

## Flight Price Prediction Project

We have 2 datasets here — training set and test set. The training set contains the features, along with the prices of the flights. It contains 10683 records, 10 input features and 1 output column — 'Price'. The test set contains 2671 records and 10 input features. The output 'Price' column needs to be predicted in this set. We will use Regression techniques here, since the predicted output will be a continuous value.

## Importing Libraries

```
In [1]: import matplotlib.pyplot as plt
import seaborn as sns
```

## Loading DataSet

```
In [5]: import pandas as pd
df=pd.read_csv("Data_Train.csv")
df.head()
```

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional_Info	Price
0	IndGo	24/03/2019	Banglore	New Delhi	BLR ? DEL	22:20	01:10 22 Mar	2h 50m	non-stop	No info	3897
1	Air India	1/05/2019	Kolkata	Banglore	CCU ? IXR ? BBI ? BLR	05:50	13:15	7h 25m	2 stops	No info	7662
2	Jet Airways	9/06/2019	Delhi	Cochin	DEL ? LKO ? BOM ? COK	09:25	04:25 10 Jun	19h	2 stops	No info	13882
3	IndGo	12/05/2019	Kolkata	Banglore	CCU ? NAG ? BLR	18:05	23:30	5h 25m	1 stop	No info	6218
4	IndGo	01/03/2019	Banglore	New Delhi	BLR ? NAG ? DEL	16:50	21:35	4h 45m	1 stop	No info	13302

```
In [6]: df
```

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional_Info	Price
0	IndGo	24/03/2019	Banglore	New Delhi	BLR ? DEL	22:20	01:10 22 Mar	2h 50m	non-stop	No info	3897
1	Air India	1/05/2019	Kolkata	Banglore	CCU ? IXR ? BBI ? BLR	05:50	13:15	7h 25m	2 stops	No info	7662
2	Jet Airways	9/06/2019	Delhi	Cochin	DEL ? LKO ? BOM ? COK	09:25	04:25 10 Jun	19h	2 stops	No info	13882
3	IndGo	12/05/2019	Kolkata	Banglore	CCU ? NAG ? BLR	18:05	23:30	5h 25m	1 stop	No info	6218
4	IndGo	01/03/2019	Banglore	New Delhi	BLR ? NAG ? DEL	16:50	21:35	4h 45m	1 stop	No info	13302
...	...	...	...	...	...	...	...	...	...	...	...
10678	Air Asia	9/04/2019	Kolkata	Banglore	CCU ? BLR	19:55	22:25	2h 30m	non-stop	No info	4107
10679	Air India	27/04/2019	Kolkata	Banglore	CCU ? BLR	20:45	23:20	2h 35m	non-stop	No info	4145
10680	Jet Airways	27/04/2019	Banglore	Delhi	BLR ? DEL	08:20	11:20	3h	non-stop	No info	7229
10681	Vistara	01/02/2019	Banglore	New Delhi	BLR ? DEL	11:30	14:10	2h 40m	non-stop	No info	12648
10682	Air India	9/05/2019	Delhi	Cochin	DEL ? GOI ? BOM ? COK	10:55	19:15	8h 20m	2 stops	No info	11753

10683 rows × 11 columns

```
In [7]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10683 entries, 0 to 10682
Data columns (total 11 columns):
# Column            Non-Null Count  Dtype
---  ---
0 Airline            10683 non-null object
1 Date_of_Journey    10683 non-null object
2 Source             10683 non-null object
3 Destination         10683 non-null object
4 Route              10682 non-null object
5 Dep_Time           10683 non-null object
6 Arrival_Time       10683 non-null object
7 Duration           10683 non-null object
8 Total_Stops        10682 non-null object
9 Additional_Info     10682 non-null object
10 Price             10683 non-null int64
dtypes: int64(1), object(10)
memory usage: 518.2+ KB
```

```
In [8]: df["Duration"].value_counts()
```

```
Out[8]: 2h 50m    559
1h 30m      396
2h 55m      337
2h 45m      337
2h 35m      329
...
33h 45m      1
47h 40m      1
47h          1
41h 20m      1
42h 45m      1
Name: Duration, Length: 368, dtype: int64
```

```
In [9]: df.dropna(inplace = True)
```

## Checking Null Value

```
In [10]: df.isnull().sum()
```

```
Out[10]: Airline      0
Date_of_Journey    0
Source              0
Destination         0
Route              0
Dep_Time           0
Arrival_Time       0
Duration           0
Total_Stops        0
Additional_Info     0
Price              0
dtype: int64
```

```
In [11]: df["Journey_day"] = pd.to_datetime(df.Date_of_Journey, format="%d/%m/%Y").dt.day
```

```
In [ ]: df["Journey_month"] = pd.to_datetime(df["Date_of_Journey"], format = "%d/%m/%Y").dt.month
```

```
In [12]: df.head()
```

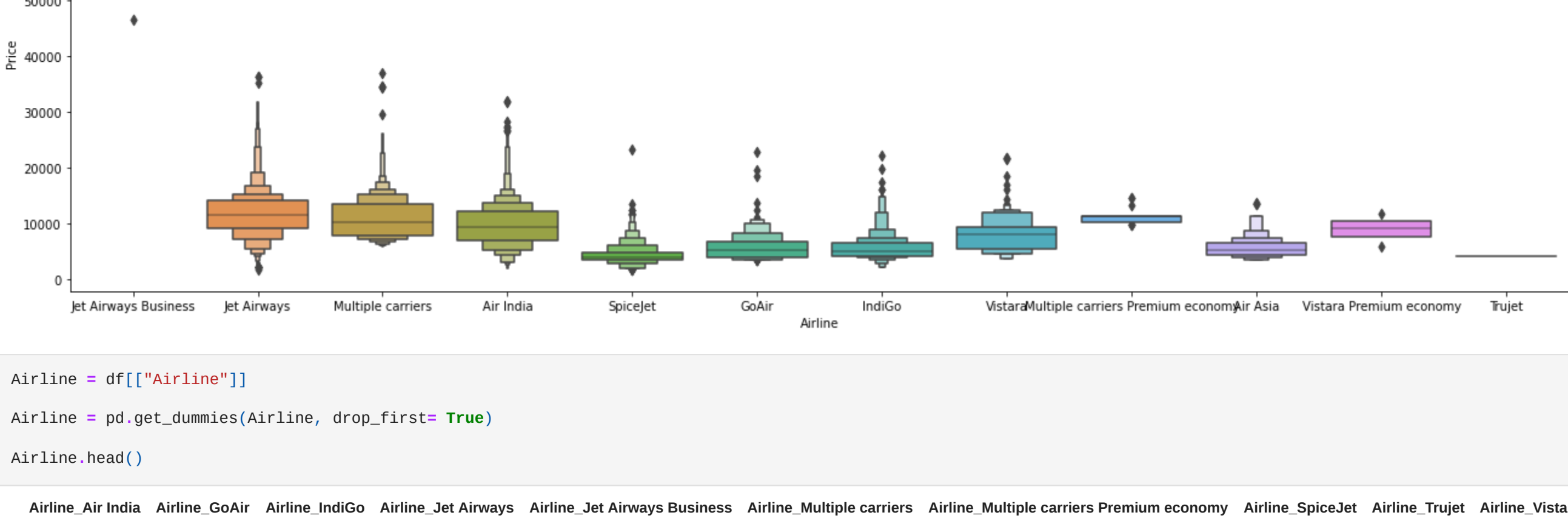
	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional_Info	Price	Journey_day
0	IndGo	24/03/2019	Banglore	New Delhi	BLR ? DEL	22:20	01:10 22 Mar	2h 50m	non-stop	No info	3897	24
1	Air India	1/05/2019	Kolkata	Banglore	CCU ? IXR ? BBI ? BLR	05:50	13:15	7h 25m	2 stops	No info	7662	1
2	Jet Airways	9/06/2019	Delhi	Cochin	DEL ? LKO ? BOM ? COK	09:25	04:25 10 Jun	19h	2 stops	No info	13882	9
3	IndGo	12/05/2019	Kolkata	Banglore	CCU ? NAG ? BLR	18:05	23:30	5h 25m	1 stop	No info	6218	12
4	IndGo	01/03/2019	Banglore	New Delhi	BLR ? NAG ? DEL	16:50	21:35	4h 45m	1 stop	No info	13302	1

```
In [ ]: df.columns
```

```
In [ ]: df.drop(["Date_of_Journey"], axis = 1, inplace = True)
```

## Data Visualization

```
In [13]: sns.catplot(y = "Price", x = "Airline", data = df.sort_values("Price", ascending = False), kind="boxen", height = 6, aspect = 3)
plt.show()
```



```
In [14]: Airline = df[["Airline"]]
Airline = pd.get_dummies(Airline, drop_first= True)
Airline.head()
```

	Airline_Air India	Airline_GoAir	Airline_IndGo	Airline_Jet Airways	Airline_Jet Airways Business	Airline_Multiple carriers	Airline_Multiple carriers Premium economy	Airline_SpiceJet	Airline_Trujet	Airline_Vistara	Airline_Vistara Premium economy
0	0	0	1	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0	0	0	0	0
3	0	0	1	0	0	0	0	0	0	0	0
4	0	0	1	0	0	0	0	0	0	0	0

```
In [15]: df["Source"].value_counts()
```

```
Out[15]: Delhi      4536
Kolkata      2871
Banglore     2197
Mumbai       697
Chennai      381
Name: Source, dtype: int64
```

```
In [16]: sns.catplot(y = "Price", x = "Source", data = df.sort_values("Price", ascending = False), kind="boxen", height = 4, aspect = 3)
plt.show()
```



```
In [17]: Source = df[["Source"]]
Source = pd.get_dummies(Source, drop_first= True)
Source.head()
```

	Source_Chennai	Source_Delhi	Source_Kolkata	Source_Mumbai
0	0	0	0	0
1	0	0	1	0
2	0	1	0	0
3	0	0	1	0
4	0	0	0	0

```
In [18]: df["Destination"].value_counts()
```

```
Out[18]: Cochin      4536
Banglore     2871
Delhi        1285
New Delhi    832
Hyderabad    697
Kolkata      381
Name: Destination, dtype: int64
```

```
In [19]: Destination = df[["Destination"]]
Destination = pd.get_dummies(Destination, drop_first = True)
Destination.head()
```

	Destination_Cochin	Destination_Delhi	Destination_Hyderabad	Destination_Kolkata	Destination_New Delhi
0	0	0	0	0	1
1	0	0	0	0	0
2	1	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	1

```
In [20]: df["Route"]
```

```
Out[20]: 0          BLR ? DEL
1      CCU ? IXR ? BBI ? BLR
2      DEL ? LKO ? BOM ? COK
3      CCU ? NAG ? BLR
4          BLR ? NAG ? DEL
...
10678      CCU ? BLR
10679      CCU ? BLR
10680      BLR ? DEL
10681      DEL ? GOI ? BOM ? COK
10682      DEL ? GOI ? BOM ? COK
Name: Route, Length: 10682, dtype: object
```

```
In [21]: df["Total_Stops"].value_counts()
```

```
Out[21]: 1 stop      5625
non-stop    3491
2 stops     1520
3 stops      45
4 stops      1
Name: Total_Stops, dtype: int64
```

```
In [ ]: df.replace(["non-stop": 0, "1 stop": 1, "2 stops": 2, "3 stops": 3, "4 stops": 4], inplace = True)
```

```
In [22]: df.head()
```

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional_Info	Price	Journey_day
0	IndGo	24/03/2019	Banglore	New Delhi	BLR ? DEL	22:20	01:10 22 Mar	2h 50m	non-stop	No info	3897	24
1	Air India	1/05/2019	Kolkata	Banglore	CCU ? IXR ? BBI ? BLR	05:50	13:15	7h 25m	2 stops	No info	7662	1
2	Jet Airways	9/06/2019	Delhi	Cochin	DEL ? LKO ? BOM ? COK	09:25	04:25 10 Jun	19h	2 stops	No info	13882	9
3	IndGo	12/05/2019	Kolkata	Banglore	CCU ? NAG ? BLR	18:05	23:30	5h 25m	1 stop	No info	6218	12
4	IndGo	01/03/2019	Banglore	New Delhi	BLR ? NAG ? DEL	16:50	21:35	4h 45m	1 stop	No info	13302	1

```
In [23]: df = pd.concat([df, Airline, Source, Destination], axis = 1)
df.head()
```

2	Jet Airways	9/06/2019	Delhi	Cochin	DEL ? LKO ? BOM ? COK	09:25	04:25 10 Jun	19h	2 stops	No info	...	0	0	1	0	0	1	0	0	0
3	IndiGo	12/05/2019	Kolkata	Banglore	CCU NAG ? BLR	18:05	23:30	5h 25m	1 stop	No info	...	0	0	0	1	0	0	0	0	0
4	IndiGo	01/03/2019	Banglore	New Delhi	BLR ? NAG ? DEL	16:50	21:35	4h 45m	1 stop	No info	...	0	0	0	0	0	0	0	0	0

5 rows × 32 columns

```
In [24]: df.shape
```

```
Out[24]: (10682, 32)
```

```
In [25]: df = pd.read_csv("Test_set.csv")
df.head(10)
```

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional_Info
0	Jet Airways	6/06/2019	Delhi	Cochin	DEL ? BOM ? COK	17:30	04:25 07 Jun	10h 55m	1 stop	No info
1	IndGo	12/05/2019	Kolkata	Banglore	CCU ? MAA ? BLR	06:20	10:20	4h	1 stop	No info
2	Jet Airways	21/05/2019	Delhi	Cochin	DEL ? BOM ? COK	19:15	19:00 22 May	23h 45m	1 stop	In-flight meal not included
3	Multiple carriers	21/05/2019	Delhi	Cochin	DEL ? BOM ? COK	08:00	21:00	13h	1 stop	No info
4	Air Asia	24/06/2019	Banglore	Delhi	BLR ? DEL	23:55	02:45 25 Jun	2h 50m	non-stop	No info
5	Jet Airways	12/06/2019	Delhi	Cochin	DEL ? BOM ? COK	18:15	12:35 13 Jun	18h 20m	1 stop	In-flight meal not included
6	Air India	12/03/2019	Banglore	New Delhi	BLR ? TRV ? DEL	07:30	22:35	15h 5m	1 stop	No info
7	IndGo	12/03/2019	Kolkata	Banglore	CCU ? HYD ? BLR	15:15	20:30	5h 15m	1 stop	No info
8	IndGo	15/05/2019	Kolkata	Banglore	CCU ? BLR	10:10	12:55	2h 45m	non-stop	No info
9	Jet Airways	18/05/2019	Kolkata	Banglore	CCU ? BOM ? BLR	16:30	22:35	6h 5m	1 stop	No info

```
In [26]: df.describe()
```

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional_Info
count	2671	2671	2671	2671	2671	2671	2671	2671	2671	2671
unique	11	44	5	6	100	199	704	320	5	6
top	Jet Airways	9/05/2019	Delhi	Cochin	DEL ? BOM ? COK	10:00	19:00	2h 50m	1 stop	No info
freq	897	144	1145	1145	624	62	113	122	1431	2148

```
In [33]: df.plot(kind='box',subplots=True,layout=(2,6),figsize=(10,10))
```

```
Out[33]: —
```

## Thank you