# Mini Project on flightpriceprediction

# Problem Statement:which model is suitable for the dataset

```
In [1]:  #importing libraries
         import numpy as np
         import pandas as pd
         from sklearn import preprocessing
         import matplotlib.pyplot as plt
         import seaborn as sns
         sns.set(style="white")
         #seaborn plots
         sns.set(style="whitegrid",color_codes=True)
         import warnings
         warnings.simplefilter (action='ignore')
```

In [3]:
```python
train_df=pd.read_csv(r"C:\Users\user\Downloads\Data_Train 1.csv")
train_df
```

Out[3]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Tota |
|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | 24/03/2019 | Banglore | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m | |
| 1 | Air India | 1/05/2019 | Kolkata | Banglore | CCU ? IXR ? BBI ? BLR | 05:50 | 13:15 | 7h 25m | |
| 2 | Jet Airways | 9/06/2019 | Delhi | Cochin | DEL ? LKO ? BOM ? COK | 09:25 | 04:25 10 Jun | 19h | |
| 3 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU ? NAG ? BLR | 18:05 | 23:30 | 5h 25m | |
| 4 | IndiGo | 01/03/2019 | Banglore | New Delhi | BLR ? NAG ? DEL | 16:50 | 21:35 | 4h 45m | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 10678 | Air Asia | 9/04/2019 | Kolkata | Banglore | CCU ? BLR | 19:55 | 22:25 | 2h 30m | |
| 10679 | Air India | 27/04/2019 | Kolkata | Banglore | CCU ? BLR | 20:45 | 23:20 | 2h 35m | |
| 10680 | Jet Airways | 27/04/2019 | Banglore | Delhi | BLR ? DEL | 08:20 | 11:20 | 3h | |
| 10681 | Vistara | 01/03/2019 | Banglore | New Delhi | BLR ? DEL | 11:30 | 14:10 | 2h 40m | |
| 10682 | Air India | 9/05/2019 | Delhi | Cochin | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m | |

10683 rows × 11 columns

In [4]:
```python
test_df=pd.read_csv(r"C:\Users\user\Downloads\Test_set.csv")
test_df
```

Out[4]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Jet Airways | 6/06/2019 | Delhi | Cochin | DEL ? BOM ? COK | 17:30 | 04:25 07 Jun | 10h 55m | |
| 1 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU ? MAA ? BLR | 06:20 | 10:20 | 4h | |
| 2 | Jet Airways | 21/05/2019 | Delhi | Cochin | DEL ? BOM ? COK | 19:15 | 19:00 22 May | 23h 45m | |
| 3 | Multiple carriers | 21/05/2019 | Delhi | Cochin | DEL ? BOM ? COK | 08:00 | 21:00 | 13h | |
| 4 | Air Asia | 24/06/2019 | Banglore | Delhi | BLR ? DEL | 23:55 | 02:45 25 Jun | 2h 50m | n |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 2666 | Air India | 6/06/2019 | Kolkata | Banglore | CCU ? DEL ? BLR | 20:30 | 20:25 07 Jun | 23h 55m | |
| 2667 | IndiGo | 27/03/2019 | Kolkata | Banglore | CCU ? BLR | 14:20 | 16:55 | 2h 35m | n |
| 2668 | Jet Airways | 6/03/2019 | Delhi | Cochin | DEL ? BOM ? COK | 21:50 | 04:25 07 Mar | 6h 35m | |
| 2669 | Air India | 6/03/2019 | Delhi | Cochin | DEL ? BOM ? COK | 04:00 | 19:15 | 15h 15m | |
| 2670 | Multiple carriers | 15/06/2019 | Delhi | Cochin | DEL ? BOM ? COK | 04:55 | 19:15 | 14h 20m | |

2671 rows × 10 columns

In [5]: `train_df.head()`

Out[5]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_St |
|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | 24/03/2019 | Banglore | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m | non-s |
| 1 | Air India | 1/05/2019 | Kolkata | Banglore | CCU ? IXR ? BBI ? BLR | 05:50 | 13:15 | 7h 25m | 2 st |
| 2 | Jet Airways | 9/06/2019 | Delhi | Cochin | DEL ? LKO ? BOM ? COK | 09:25 | 04:25 10 Jun | 19h | 2 st |
| 3 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU ? NAG ? BLR | 18:05 | 23:30 | 5h 25m | 1 s |
| 4 | IndiGo | 01/03/2019 | Banglore | New Delhi | BLR ? NAG ? DEL | 16:50 | 21:35 | 4h 45m | 1 s |

In [6]: `test_df.head()`

Out[6]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_St |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Jet Airways | 6/06/2019 | Delhi | Cochin | DEL ? BOM ? COK | 17:30 | 04:25 07 Jun | 10h 55m | 1 s |
| 1 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU ? MAA ? BLR | 06:20 | 10:20 | 4h | 1 s |
| 2 | Jet Airways | 21/05/2019 | Delhi | Cochin | DEL ? BOM ? COK | 19:15 | 19:00 22 May | 23h 45m | 1 s |
| 3 | Multiple carriers | 21/05/2019 | Delhi | Cochin | DEL ? BOM ? COK | 08:00 | 21:00 | 13h | 1 s |
| 4 | Air Asia | 24/06/2019 | Banglore | Delhi | BLR ? DEL | 23:55 | 02:45 25 Jun | 2h 50m | non-s |

In [7]: `train_df.tail()`

Out[7]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Tota |
|---|---|---|---|---|---|---|---|---|---|
| 10678 | Air Asia | 9/04/2019 | Kolkata | Banglore | CCU ? BLR | 19:55 | 22:25 | 2h 30m | |
| 10679 | Air India | 27/04/2019 | Kolkata | Banglore | CCU ? BLR | 20:45 | 23:20 | 2h 35m | |
| 10680 | Jet Airways | 27/04/2019 | Banglore | Delhi | BLR ? DEL | 08:20 | 11:20 | 3h | |
| 10681 | Vistara | 01/03/2019 | Banglore | New Delhi | BLR ? DEL | 11:30 | 14:10 | 2h 40m | |
| 10682 | Air India | 9/05/2019 | Delhi | Cochin | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m | |

In [8]: `test_df.tail()`

Out[8]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_ |
|---|---|---|---|---|---|---|---|---|---|
| 2666 | Air India | 6/06/2019 | Kolkata | Banglore | CCU ? DEL ? BLR | 20:30 | 20:25 07 Jun | 23h 55m | |
| 2667 | IndiGo | 27/03/2019 | Kolkata | Banglore | CCU ? BLR | 14:20 | 16:55 | 2h 35m | no |
| 2668 | Jet Airways | 6/03/2019 | Delhi | Cochin | DEL ? BOM ? COK | 21:50 | 04:25 07 Mar | 6h 35m | |
| 2669 | Air India | 6/03/2019 | Delhi | Cochin | DEL ? BOM ? COK | 04:00 | 19:15 | 15h 15m | |
| 2670 | Multiple carriers | 15/06/2019 | Delhi | Cochin | DEL ? BOM ? COK | 04:55 | 19:15 | 14h 20m | |

In [9]: `train_df.shape`

Out[9]: `(10683, 11)`

In [10]: `test_df.shape`

Out[10]: `(2671, 10)`

In [11]: `train_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  10683 non-null  object
 10  Price            10683 non-null  int64
dtypes: int64(1), object(10)
memory usage: 918.2+ KB
```

```
In [12]: test_df.info
```

```
Out[12]: <bound method DataFrame.info of                   Airline Date_of_Journey    Source
         Destination
         0              Jet Airways       6/06/2019     Delhi       Cochin  \
         1                   IndiGo      12/05/2019   Kolkata     Banglore
         2              Jet Airways      21/05/2019     Delhi       Cochin
         3         Multiple carriers     21/05/2019     Delhi       Cochin
         4                 Air Asia      24/06/2019  Banglore        Delhi
         ...                    ...             ...       ...          ...
         2666             Air India       6/06/2019   Kolkata     Banglore
         2667               IndiGo      27/03/2019   Kolkata     Banglore
         2668           Jet Airways       6/03/2019     Delhi       Cochin
         2669             Air India       6/03/2019     Delhi       Cochin
         2670     Multiple carriers     15/06/2019     Delhi       Cochin

                          Route Dep_Time  Arrival_Time Duration Total_Stops
         0        DEL ? BOM ? COK    17:30  04:25 07 Jun  10h 55m     1 stop  \
         1        CCU ? MAA ? BLR    06:20         10:20       4h     1 stop
         2        DEL ? BOM ? COK    19:15  19:00 22 May  23h 45m     1 stop
         3        DEL ? BOM ? COK    08:00         21:00      13h     1 stop
         4              BLR ? DEL    23:55  02:45 25 Jun   2h 50m    non-stop
         ...                ...      ...           ...      ...         ...
         2666     CCU ? DEL ? BLR    20:30  20:25 07 Jun  23h 55m     1 stop
         2667           CCU ? BLR    14:20         16:55   2h 35m    non-stop
         2668     DEL ? BOM ? COK    21:50  04:25 07 Mar   6h 35m     1 stop
         2669     DEL ? BOM ? COK    04:00         19:15  15h 15m     1 stop
         2670     DEL ? BOM ? COK    04:55         19:15  14h 20m     1 stop

                          Additional_Info
         0                        No info
         1                        No info
         2        In-flight meal not included
         3                        No info
         4                        No info
         ...                          ...
         2666                     No info
         2667                     No info
         2668                     No info
         2669                     No info
         2670                     No info

         [2671 rows x 10 columns]>
```

In [13]: `train_df.describe()`

Out[13]:

|        | Price          |
|--------|----------------|
| count  | 10683.000000   |
| mean   | 9087.064121    |
| std    | 4611.359167    |
| min    | 1759.000000    |
| 25%    | 5277.000000    |
| 50%    | 8372.000000    |
| 75%    | 12373.000000   |
| max    | 79512.000000   |

In [14]: `test_df.describe()`

Out[14]:

|        | Airline      | Date_of_Journey | Source | Destination | Route              | Dep_Time | Arrival_Time | Duration | Tota |
|--------|--------------|-----------------|--------|-------------|--------------------|----------|--------------|----------|------|
| count  | 2671         | 2671            | 2671   | 2671        | 2671               | 2671     | 2671         | 2671     |      |
| unique | 11           | 44              | 5      | 6           | 100                | 199      | 704          | 320      |      |
| top    | Jet Airways  | 9/05/2019       | Delhi  | Cochin      | DEL ? BOM ? COK    | 10:00    | 19:00        | 2h 50m   |      |
| freq   | 897          | 144             | 1145   | 1145        | 624                | 62       | 113          | 122      |      |

In [15]: `train_df.isnull().sum()`

Out[15]:
```
Airline             0
Date_of_Journey     0
Source              0
Destination         0
Route               1
Dep_Time            0
Arrival_Time        0
Duration            0
Total_Stops         1
Additional_Info     0
Price               0
dtype: int64
```

In [17]: `test_df.isnull().sum()`

Out[17]: 
```
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
dtype: int64
```

In [18]: `train_df.dropna(inplace=True)`

In [19]: `train_df['Airline'].value_counts()`

Out[19]: 
```
Airline
Jet Airways                          3849
IndiGo                               2053
Air India                            1751
Multiple carriers                    1196
SpiceJet                              818
Vistara                              479
Air Asia                             319
GoAir                                194
Multiple carriers Premium economy     13
Jet Airways Business                   6
Vistara Premium economy                3
Trujet                                 1
Name: count, dtype: int64
```

In [20]: `train_df['Source'].value_counts()`

Out[20]: 
```
Source
Delhi      4536
Kolkata    2871
Banglore   2197
Mumbai      697
Chennai     381
Name: count, dtype: int64
```

```
In [21]: airline={"Airline":{"Jet Airways":0,"IndiGo":1,"Air India":2,"Multiple carriers":3,
         "SpiceJet":4,"Vistara":5,"Air Asia":6,"GoAir":7,
         "Multiple carriers Premium economy":8,
         "Jet Airways Business":9,"Vistara Premium economy":10,"Trujet":11}}
         train_df=train_df.replace(airline)
         train_df
```

Out[21]:

|  | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Tota |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 24/03/2019 | Banglore | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m | r |
| 1 | 2 | 1/05/2019 | Kolkata | Banglore | CCU ? IXR ? BBI ? BLR | 05:50 | 13:15 | 7h 25m |  |
| 2 | 0 | 9/06/2019 | Delhi | Cochin | DEL ? LKO ? BOM ? COK | 09:25 | 04:25 10 Jun | 19h |  |
| 3 | 1 | 12/05/2019 | Kolkata | Banglore | CCU ? NAG ? BLR | 18:05 | 23:30 | 5h 25m |  |
| 4 | 1 | 01/03/2019 | Banglore | New Delhi | BLR ? NAG ? DEL | 16:50 | 21:35 | 4h 45m |  |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |  |
| 10678 | 6 | 9/04/2019 | Kolkata | Banglore | CCU ? BLR | 19:55 | 22:25 | 2h 30m | r |
| 10679 | 2 | 27/04/2019 | Kolkata | Banglore | CCU ? BLR | 20:45 | 23:20 | 2h 35m | r |
| 10680 | 0 | 27/04/2019 | Banglore | Delhi | BLR ? DEL | 08:20 | 11:20 | 3h | r |
| 10681 | 5 | 01/03/2019 | Banglore | New Delhi | BLR ? DEL | 11:30 | 14:10 | 2h 40m | r |
| 10682 | 2 | 9/05/2019 | Delhi | Cochin | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m |  |

10682 rows × 11 columns

```
In [22]: city={"Source":{"Delhi":0,"Kolkata":1,"Banglore":2,
         "Mumbai":3,"Chennai":4}}
         train_df=train_df.replace(city)
         train_df
```

Out[22]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 24/03/2019 | 2 | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m | nc |
| 1 | 2 | 1/05/2019 | 1 | Banglore | CCU ? IXR ? BBI ? BLR | 05:50 | 13:15 | 7h 25m | 2 |
| 2 | 0 | 9/06/2019 | 0 | Cochin | DEL ? LKO ? BOM ? COK | 09:25 | 04:25 10 Jun | 19h | 2 |
| 3 | 1 | 12/05/2019 | 1 | Banglore | CCU ? NAG ? BLR | 18:05 | 23:30 | 5h 25m | |
| 4 | 1 | 01/03/2019 | 2 | New Delhi | BLR ? NAG ? DEL | 16:50 | 21:35 | 4h 45m | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 10678 | 6 | 9/04/2019 | 1 | Banglore | CCU ? BLR | 19:55 | 22:25 | 2h 30m | nc |
| 10679 | 2 | 27/04/2019 | 1 | Banglore | CCU ? BLR | 20:45 | 23:20 | 2h 35m | nc |
| 10680 | 0 | 27/04/2019 | 2 | Delhi | BLR ? DEL | 08:20 | 11:20 | 3h | nc |
| 10681 | 5 | 01/03/2019 | 2 | New Delhi | BLR ? DEL | 11:30 | 14:10 | 2h 40m | nc |
| 10682 | 2 | 9/05/2019 | 0 | Cochin | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m | 2 |

10682 rows × 11 columns

In [23]:
```python
dest={"Destination":{"Cochin":0,"Banglore":1,"Delhi":2,
"New Delhi":3,"Hyderabad":4,"Kolkata":5}}
train_df=train_df.replace(dest)
train_df
```

Out[23]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 24/03/2019 | 2 | 3 | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m | no |
| 1 | 2 | 1/05/2019 | 1 | 1 | CCU ? IXR ? BBI ? BLR | 05:50 | 13:15 | 7h 25m | 2 |
| 2 | 0 | 9/06/2019 | 0 | 0 | DEL ? LKO ? BOM ? COK | 09:25 | 04:25 10 Jun | 19h | 2 |
| 3 | 1 | 12/05/2019 | 1 | 1 | CCU ? NAG ? BLR | 18:05 | 23:30 | 5h 25m | |
| 4 | 1 | 01/03/2019 | 2 | 3 | BLR ? NAG ? DEL | 16:50 | 21:35 | 4h 45m | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 10678 | 6 | 9/04/2019 | 1 | 1 | CCU ? BLR | 19:55 | 22:25 | 2h 30m | no |
| 10679 | 2 | 27/04/2019 | 1 | 1 | CCU ? BLR | 20:45 | 23:20 | 2h 35m | no |
| 10680 | 0 | 27/04/2019 | 2 | 2 | BLR ? DEL | 08:20 | 11:20 | 3h | no |
| 10681 | 5 | 01/03/2019 | 2 | 3 | BLR ? DEL | 11:30 | 14:10 | 2h 40m | no |
| 10682 | 2 | 9/05/2019 | 0 | 0 | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m | 2 |

10682 rows × 11 columns

In [24]:
```python
stops={"Total_Stops":{"non-stop":0,"1 stop":1,"2 stops":2,
"3 stops":3,"4 stops":4}}
train_df=train_df.replace(stops)
train_df
```

Out[24]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_ |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 24/03/2019 | 2 | 3 | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m | |
| **1** | 2 | 1/05/2019 | 1 | 1 | CCU ? IXR ? BBI ? BLR | 05:50 | 13:15 | 7h 25m | |
| **2** | 0 | 9/06/2019 | 0 | 0 | DEL ? LKO ? BOM ? COK | 09:25 | 04:25 10 Jun | 19h | |
| **3** | 1 | 12/05/2019 | 1 | 1 | CCU ? NAG ? BLR | 18:05 | 23:30 | 5h 25m | |
| **4** | 1 | 01/03/2019 | 2 | 3 | BLR ? NAG ? DEL | 16:50 | 21:35 | 4h 45m | |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | |
| **10678** | 6 | 9/04/2019 | 1 | 1 | CCU ? BLR | 19:55 | 22:25 | 2h 30m | |
| **10679** | 2 | 27/04/2019 | 1 | 1 | CCU ? BLR | 20:45 | 23:20 | 2h 35m | |
| **10680** | 0 | 27/04/2019 | 2 | 2 | BLR ? DEL | 08:20 | 11:20 | 3h | |
| **10681** | 5 | 01/03/2019 | 2 | 3 | BLR ? DEL | 11:30 | 14:10 | 2h 40m | |
| **10682** | 2 | 9/05/2019 | 0 | 0 | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m | |

10682 rows × 11 columns

# Data Visualization:

```
In [25]: fdf=train_df[['Airline','Source','Destination','Total_Stops','Price']]
         sns.heatmap(fdf.corr(),annot=True)
```

Out[25]: <Axes: >



# Feature Scaling : To Split the data into train and test data

```
In [26]: x=fdf[['Airline','Source','Destination','Total_Stops']]
         y=fdf['Price']
```

```
In [27]: from sklearn.model_selection import train_test_split
         X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=100)
```

# Linear Regression

In [28]:
```python
from sklearn.linear_model import LinearRegression
regr=LinearRegression()
regr.fit(X_train,y_train)
print(regr.intercept_)
coeff_df=pd.DataFrame(regr.coef_,x.columns,columns=['coefficient'])
coeff_df
```

7211.098088897498

Out[28]:

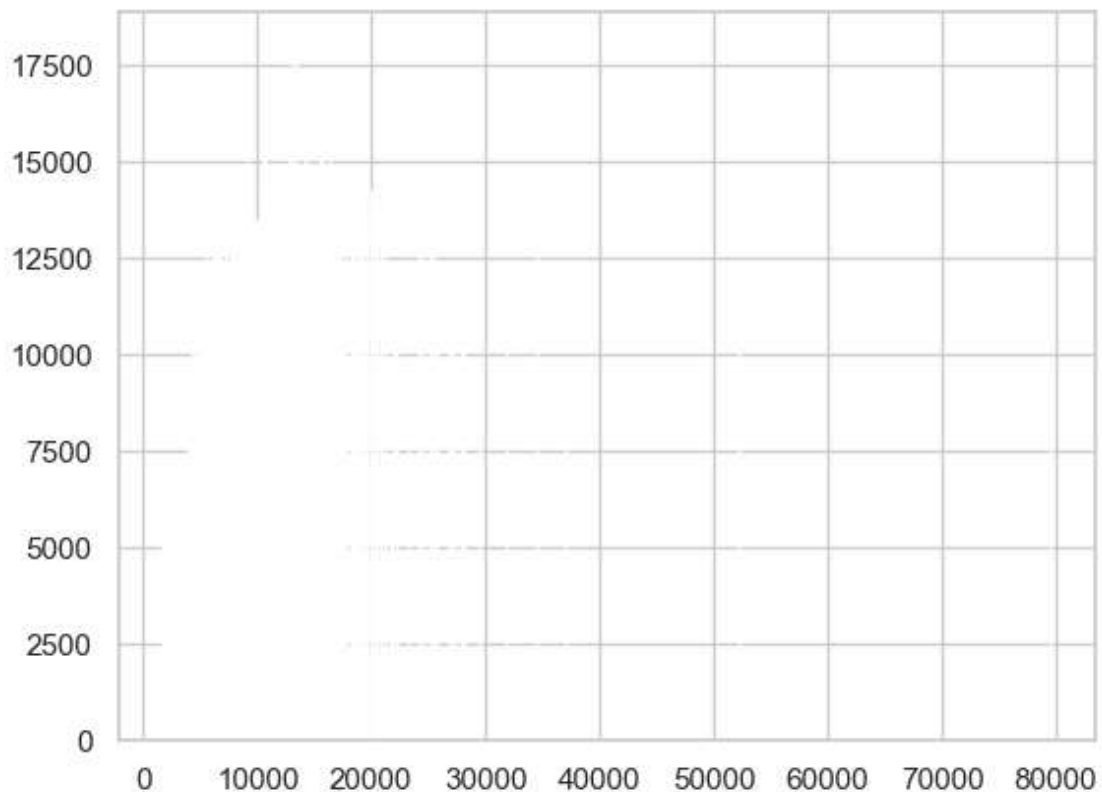|  | coefficient |
|---|---|
| **Airline** | -418.483922 |
| **Source** | -3275.073380 |
| **Destination** | 2505.480291 |
| **Total_Stops** | 3541.798053 |

In [29]:
```python
score=regr.score(X_test,y_test)
print(score)
```

0.4108304890928346

In [30]:
```python
predictions=regr.predict(X_test)
```

In [34]:
```python
plt.bar(y_test,predictions)
```

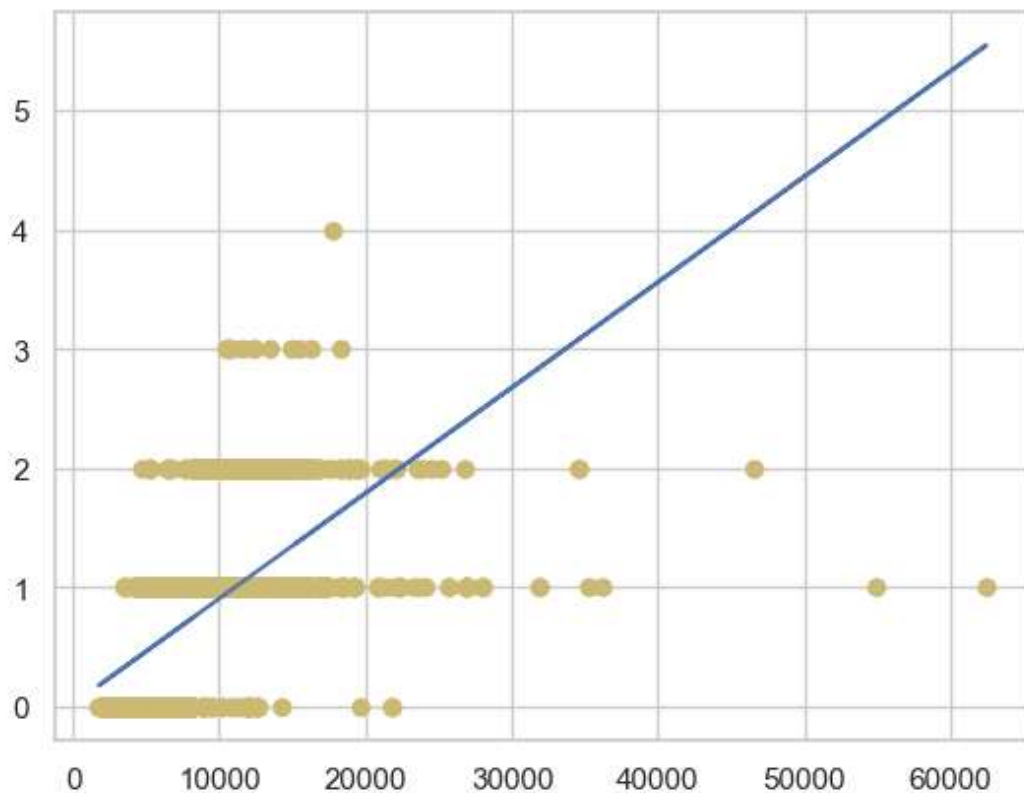Out[34]: <BarContainer object of 3205 artists>

```
In [35]: x=np.array(fdf['Price']).reshape(-1,1)
         y=np.array(fdf['Total_Stops']).reshape(-1,1)
         fdf.dropna(inplace=True)
```

```
In [37]: X_train,X_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
         regr.fit(X_train,y_train)
         regr.fit(X_train,y_train)
```

```
Out[37]:  ▾ LinearRegression

          LinearRegression()
```

```
In [38]: y_pred=regr.predict(X_test)
         plt.scatter(X_test,y_test,color='y')
         plt.plot(X_test,y_pred,color='b')
         plt.show()
```



# Logistic Regression

```
In [39]: x=np.array(fdf['Price']).reshape(-1,1)
         y=np.array(fdf['Total_Stops']).reshape(-1,1)
         fdf.dropna(inplace=True)
         x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=1)
         from sklearn.linear_model import LogisticRegression
         lr=LogisticRegression(max_iter=10000)
```

In [40]: `lr.fit(x_train,y_train)`
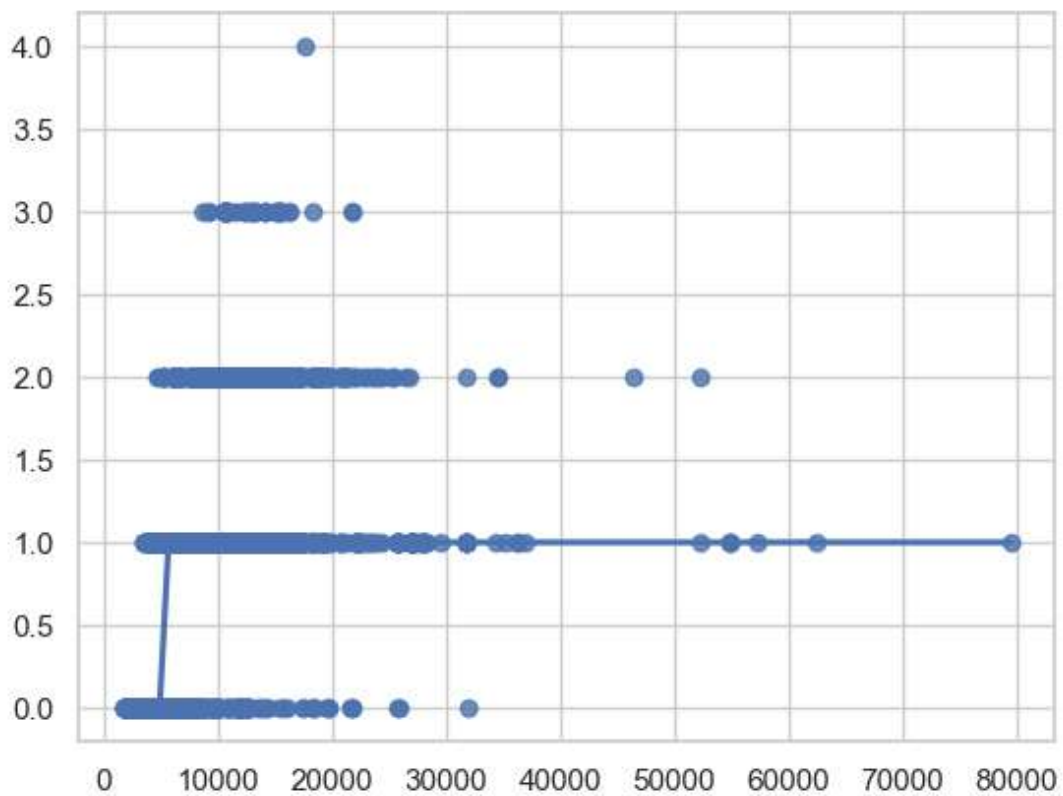
Out[40]:
```
        ▼              LogisticRegression
        LogisticRegression(max_iter=10000)
```

In [41]:
```
score=lr.score(x_test,y_test)
print(score)
```

```
0.7160686427457098
```

In [42]: `sns.regplot(x=x,y=y,data=fdf,logistic=True,ci=None)`

Out[42]: `<Axes: >`



# Decision tree

In [43]:
```
from sklearn.tree import DecisionTreeClassifier
clf=DecisionTreeClassifier(random_state=0)
clf.fit(x_train,y_train)
```

Out[43]:
```
        ▼            DecisionTreeClassifier
        DecisionTreeClassifier(random_state=0)
```

```
In [44]: score=clf.score(x_test,y_test)
         print(score)
```

```
0.9369734789391576
```

# Random Forest

```
In [45]: from sklearn.ensemble import RandomForestClassifier
         rfc=RandomForestClassifier()
         rfc.fit(X_train,y_train)
```

Out[45]:
```
▼ RandomForestClassifier

RandomForestClassifier()
```

```
In [46]: params={'max_depth':[2,3,5,10,20],
         'min_samples_leaf':[5,10,20,50,100,200],'n_estimators':[10,25,30,50,100,200]}
```

```
In [47]: from sklearn.model_selection import GridSearchCV
         grid_search=GridSearchCV(estimator=rfc,param_grid=params,cv=2,scoring="accuracy")
```

```
In [48]: grid_search.fit(X_train,y_train)
```

Out[48]:
```
▸            GridSearchCV

▸ estimator: RandomForestClassifier

     ▸ RandomForestClassifier
```
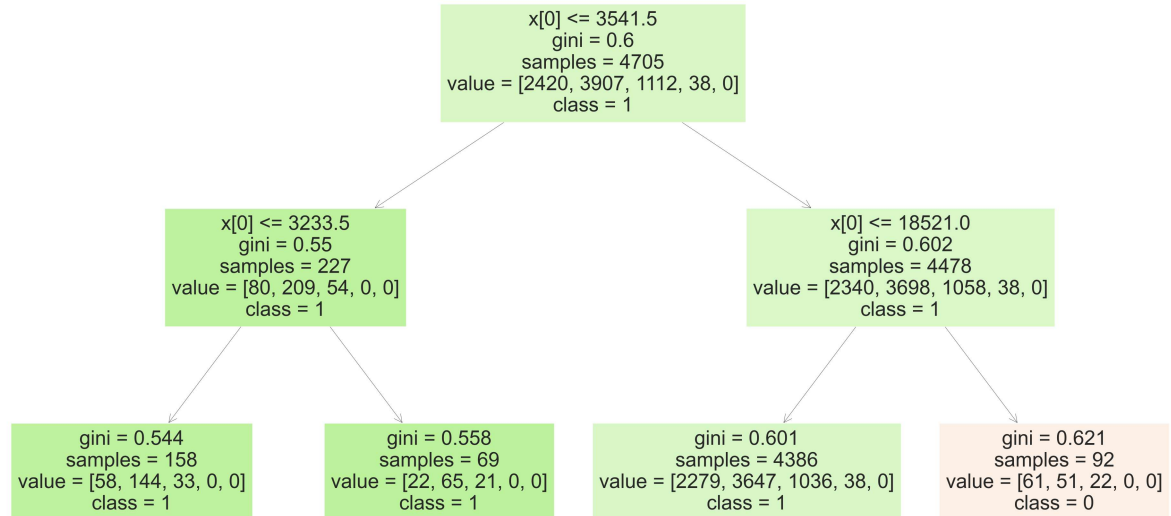
```
In [49]: grid_search.best_score_
```

Out[49]: 0.523605715699528

```
In [50]: rf_best=grid_search.best_estimator_
         rf_best
```

Out[50]:
```
▼                        RandomForestClassifier

RandomForestClassifier(max_depth=2, min_samples_leaf=50, n_estimators=25)
```

In [51]:
```python
from sklearn.tree import plot_tree
plt.figure(figsize=(80,40))
plot_tree(rf_best.estimators_[4],class_names=['0','1','2','3','4'],filled=True);
```

```
                              x[0] <= 3541.5
                               gini = 0.6
                             samples = 4705
                      value = [2420, 3907, 1112, 38, 0]
                               class = 1


          x[0] <= 3233.5                              x[0] <= 18521.0
           gini = 0.55                                 gini = 0.602
         samples = 227                               samples = 4478
   value = [80, 209, 54, 0, 0]                 value = [2340, 3698, 1058, 38, 0]
           class = 1                                    class = 1


  gini = 0.544        gini = 0.558        gini = 0.601              gini = 0.621
 samples = 158       samples = 69        samples = 4386            samples = 92
value = [58, 144,   value = [22, 65,   value = [2279, 3647,   value = [61, 51,
  33, 0, 0]           21, 0, 0]           1036, 38, 0]            22, 0, 0]
  class = 1           class = 1            class = 1              class = 0
```

In [52]:
```python
score=rfc.score(x_test,y_test)
print(score)
```

```
0.48174726989079564
```

# Conclusion:The above implemented models "Decision Tree" is high accuracy score.So it is the best model

In [ ]: