'alpha': 0.01, 'gamma': 0.01}`.
- Model Evaluation:
  - On the validation set, the model achieved an MSE of 80,250.18 and an $R^2$ of 0.105.
  - On the test set, the model had an MSE of 84,477.87 and an $R^2$ of 0.098.
- Computational Details:
  - The grid search involved 3-fold cross-validation for each of 6 parameter combinations, totaling 18 fits.
- Observations:
  - The increase in sample size to 10,000 data points led to higher MSE and significantly lower $R^2$ values compared to the initial results with a smaller sample size. This indicates a decrease in model performance and suggests overfitting with the smaller sample or that the model may not generalize well to larger datasets.

# Kernel Regression

| Set | Mse | R^2 Value |
|---|---|---|
| Validation | 80,250.18 | 0.105 |
| Test | 84,477.87 | 0.098 |

# Comparison of Regression Models

| Method | $R^2$ Value - Test | Mean Squared Error - Test |
|---|---|---|
| Ordinary Linear Regression | 0.911 | $7.94 \times 10^3$ |
| LASSO Regression | 0.911 | $7.94 \times 10^3$ |
| Ridge Regression | 0.911 | $7.94 \times 10^3$ |
| Decision Tree Regression | 0.977 | $2.07 \times 10^3$ |
| Kernel Regression | 0.098 | $84.5 \times 10^3$ |

# References

Dataset: https://www.kaggle.com/datasets/shubhambathwal/flight-price-prediction/data

INR to USD Conversion: https://www.exchangerates.org.uk/

# Questions?