



FLIGHT PRICE PREDICTION PROJECT

Submitted by:

ANISH ANTONY

ABSTRACT:

Traveling in an airplane is one of the most exciting and expensive modes of travel. Since it is expensive, we must plan well accordingly. In this project we will see how the ticket vary based on advance booking of the flight and also how the price varies as the dates come closer. In this project we will be predicting the flight price from machine learning models. To predict this model, we need to scrape the flight booking details for two months from online websites using selenium web driver, and the data have preprocessed, trained and tested using the regression algorithms. Then it is hypertuned and best algorithm with best parameters is obtained and finally the ticket price of the flight is predicted.

Keywords: Flight ticket ,Selenium ,Data cleaning ,Ticket Price ,Regression

The features used in the dataset are:

- ❖ Flight Company
- ❖ Flight No
- ❖ Flight Class
- ❖ Flight Model
- ❖ Travel
- ❖ Departure Time
- ❖ Departure Day
- ❖ Departure Date
- ❖ Departure Month
- ❖ Departure Duration
- ❖ Arrival Time
- ❖ Arrival Day
- ❖ Meal
- ❖ Flight Ticket Price
- ❖ Luggage
- ❖ Connecting Planes
- ❖ No of days

Data Cleaning:

- ❖ Flight Company – Grouped as str
- ❖ Flight No – Grouped as str
- ❖ Flight Class – Grouped as str
- ❖ Flight Model – Grouped as str
- ❖ Travel – Combined source and destination grouped as str
- ❖ Departure Time – Converted to hours
- ❖ Departure Day – Grouped as str
- ❖ Departure Date – Grouped as str
- ❖ Departure Month – Grouped as str
- ❖ Departure Duration – Converted to minutes as int
- ❖ Arrival Time – Grouped as str
- ❖ Arrival Day – Grouped as str
- ❖ Meal – Grouped as str
- ❖ Flight Ticket Price – Grouped as int
- ❖ Luggage – Grouped as str
- ❖ Connecting Planes – Grouped as str
- ❖ No of days – Grouped as int

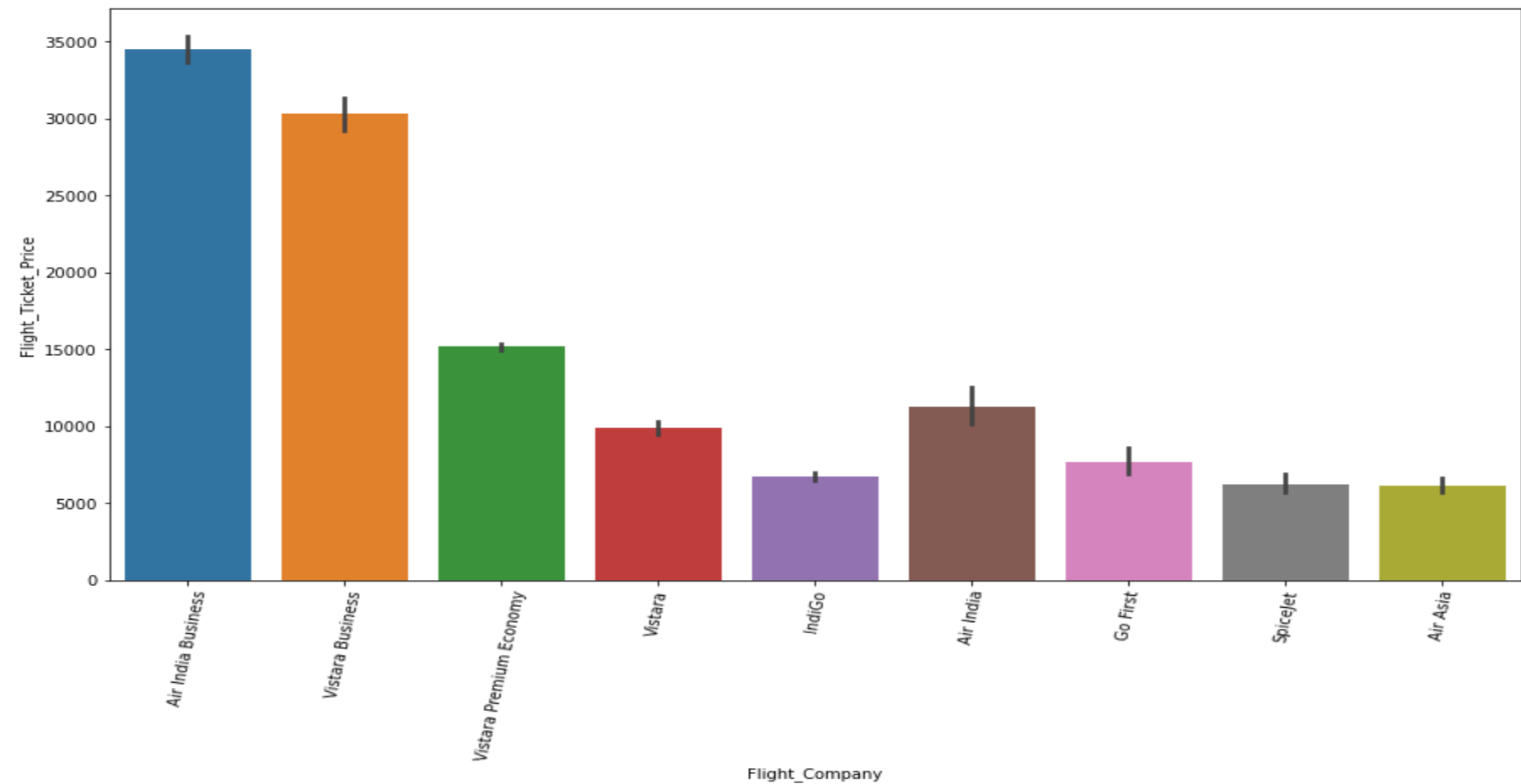
Dataset Description:

Columns	Datatype	Unique values	Mode Mean	Values
1	Flight Company	object	9	Vistara Premium Economy
2	Flight No	object	446	UK832
3	Flight Class	object	3	Economy
4	Flight Model	object	19	Vistara, Airbus A32-100
5	Travel	object	6	Bengaluru to Mumbai
6	Departure Time	object	22	7
7	Departure Day	object	7	Sat
8	Departure Date	object	28	29
9	Departure Month	object	2	Aug
10	Departure Duration	int32	68	126.74935
11	Arrival Time	int32	23	15.164451
12	Arrival Day	object	7	Sat
13	Meal	object	2	Free Meal
14	Flight Ticket Price	int64	652	19838.348439
15	Luggage	object	5	15
16	Connecting Planes	object	3	Has one connecting plane
17	No of days	int32	5	19.52242

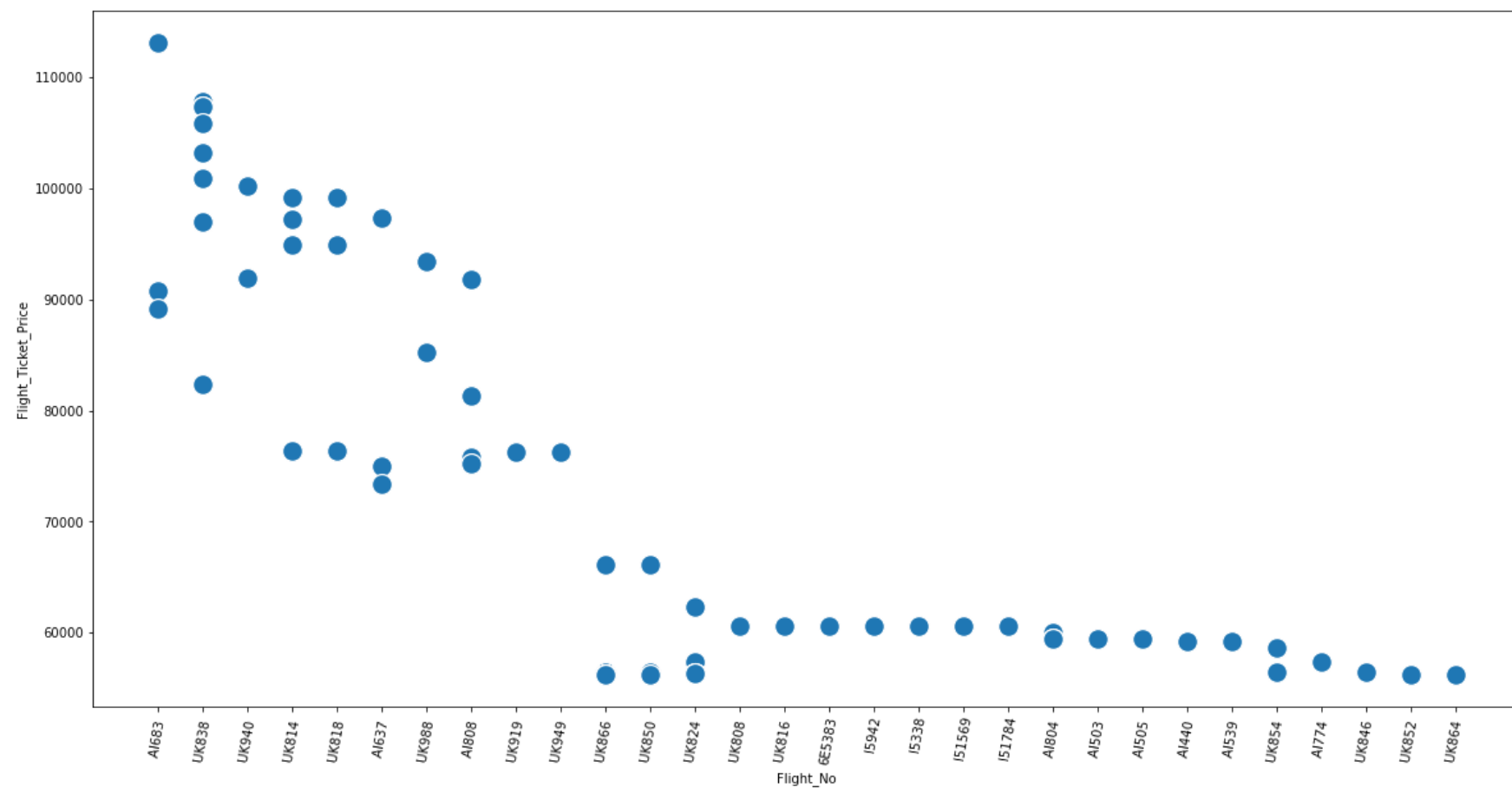
Hardware and Software Requirements and Tools Used:

- Hardware – PC Windows 10, 4 GB Ram
- Software – Google chrome, MS Excel, Python, Selenium webdriver
- Libraries – Pandas, NumPy, Matplotlib, Seaborn, sklearn, SciPy. Stats
 - ☐ Browsing – Google Chrome
 - ☐ Webscraping – Python, Selenium webdriver
 - ☐ Data cleaning – Python, Pandas, NumPy & SciPy. Stats
 - ☐ Data visualization – Matplotlib & Seaborn
 - ☐ Machine learning – Sklearn

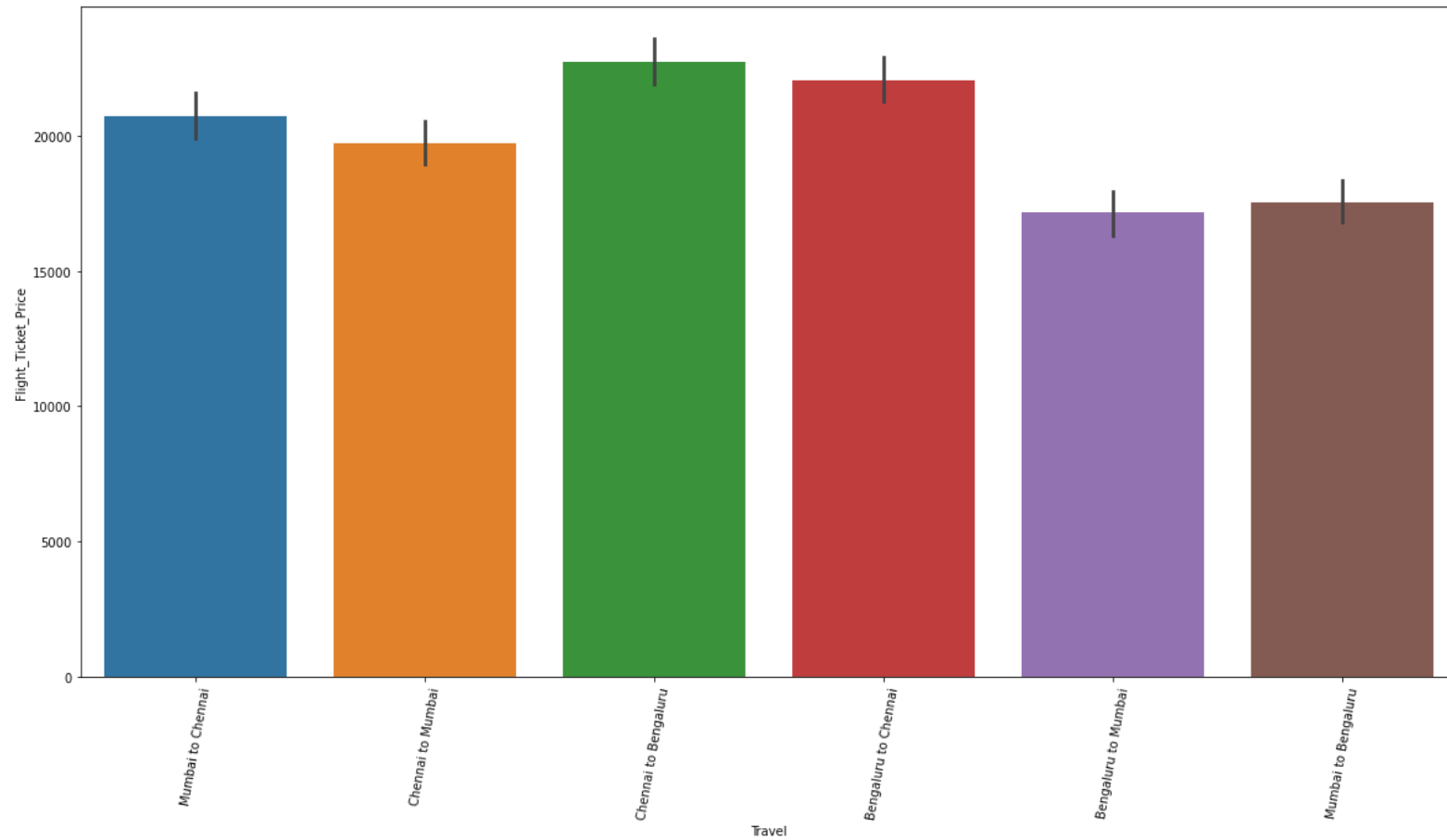
Flight Company :



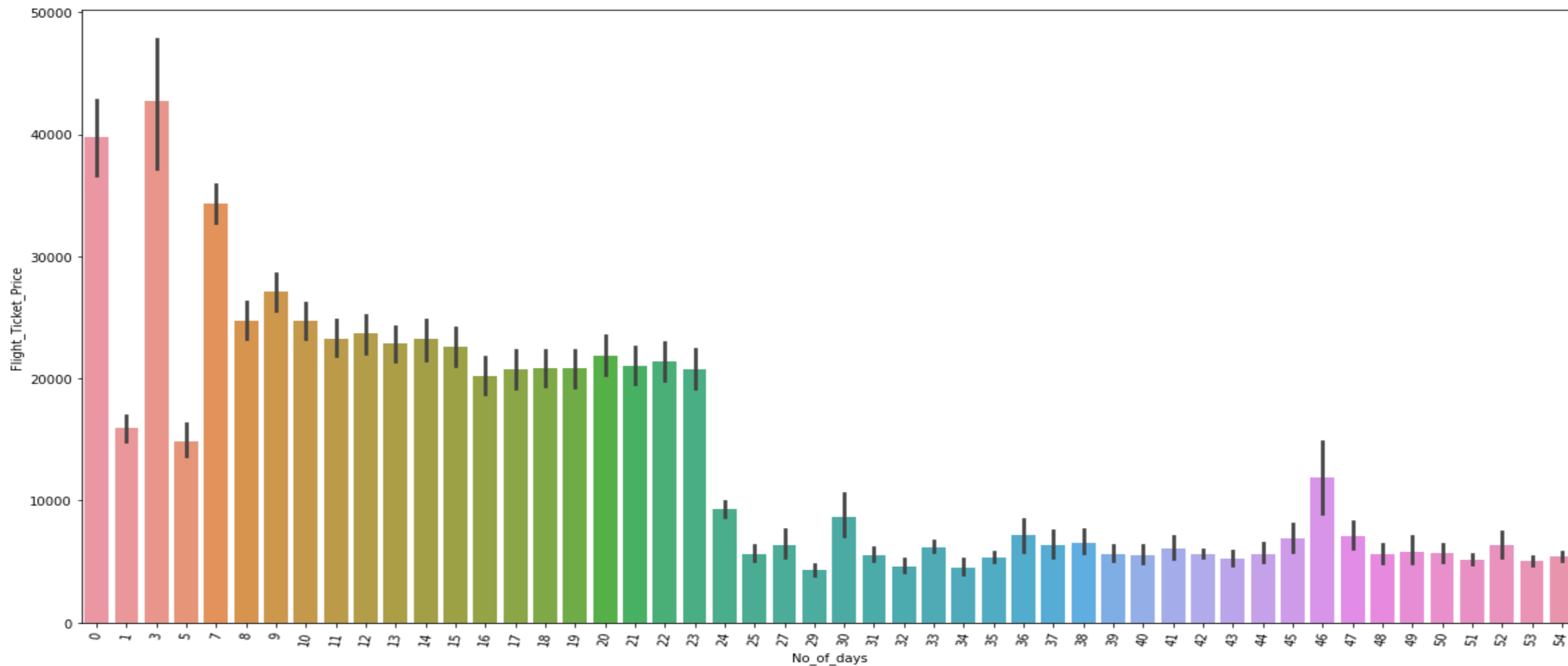
Flight No:



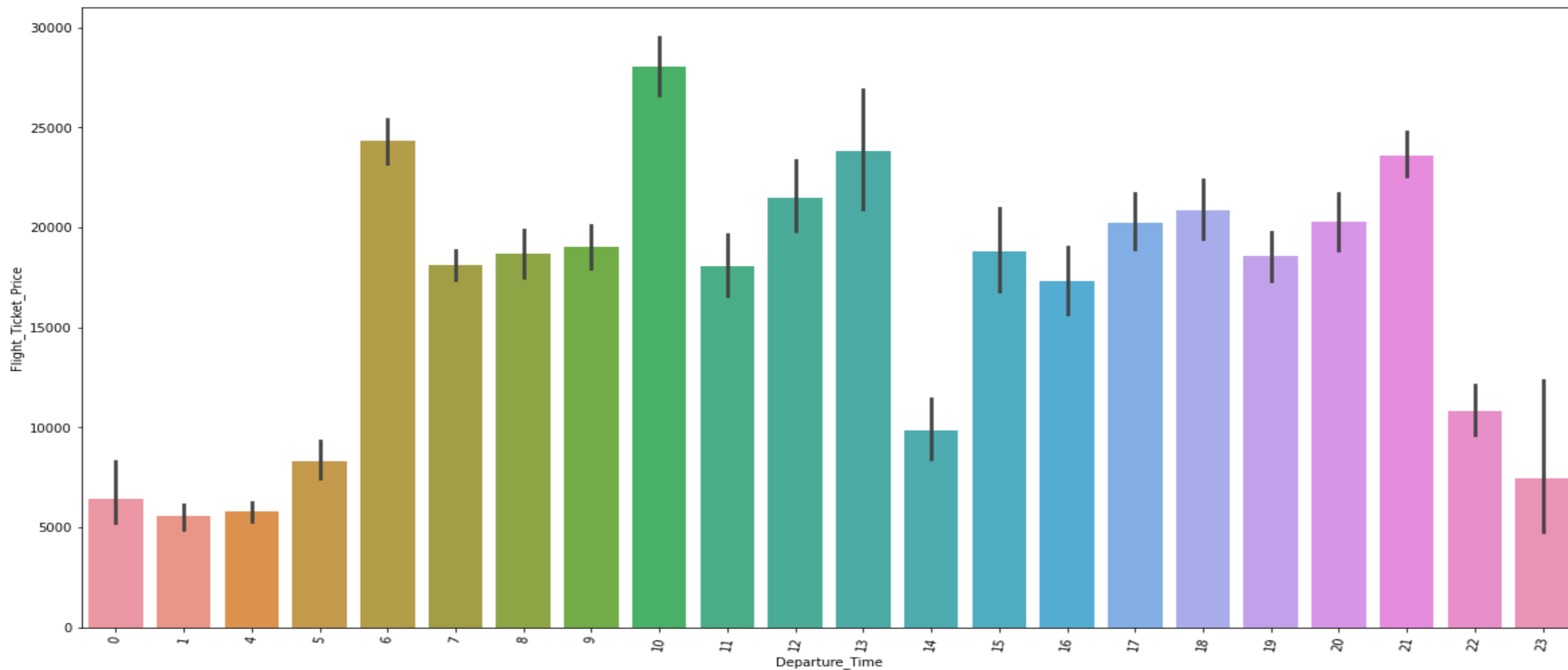
Travel:



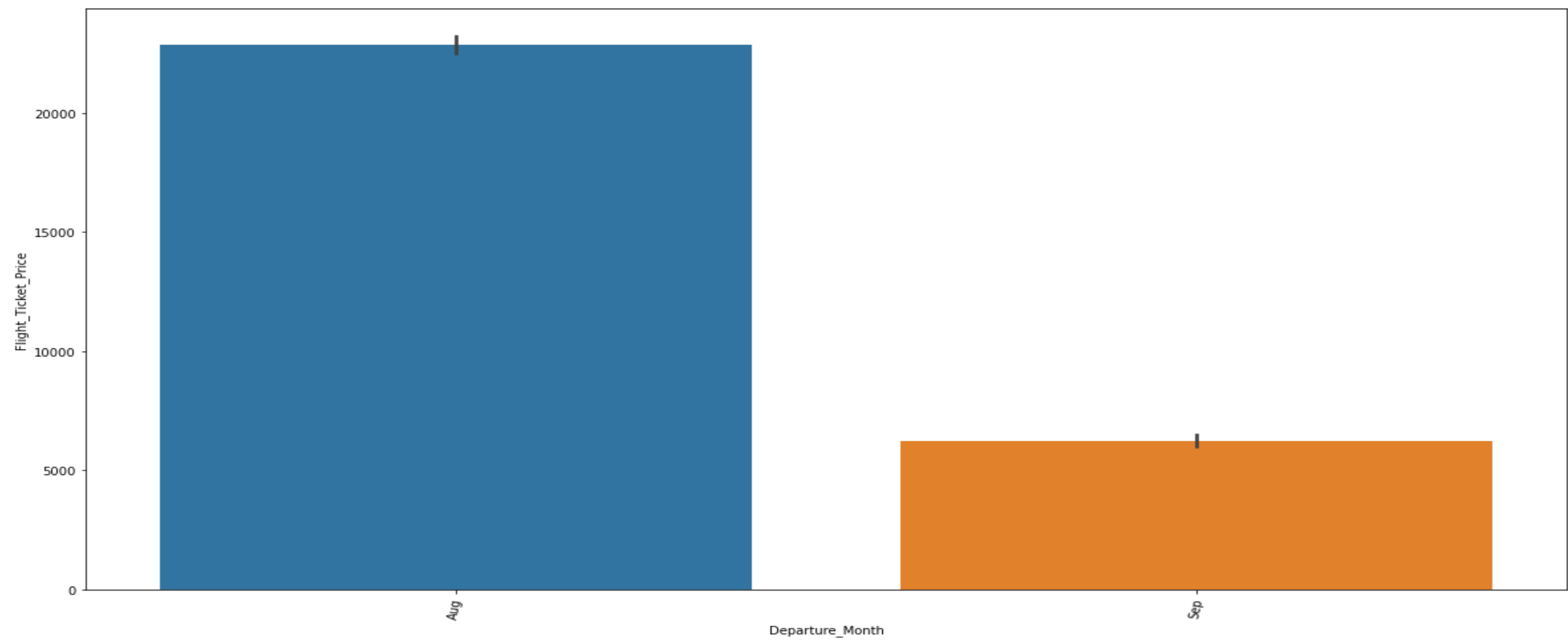
No of days:



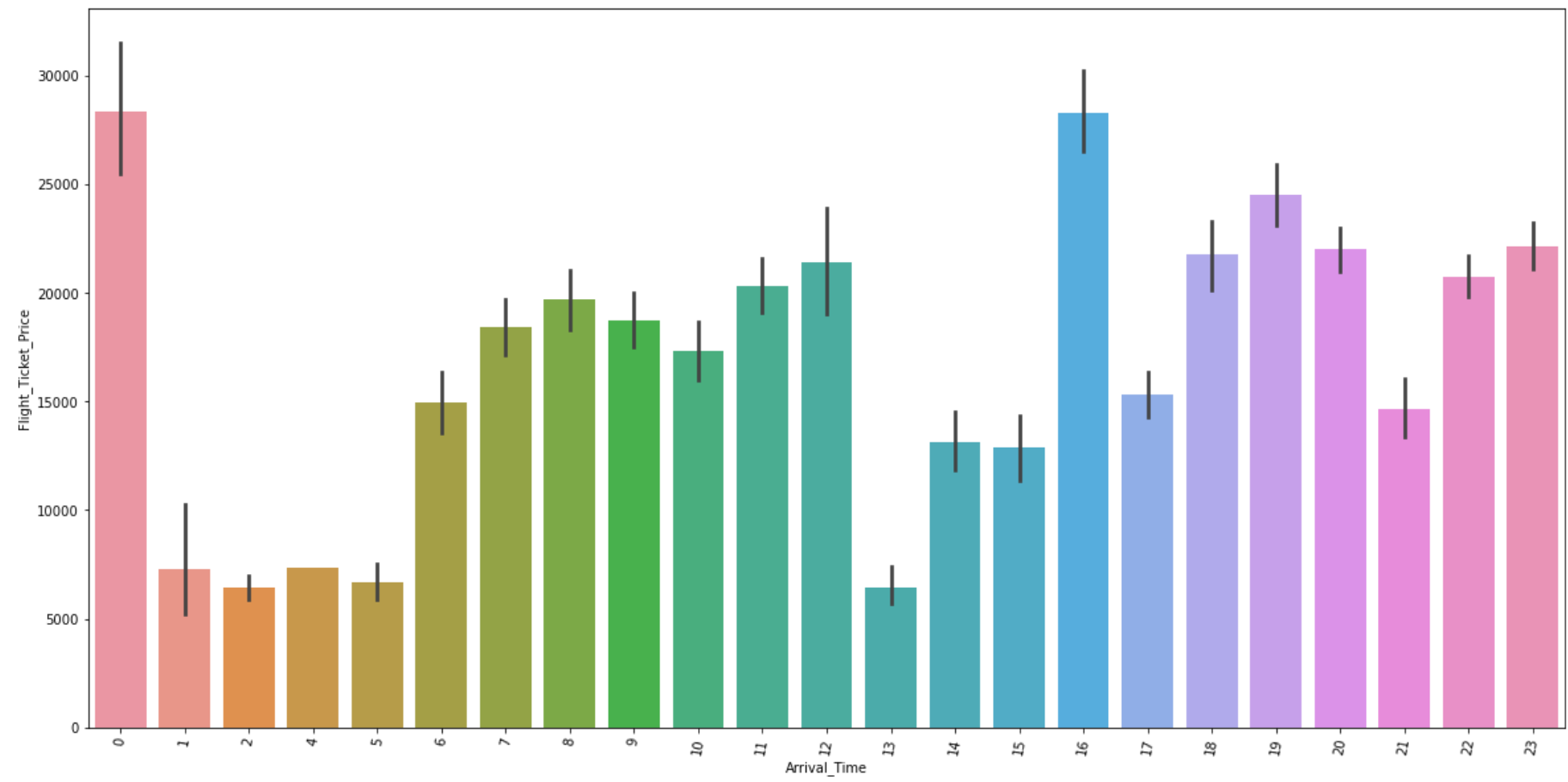
Departure Time:



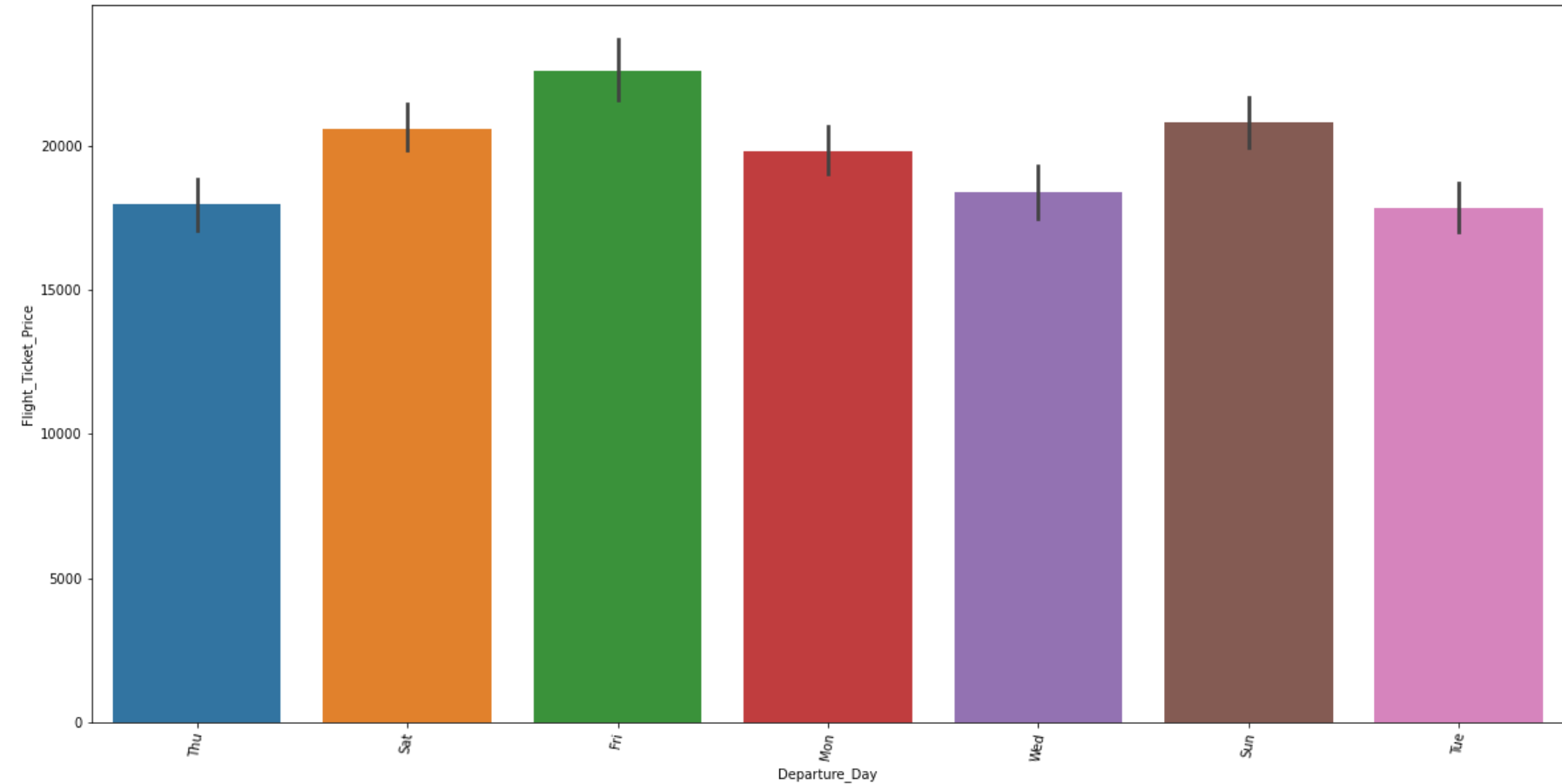
Departure Month:



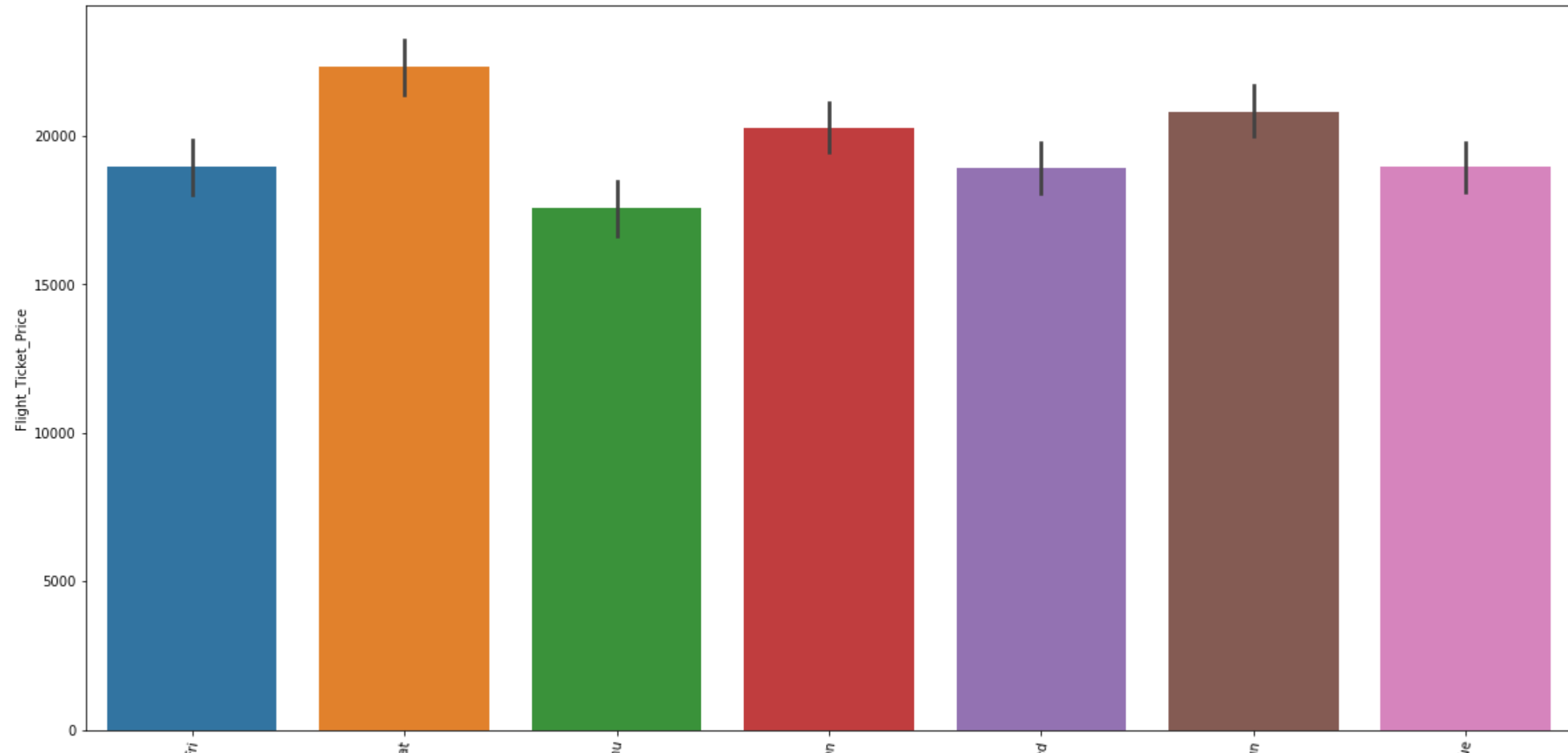
Arrival Time:



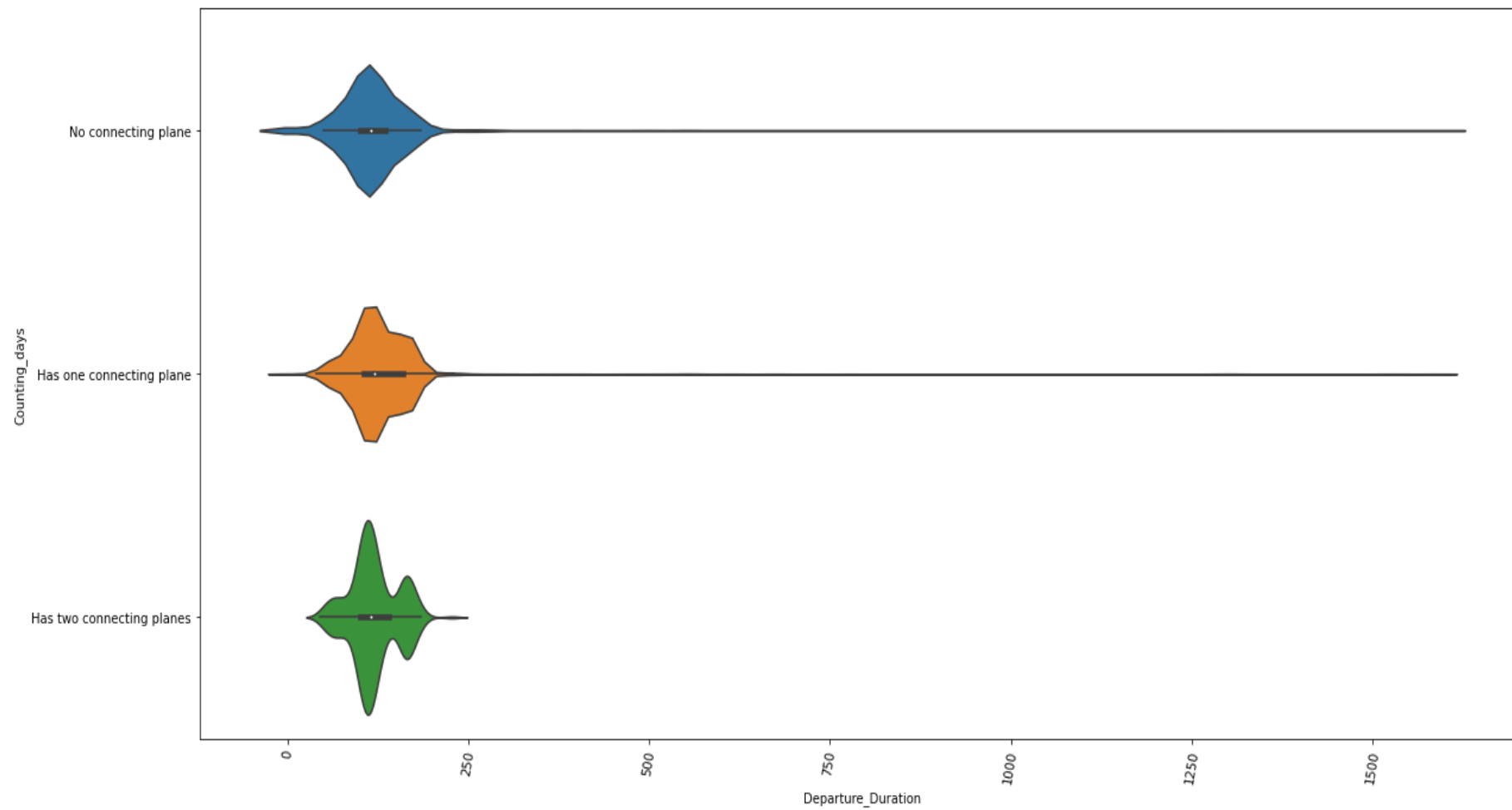
Departure Day:



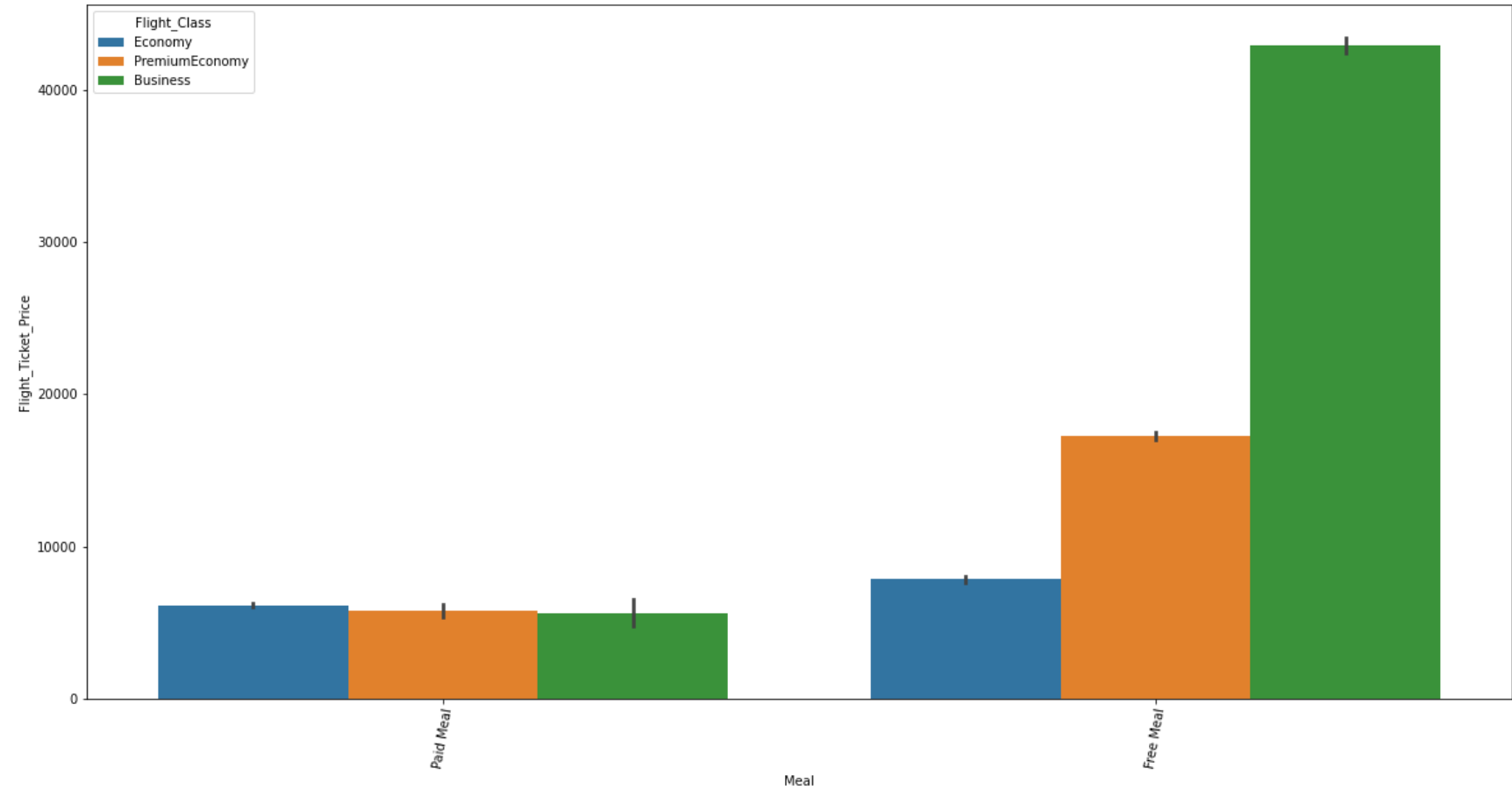
Arrival Day:



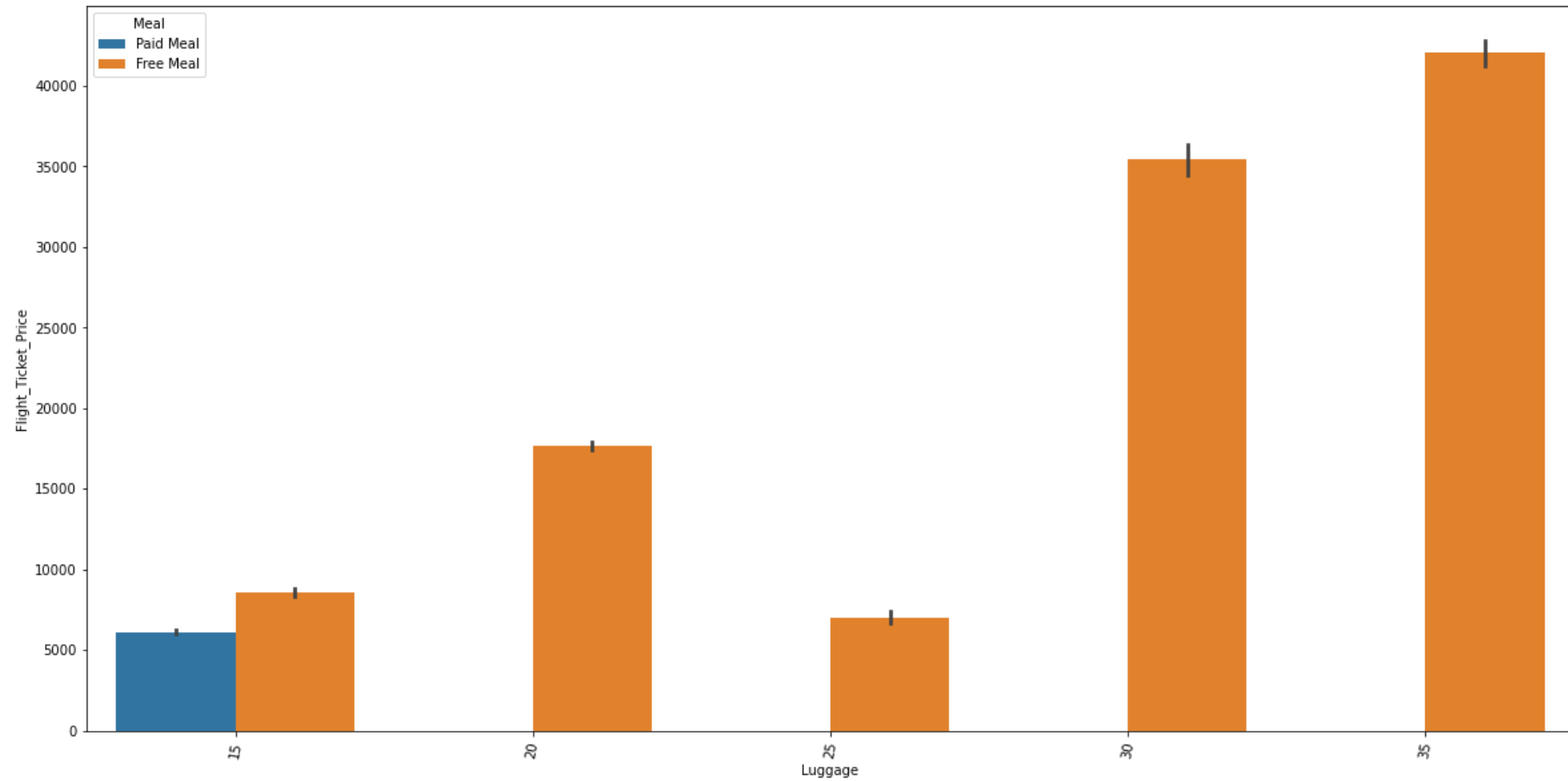
Duration:



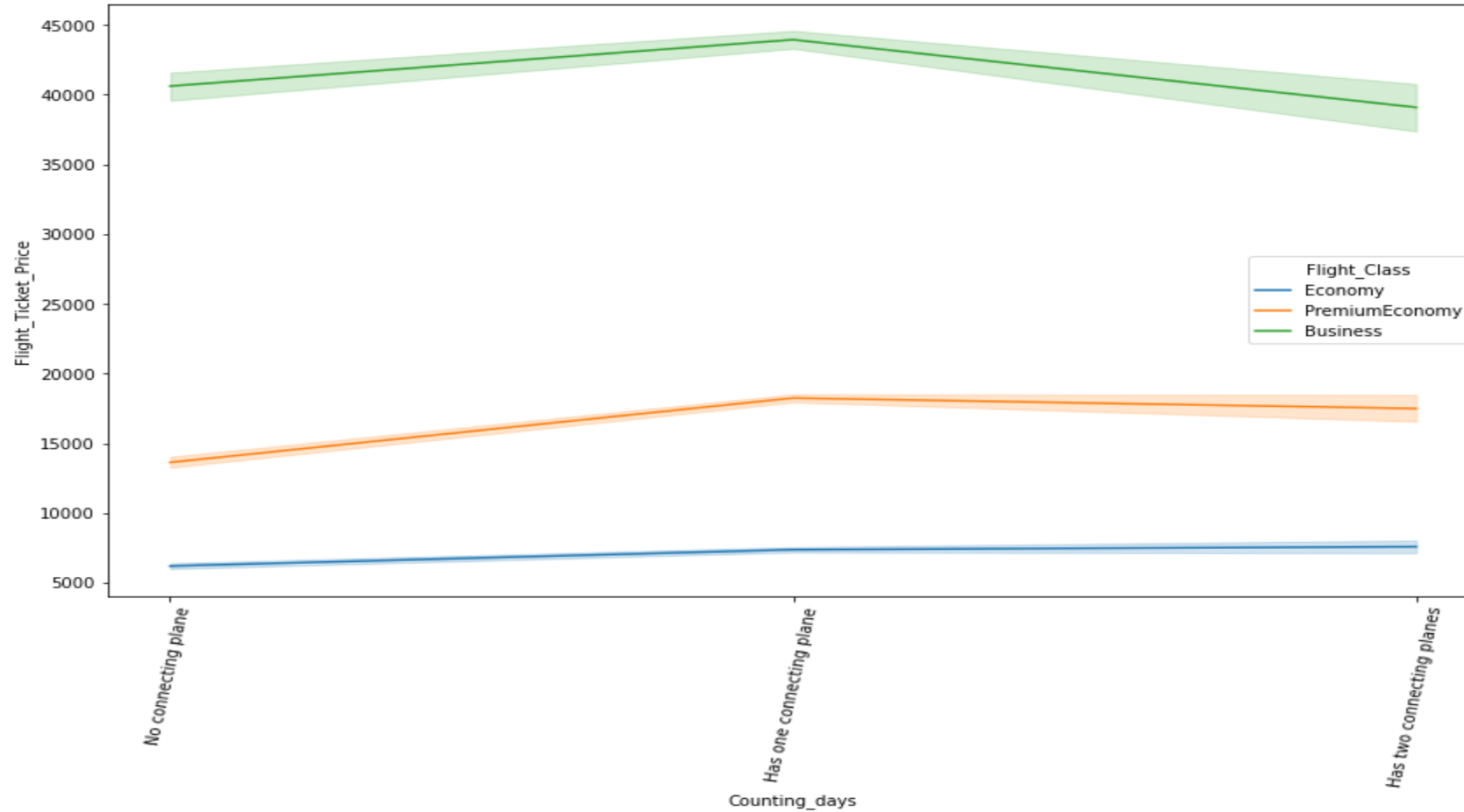
Meal:



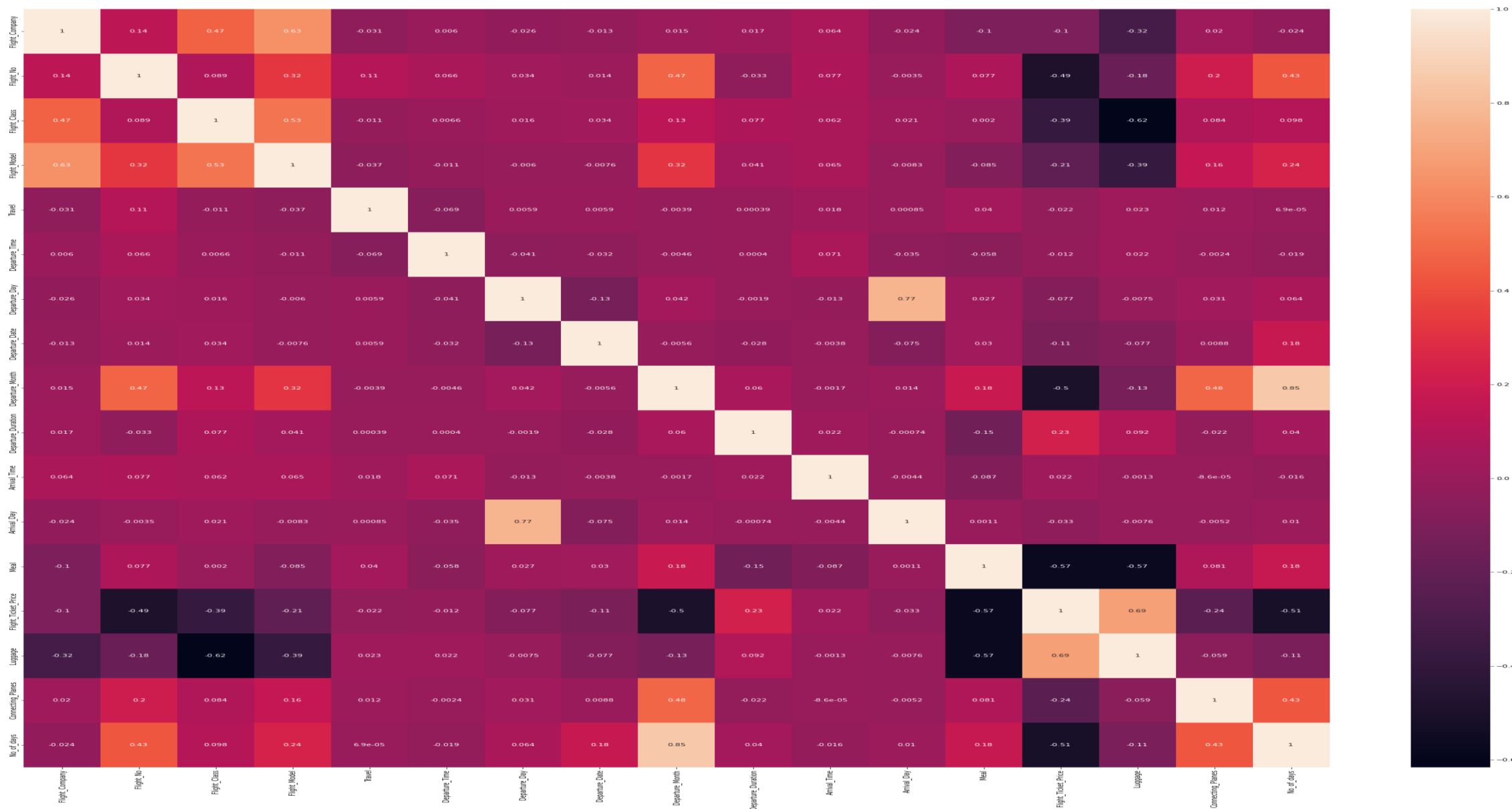
Luggage:



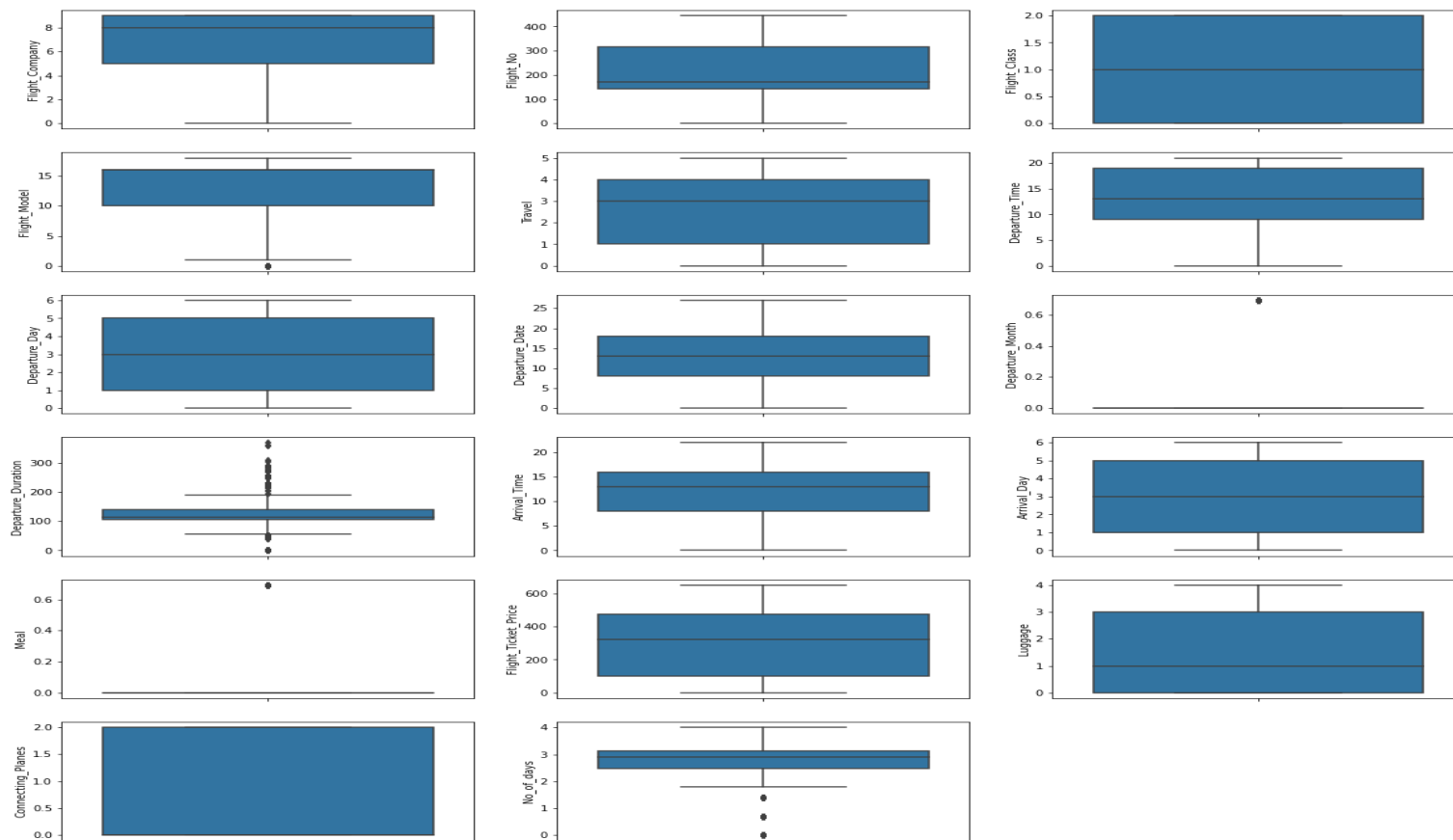
Connecting Flights:



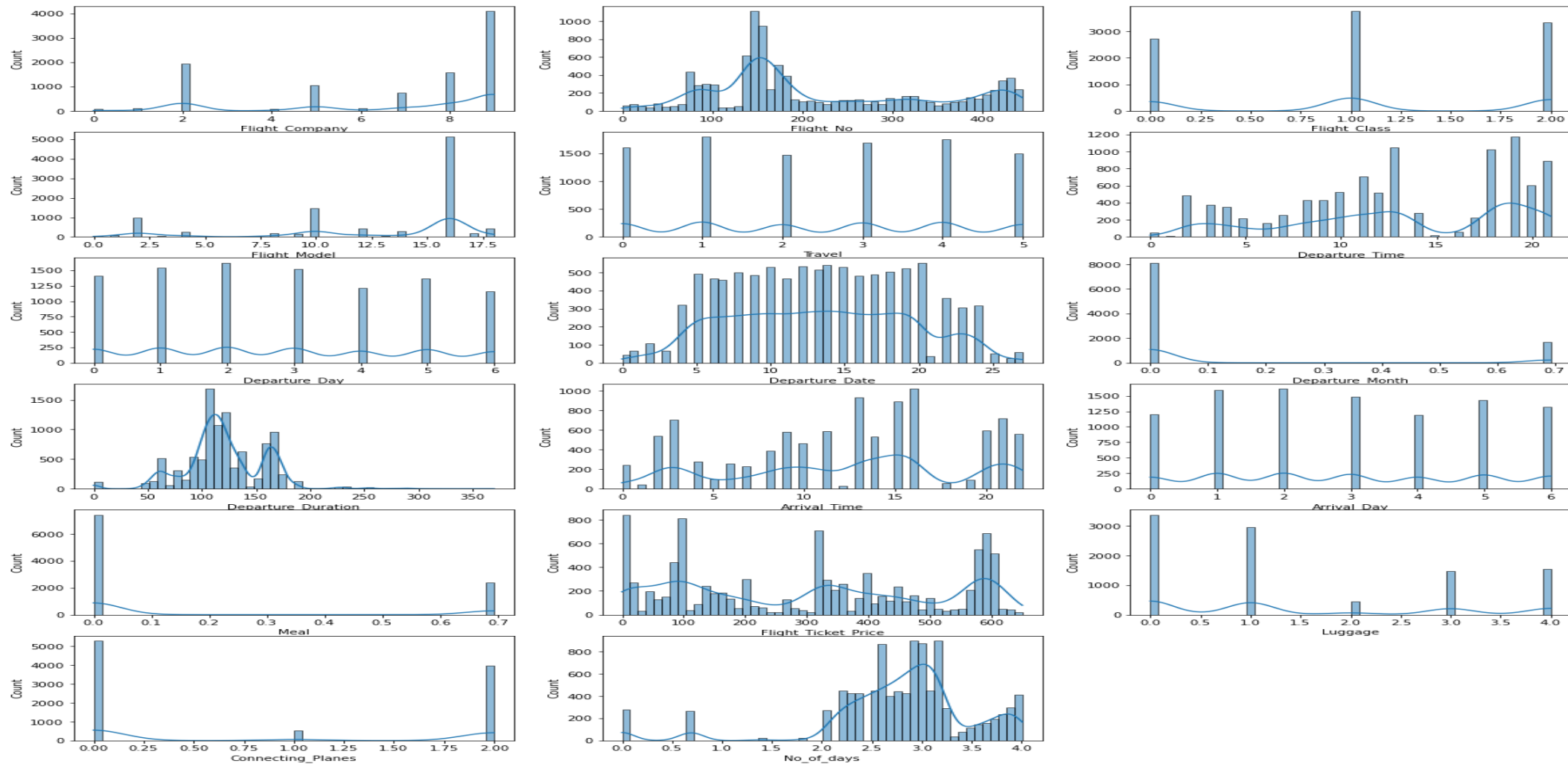
Heatmap:



Boxplot:



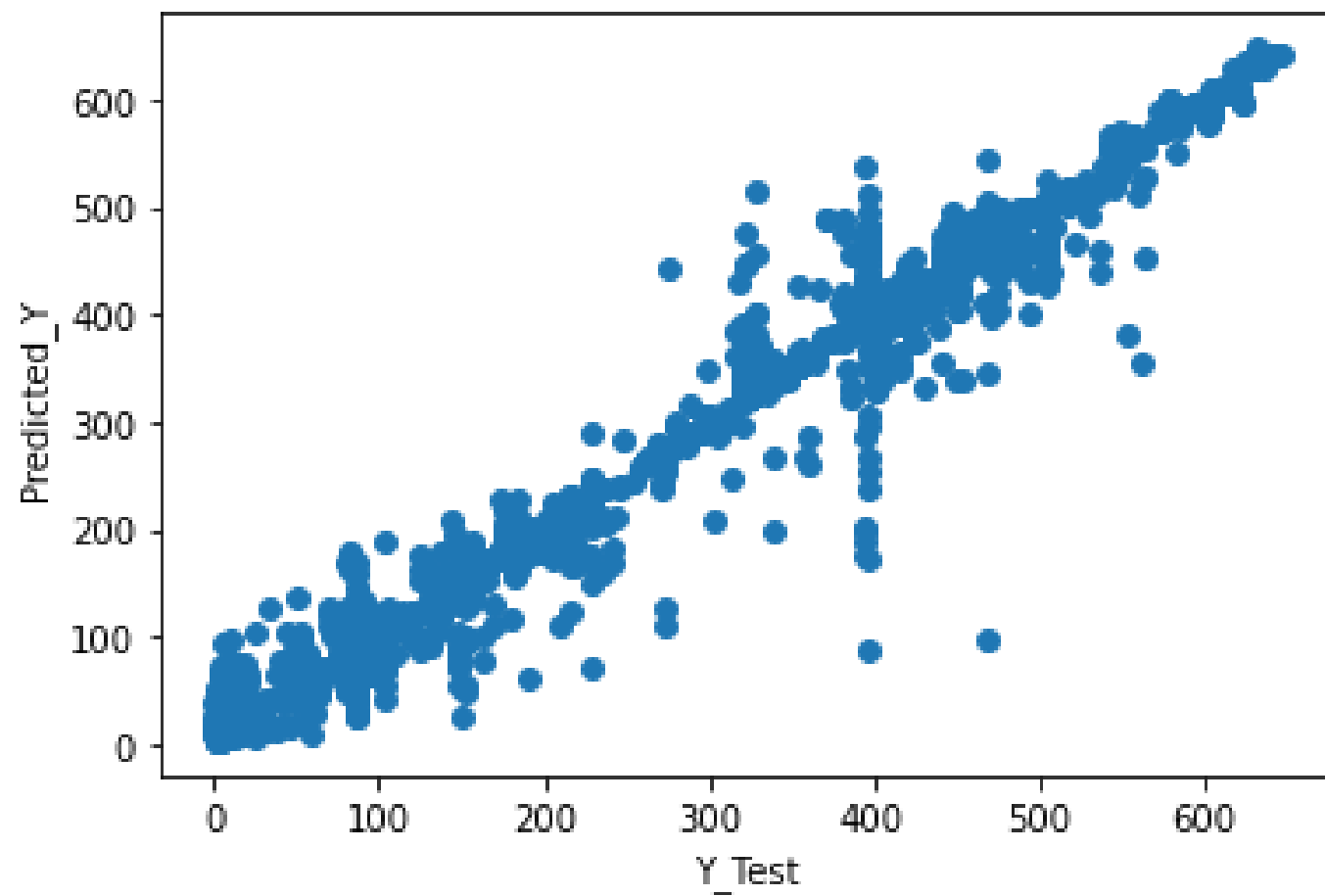
Histogram:



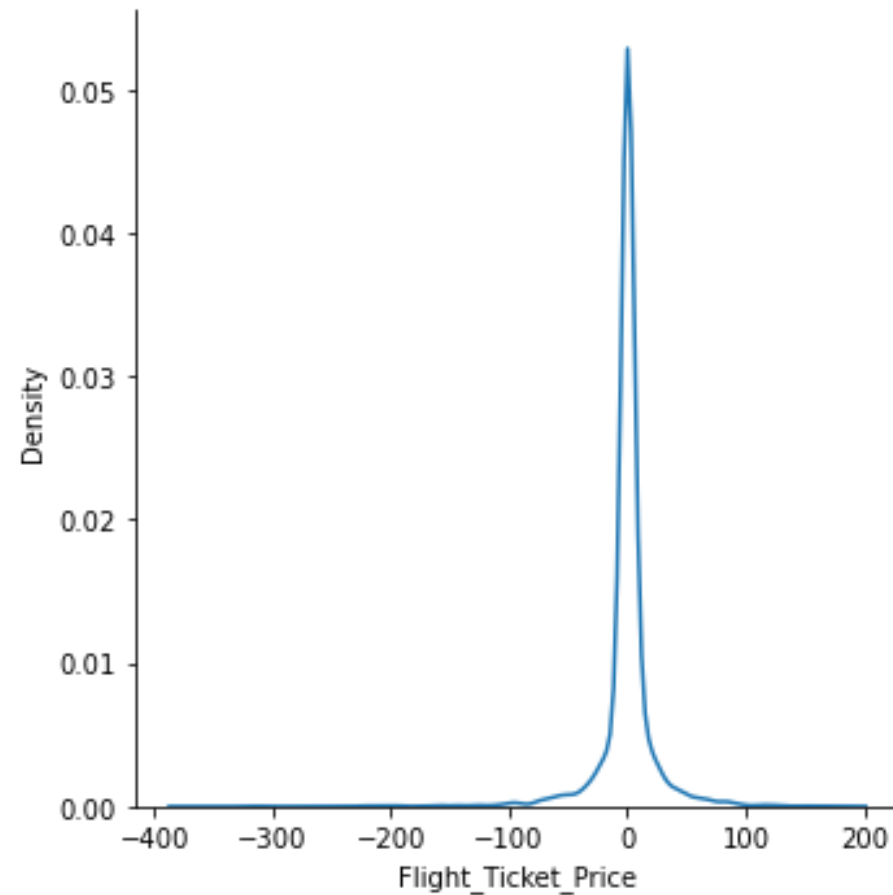
Testing of Identified Approaches (Algorithms):

- Linear Regression
- Gradient Boosting Regressor
- AdaBoost Regressor
- Decision Tree Regressor
- KNeighbors Regressor
- Extra Trees Regressor
- Random Forest Regressor

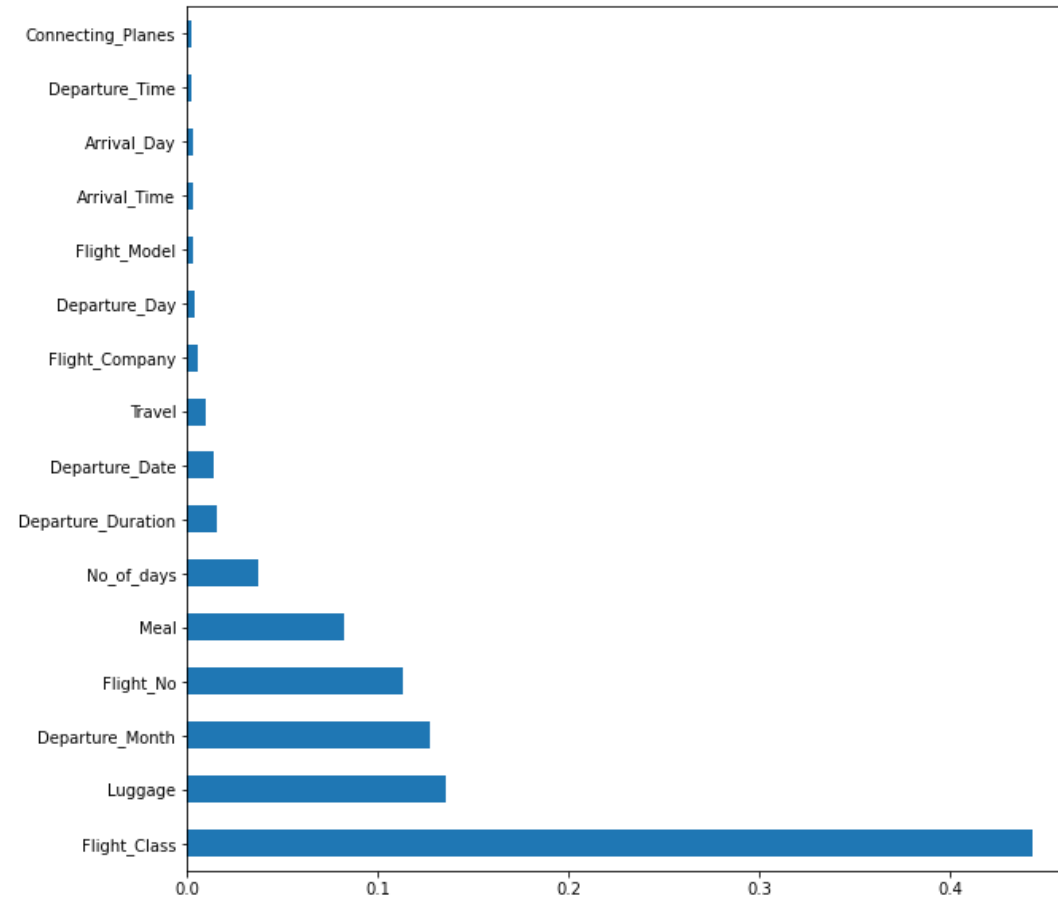
Y test plot:



Distribution of Predicted Values:



Feature importance's



Accuracy Parameter:

- R2 Score: 98.48414998011373
- Mean Absolute Error: 10.476769408502774
- Extra Forest Regression gives an R2 score 0.98
- Flight class dominates the flight price more.

Key Findings and Conclusions of the Study:

- This dataset has been taken from 2 websites of this, Yatra and MakeMyTrip constitutes the majority of data
- Since the target feature is continuous data, this problem can be solved by regression algorithms
- Extra Forest Regression gives an R^2 score 0.98
- Flight class dominates the flight price more.