

Flight Price Prediction



Team 6

Abdallah Hussien

David Raafat

François Adham

Karim Mohamed

Submitted To:

Eng. **Omar Samir**

Agenda

1. Business Part
 - a. Motivation
 - b. Added Value
2. Technical Part
 - a. Preprocessing
 - b. Visualization
 - c. Insights
 - d. MapReduce
 - e. Results

The background is a solid teal color. In the top-left corner, there are three vertical bars of varying heights, each composed of three overlapping circles. In the bottom-right corner, there are four vertical bars of increasing height, each composed of three overlapping circles.

Business Part



Motivation

As stated by Harvard Business Review (HBR) travelling is one of the most common hobbies around the world as well as being essential for some careers, so we thought about making a model to book the best flight an easier decision.





Added Value

Our model is to predict flight prices knowing some data like destination, airline and class.

This model can be used in many applications like :

1. Airlines can get use of it by predicting competitors' prices and make offers to get a bigger market share.
2. Travellers can use it to know estimated flight prices for a future trip to pick the best airline and save money.



Technical Part



Preprocessing

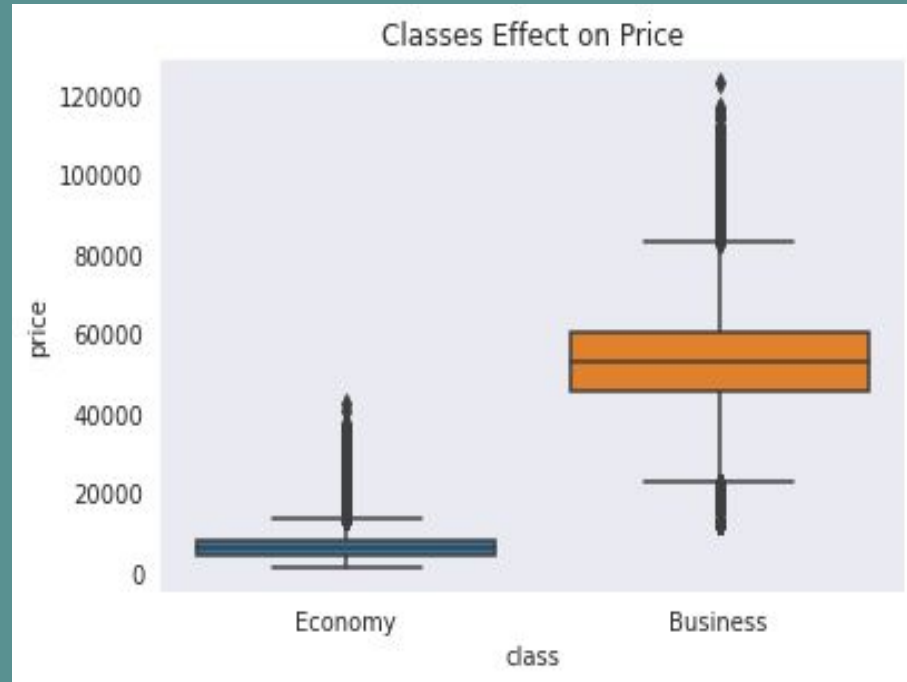
1. Explore the dataset
 - a. Check features
 - b. Check number of records
2. Check existing NAN values
3. Check how many “Categorical” feature and what are they - if exist.
4. Check how many unique values in these categorical features.
5. Drop useless features like “flight” which is the flight ID.



Visualizations

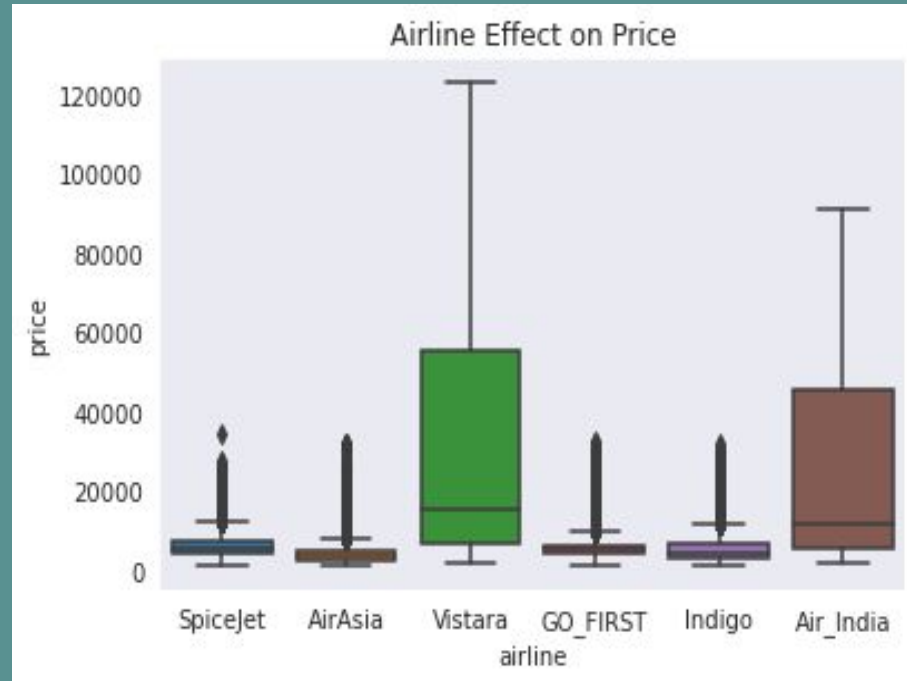


Flight class correlation with price



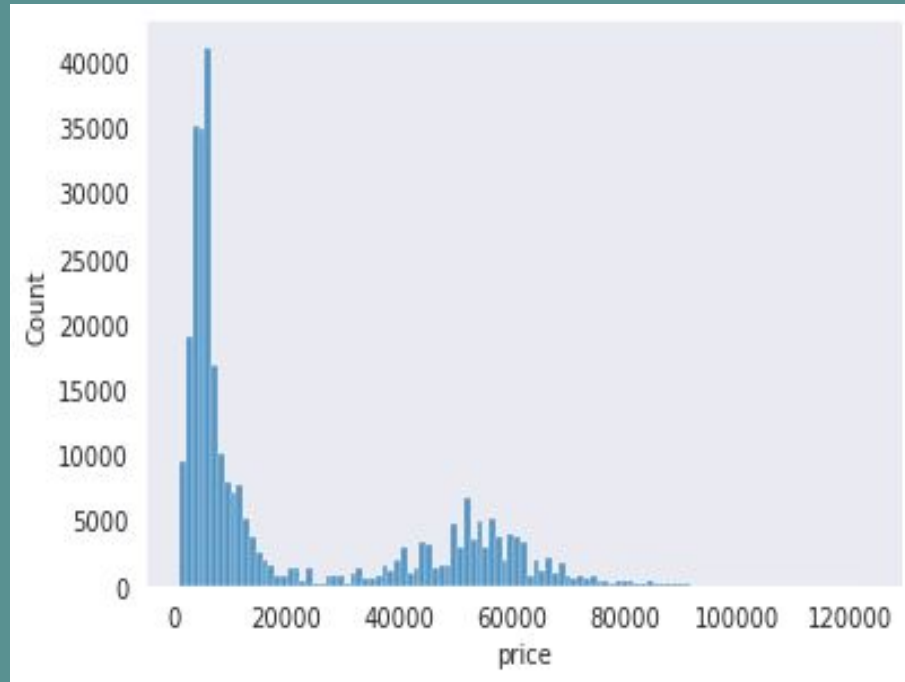


Airline correlation with price





Distribution of flights



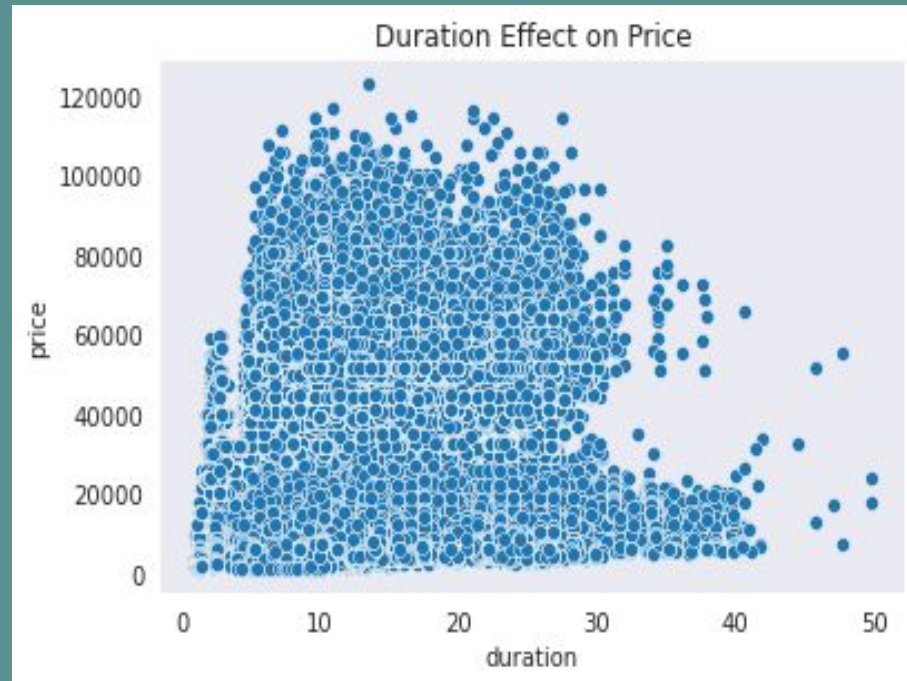


Days left correlation with price

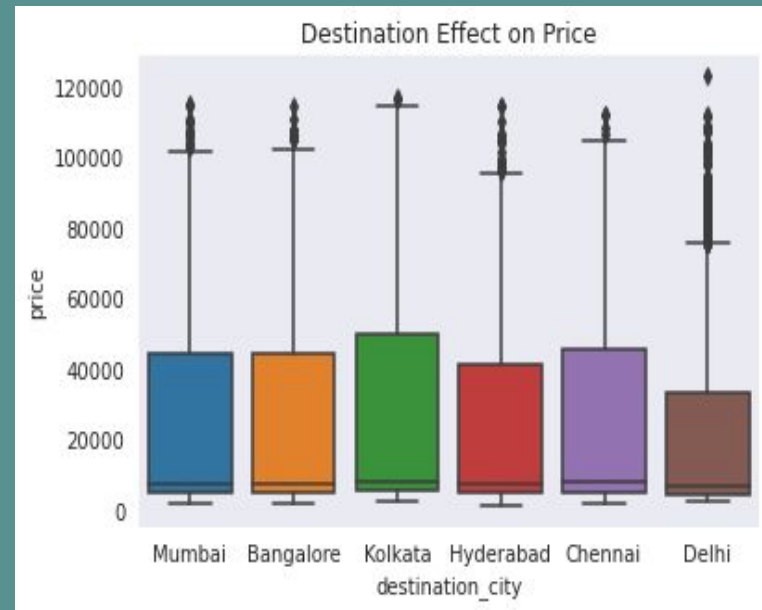
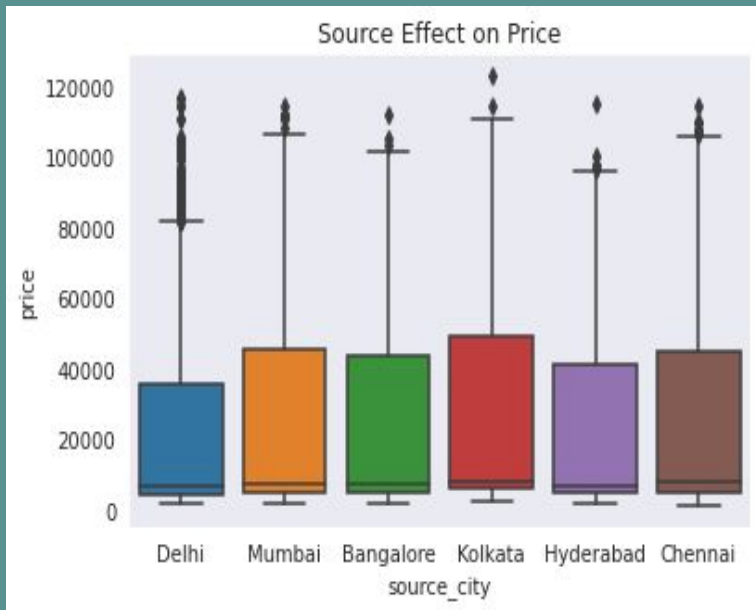




Duration effect on price

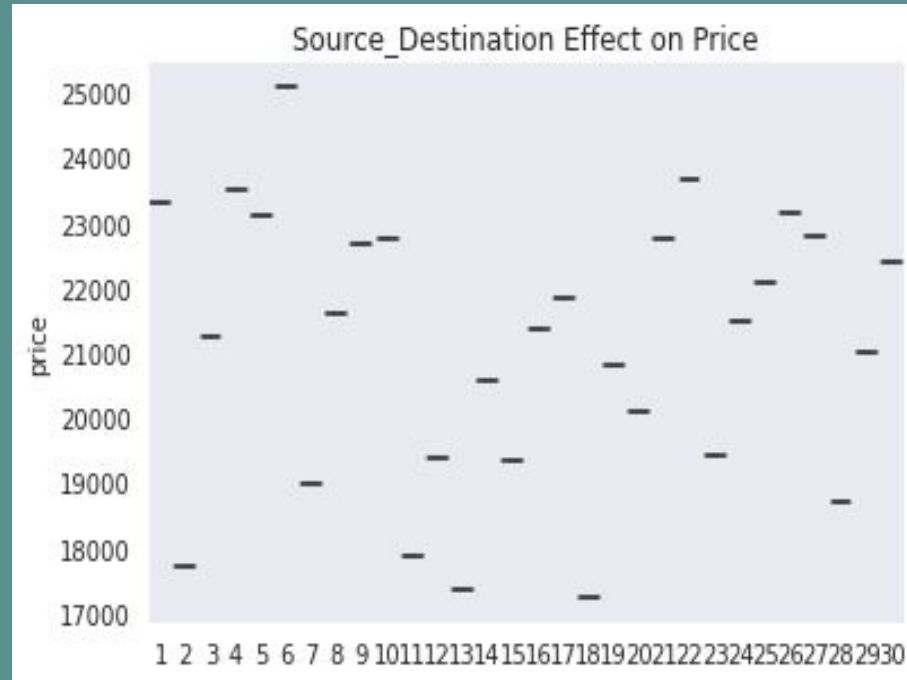


Source and Destination effect on price independently



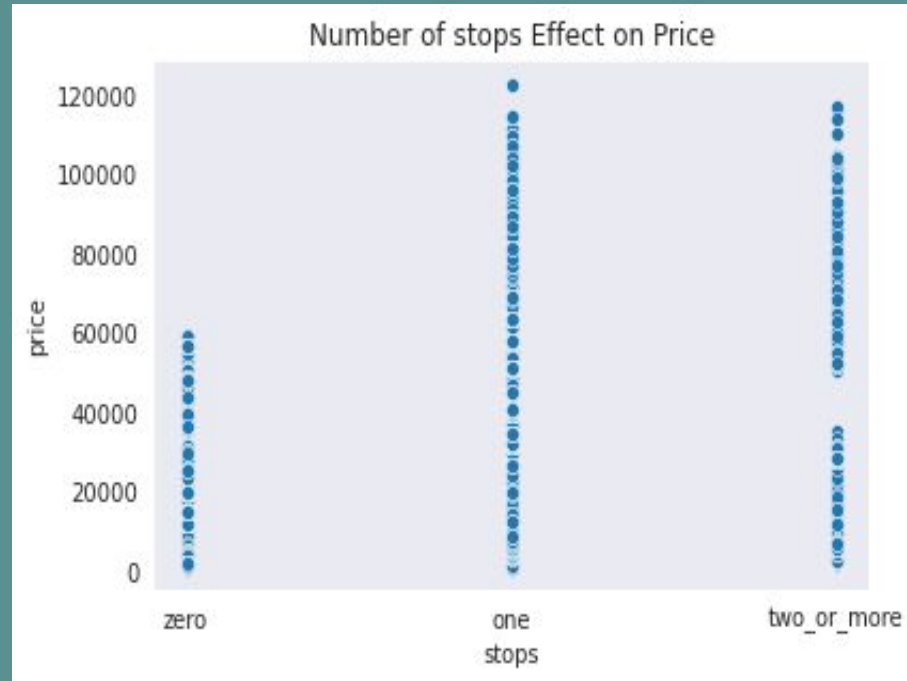


Source_Destination effect on price





Number of stops effect on price





Insights

1. Business class price is higher than economy class.
2. Number of cheap flights is higher than number of expensive flights.
3. Less “days_left” means higher price.
4. More “stops” means higher price.
5. Source and destination do not have a significant effect on the price, but pairs of (source, destination) has a higher effect.



MapReduce

Split the dataset to m splits

a. Map phase:

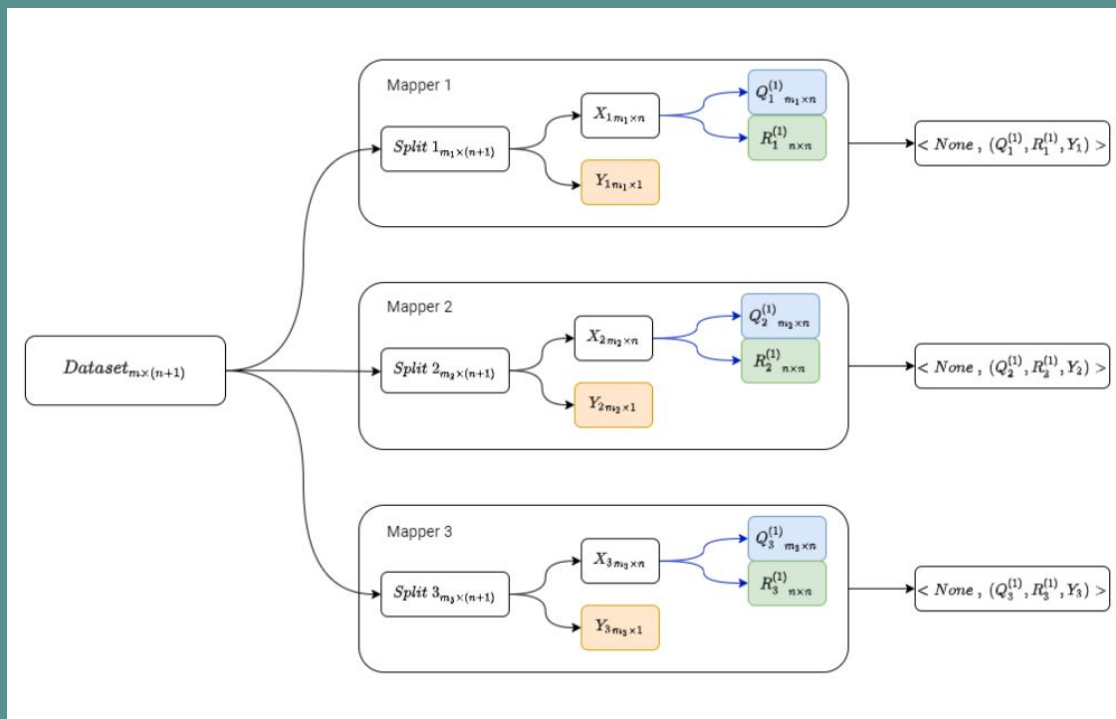
- i. Each mapper will decompose the feature vector using QR decomposition i.e. $X_j = Q_j^{(1)}R_j^{(1)}$ where j is the block number.
- ii. Map phase output will be a key-value pair where key will be NONE and the value will be (Q, R)

b. Reduce phase:

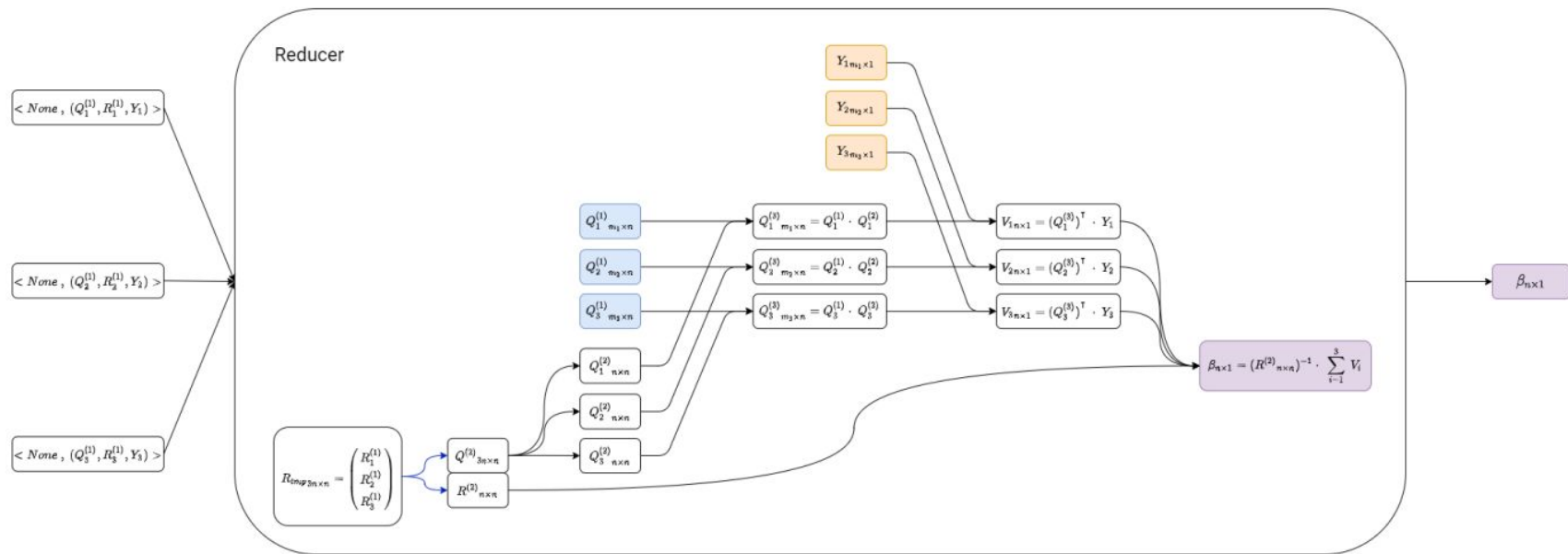
- i. Decompose each R matrix using QR decomposition i.e. $R_j^{(1)} = Q_j^{(2)}R_j^{(2)}$
- ii. Compute $Q^{(3)}$ where $Q_j^{(3)} = Q_j^{(2)}Q_j^{(1)}$
- iii. Compute V_j where $V_j = Q_j^{(3)}y_j$ where y_j is a vector of the predicted values
- iv. Compute the weights for each feature $\beta = [R^{(2)}]^{-1} \sum V_j$



Mapper



Reducer





Results

We used the coefficient of determination equation to compute our model score which is

Coefficient of Determination (R Square)

$$R^2 = \frac{SSR}{SST}$$

$$SSR = \sum_i (\hat{y}_i - \bar{y})^2$$

$$SST = \sum_i (y_i - \bar{y})^2$$

Where,

- SSR is Sum of Squared Regression also known as variation explained by the model
- SST is Total variation in the data also known as sum of squared total
- y_i is the y value for observation i
- \bar{y} is the mean of y value
- \hat{y}_i is predicted value of y for observation i



Results contd.

We tried to remove the duration feature but our model's score dropped a bit so we returned it back.

Dropping source and destination affected a lot our model's score so we returned them back.

Final results:

- Train data: **90.97%**
- Test data: **91.09%**



Thank You