# Problem statement:

To predict the best model for the given dataset based on accuracy.

In [38]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings("ignore")
```

# Data collection

In [3]:

```
1  train_df=pd.read_csv(r"C:\Users\91949\Desktop\train_data.csv")
2  train_df
```

Out[3]:

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

10683 rows × 11 columns

In [4]:

```python
test_df=pd.read_csv(r"C:\Users\91949\Desktop\test_data.csv")
test_df
```

Out[4]:

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

2671 rows × 10 columns

# Data Cleaning and Preprocessing

In [6]:

```
1  train_df.head()
```

Out[6]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration |
|---|---|---|---|---|---|---|---|---|
| 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 |

In [7]:

```
1  train_df.tail()
```

Out[7]:

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

In [8]:

```
1 test_df.head()
```

Out[8]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration |
|---|---|---|---|---|---|---|---|---|
| 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 |

In [9]:

```
1 test_df.tail()
```

Out[9]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Durati |
|---|---|---|---|---|---|---|---|---|
| 2666 | Air India | 6/06/2019 | Kolkata | Banglore | CCU ? DEL ? BLR | 20:30 | 20:25 07 Jun | 23h 5! |
| 2667 | IndiGo | 27/03/2019 | Kolkata | Banglore | CCU ? BLR | 14:20 | 16:55 | 2h 3! |
| 2668 | Jet Airways | 6/03/2019 | Delhi | Cochin | DEL ? BOM ? COK | 21:50 | 04:25 07 Mar | 6h 3! |
| 2669 | Air India | 6/03/2019 | Delhi | Cochin | DEL ? BOM ? COK | 04:00 | 19:15 | 15h 1! |
| 2670 | Multiple carriers | 15/06/2019 | Delhi | Cochin | DEL ? BOM ? COK | 04:55 | 19:15 | 14h 2( |

In [10]:

```
1  train_df.describe()
```

Out[10]:

|       | 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 [11]:

```
1  test_df.describe()
```

Out[11]:

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

In [12]:

```
1  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 [13]:

```
1  test_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2671 entries, 0 to 2670
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   Airline          2671 non-null   object
 1   Date_of_Journey  2671 non-null   object
 2   Source           2671 non-null   object
 3   Destination      2671 non-null   object
 4   Route            2671 non-null   object
 5   Dep_Time         2671 non-null   object
 6   Arrival_Time     2671 non-null   object
 7   Duration         2671 non-null   object
 8   Total_Stops      2671 non-null   object
 9   Additional_Info  2671 non-null   object
dtypes: object(10)
memory usage: 208.8+ KB
```

In [14]:

```
1  train_df.size
```

Out[14]:

117513

In [15]:

```
1  test_df.size
```

Out[15]:

26710

In [16]:

```
1  train_df.shape
```

Out[16]:

(10683, 11)

In [17]:

```
1  test_df.shape
```

Out[17]:

(2671, 10)

# Exploratory Data Analysis

In [18]:

```
1  train_df.isnull().sum()
```

Out[18]:

```
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 [19]:

```
1  test_df.isnull().sum()
```

Out[19]:

```
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 [20]:

```
1  train_df.duplicated().sum()
```

Out[20]:

```
220
```

In [21]:

```
1  test_df.duplicated().sum()
```

Out[21]:

```
26
```

In [22]:

```
1  train_df.dropna(inplace=True)
```

In [23]:

```
1  train_df.isnull().sum()
```

Out[23]:

```
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 [24]:

```
1  train_df['Airline'].value_counts()
```

Out[24]:

```
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: Airline, dtype: int64
```

In [25]:

```
1  train_df['Source'].value_counts()
```

Out[25]:

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

In [26]:

```
1  train_df['Destination'].value_counts()
```

Out[26]:

```
Cochin      4536
Banglore    2871
Delhi       1265
New Delhi    932
Hyderabad    697
Kolkata      381
Name: Destination, dtype: int64
```

In [27]:

```
1  train_df['Total_Stops'].value_counts()
```

Out[27]:

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

In [28]:

```
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[28]:

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

10682 rows × 11 columns

In [29]:

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

Out[29]:

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

10682 rows × 11 columns

In [30]:

```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[30]:

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

10682 rows × 11 columns

In [31]:

```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[31]:

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

10682 rows × 11 columns

# Data visualization

In [32]:

```python
import seaborn as sns
df=train_df[['Airline','Source','Destination','Total_Stops','Price']]
sns.heatmap(df.corr(),annot=True)
```

Out[32]:

```
<Axes: >
```



In [33]:

```python
x=df[['Airline','Source','Destination','Total_Stops']]
y=df['Price']
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)
```

In [34]:

```python
#Data prediction and Evaluation
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
```

7211.098088897488

Out[34]:

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

In [35]:

```python
score=regr.score(x_test,y_test)
print(score)
```

0.4108304890928348

In [36]:

```
1  predictions=regr.predict(x_test)
2  plt.scatter(y_test,predictions)
```

Out[36]:

```
<matplotlib.collections.PathCollection at 0x23f24a67710>
```



In [39]:

```
1  x=np.array(df['Price']).reshape(-1,1)
2  y=np.array(df['Total_Stops']).reshape(-1,1)
3  df.dropna(inplace=True)
```

In [40]:

```
1  x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
2  regr.fit(x_train,y_train)
3  regr.fit(x_train,y_train)
```
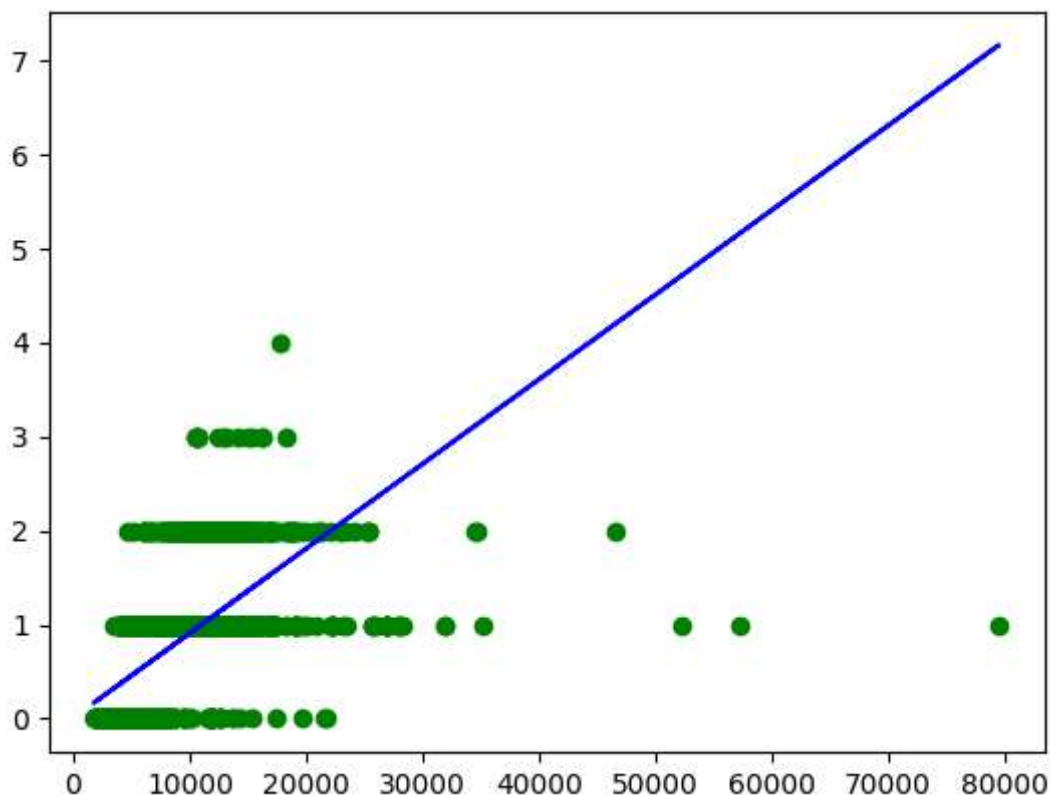
Out[40]:

```
▼ LinearRegression
LinearRegression()
```

In [41]:

```python
y_pred=regr.predict(x_test)
plt.scatter(x_test,y_test,color='g')
plt.plot(x_test,y_pred,color='b')
plt.show()
```



## Logistic Regression

In [42]:

```python
x=np.array(df['Price']).reshape(-1,1)
y=np.array(df['Total_Stops']).reshape(-1,1)
df.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)
import warnings
warnings.simplefilter(action='ignore')
```

In [43]:

```python
lr.fit(x_train,y_train)
```

Out[43]:

```
▼       LogisticRegression
LogisticRegression(max_iter=10000)
```

In [44]:

```
1  score=lr.score(x_test,y_test)
2  print(score)
```
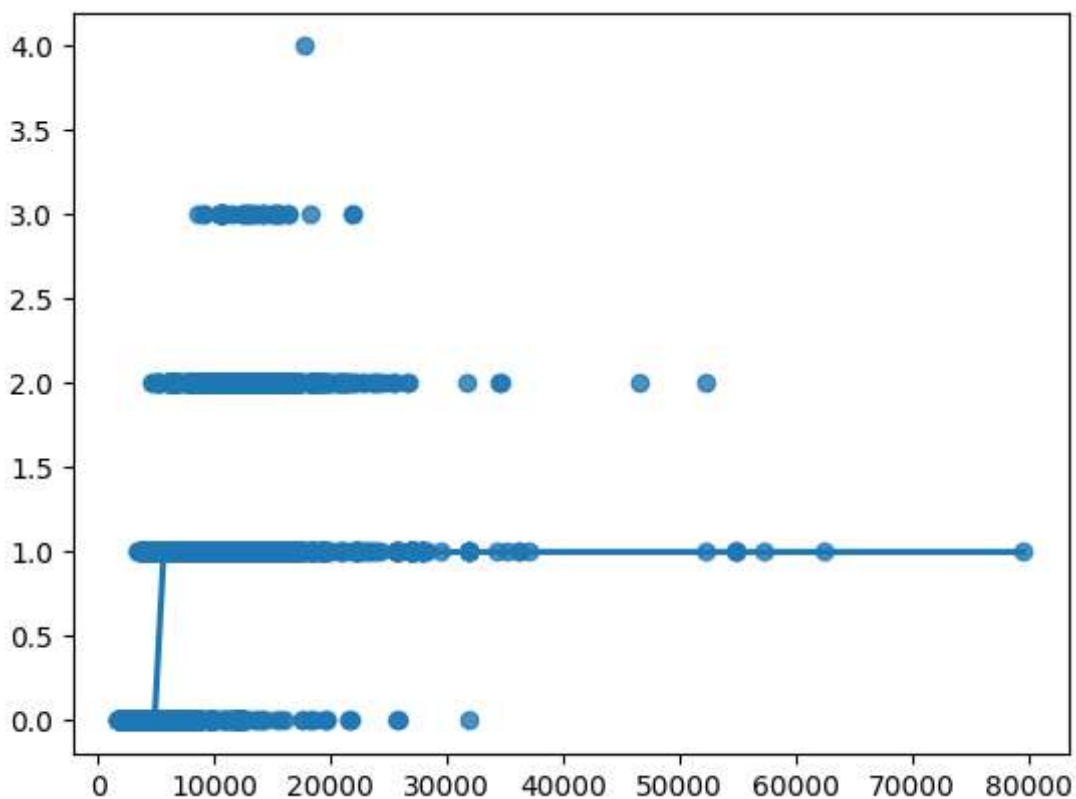
0.7160686427457098

In [45]:

```
1  sns.regplot(x=x,y=y,data=df,logistic=True,ci=None)
2  plt.show
```

Out[45]:

<function matplotlib.pyplot.show(close=None, block=None)>



## Decision Tree

In [46]:

```
1  from sklearn.tree import DecisionTreeClassifier
2  clf=DecisionTreeClassifier(random_state=0)
3  clf.fit(x_train,y_train)
```

Out[46]:

```
▼          DecisionTreeClassifier
DecisionTreeClassifier(random_state=0)
```

In [47]:

```python
score=clf.score(x_test,y_test)
print(score)
```

0.9369734789391576

## Random Forest

In [48]:

```python
from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[48]:

```
▾ RandomForestClassifier
RandomForestClassifier()
```

In [49]:

```python
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 [50]:

```python
from sklearn.model_selection import GridSearchCV
grid_search=GridSearchCV(estimator=rfc,param_grid=params,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

Out[50]:

```
▸           GridSearchCV
▸ estimator: RandomForestClassifier
      ▸ RandomForestClassifier
```

In [51]:

```python
grid_search.best_score_
```

Out[51]:

0.873745472898494

In [52]:

```python
1  rf_best=grid_search.best_estimator_
2  rf_best
```

Out[52]:

```
▼                    RandomForestClassifier
RandomForestClassifier(max_depth=20, min_samples_leaf=5, n_estimators=5
0)
```

## Conclusion:

we conclude that "Decision Tree" is the best model for Flight Price Prediction dataset,because it got highest accuracy compared to other models.