

ABOUT DATASET

INTRODUCTION

The objective of the study is to analyse the flight booking dataset obtained from “Ease My Trip” website and to conduct various statistical hypothesis tests in order to get meaningful information from it. The 'Linear Regression' statistical algorithm would be used to train the dataset and predict a continuous target variable. 'Easemytrip' is an internet platform for booking flight tickets, and hence a platform that potential passengers use to buy tickets. A thorough study of the data will aid in the discovery of valuable insights that will be of enormous value to passengers.

RESEARCH QUESTIONS

The aim of our study is to answer the below research questions:

- a) Does price vary with Airlines?
- b) How is the price affected when tickets are bought just 1 or 2 days before departure?
- c) Does ticket price change based on the departure time and arrival time?
- d) How does the price change with change in Source and Destination?
- e) How does the ticket price vary between Economy and Business class?

DATA COLLECTION AND METHODOLOGY

Octoparse scraping tool was used to extract data from the website. Data was collected in two parts: one for economy class tickets and another for business class tickets. A total of 300261 distinct flight booking options were extracted from the site. Data was collected for 50 days, from February 11th to March 31st, 2022.

Data source was secondary data and was collected from Ease my trip website.

DATASET

Dataset contains information about flight booking options from the website Easemytrip for flight travel between India's top 6 metro cities. There are 300261 data points and 11 features in the cleaned dataset.

FEATURES

The various features of the cleaned dataset are explained below:

- 1) Airline:** The name of the airline company is stored in the airline column. It is a categorical feature having 6 different airlines.
- 2) Flight:** Flight stores information regarding the plane's flight code. It is a categorical feature.

3) Source City: City from which the flight takes off. It is a categorical feature having 6 unique cities.

4) Departure Time: This is a derived categorical feature created by grouping time periods into bins. It stores information about the departure time and has 6 unique time labels.

5) Stops: A categorical feature with 3 distinct values that stores the number of stops between the source and destination cities.

6) Arrival Time: This is a derived categorical feature created by grouping time intervals into bins. It has six distinct time labels and keeps information about the arrival time.

7) Destination City: City where the flight will land. It is a categorical feature having 6 unique cities.

8) Class: A categorical feature that contains information on seat class; it has two distinct values: Business and Economy.

9) Duration: A continuous feature that displays the overall amount of time it takes to travel between cities in hours.

10) Days Left: This is a derived characteristic that is calculated by subtracting the trip date by the booking date.

11) Price: Target variable stores information of the ticket price.