

# Problem Statement

we have to perform the models. And predict the which model is best fit and best accuracy for flightprice dataset

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn import metrics
```

In [2]:

```
train_df=pd.read_csv(r"C:\Users\shaha\OneDrive\Desktop\Excel\Cop of Data_Train.csv")
train_df
```

Out[2]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Dura
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



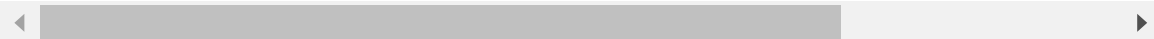
In [3]:

```
test_df=pd.read_csv(r"C:\Users\shaha\OneDrive\Desktop\202U1A0597\Test_set.csv")
test_df
```

Out[3]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Durat
0	Jet Airways	6/06/2019	Delhi	Cochin	DEL ? BOM ? COK	17:30	04:25 07 Jun	10h 5
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 4
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 5
...	...	...	...	...	...	...	...	...
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

2671 rows × 10 columns



In [4]:

```
train_df.head()
```

Out[4]:

	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 [5]:

```
train_df.tail()
```

Out[5]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Dura
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 [6]:

```
test_df.head()
```

Out[6]:

	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 [7]:

```
test_df.tail()
```

Out[7]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duratic
2666	Air India	6/06/2019	Kolkata	Banglore	CCU ? DEL ? BLR	20:30	20:25 07 Jun	23h 55
2667	IndiGo	27/03/2019	Kolkata	Banglore	CCU ? BLR	14:20	16:55	2h 35
2668	Jet Airways	6/03/2019	Delhi	Cochin	DEL ? BOM ? COK	21:50	04:25 07 Mar	6h 35
2669	Air India	6/03/2019	Delhi	Cochin	DEL ? BOM ? COK	04:00	19:15	15h 15
2670	Multiple carriers	15/06/2019	Delhi	Cochin	DEL ? BOM ? COK	04:55	19:15	14h 20

In [8]:

```
train_df.shape
```

Out[8]:

(10683, 11)

In [9]:

```
test_df.shape
```

Out[9]:

(2671, 10)

In [10]:

```
train_df.describe()
```

Out[10]:

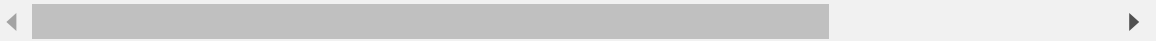
	Price
count	10683.000000
mean	9087.064121
std	4611.359167
min	1759.000000
25%	5277.000000
50%	8372.000000
75%	12373.000000
max	79512.000000

In [11]:

```
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	2
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]:

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]:

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]:

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

Out[14]:

```
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 [15]:

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

Out[15]:

```
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
```

Removing null values in Train dataset

In [16]:

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



In [17]:

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

Out[17]:

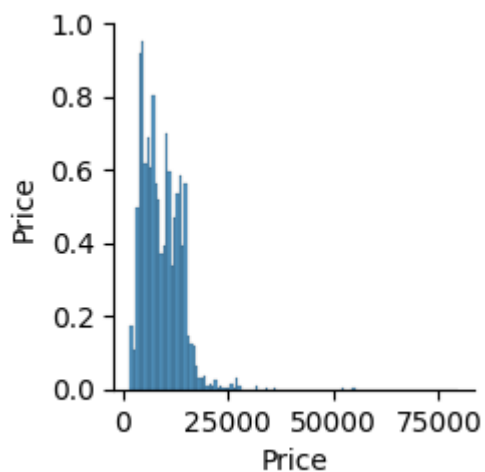
```
Airline           0
Date_of_Journey   0
Source            0
Destination       0
Route             0
Dep_Time          0
Arrival_Time      0
Duration          0
Total_Stops       0
Additional_Info    0
Price             0
dtype: int64
```

In [18]:

```
sns.pairplot(train_df)
```

Out[18]:

<seaborn.axisgrid.PairGrid at 0x125ee37afe0>



Converting string into numerical values

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]:

```
A={"Airline":{"Jet Airways":0,"IndiGo":1,"Air India":2,"Multiple carriers":3,"SpiceJet":4,"Multiple carriers Premium economy":8,"Jet Airways Business":9,"Vistara Premium":10}}
train_df=train_df.replace(A)
train_df
```

Out[20]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Durat
0	1	24/03/2019	Banglore	New Delhi	BLR ? DEL	22:20	01:10 22 Mar	2h 5
1	2	1/05/2019	Kolkata	Banglore	CCU ? IXR ? BBI ? BLR	05:50	13:15	7h 2
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 2
4	1	01/03/2019	Banglore	New Delhi	BLR ? NAG ? DEL	16:50	21:35	4h 4
...	...	...	...	...	...	...	...	
10678	6	9/04/2019	Kolkata	Banglore	CCU ? BLR	19:55	22:25	2h 3
10679	2	27/04/2019	Kolkata	Banglore	CCU ? BLR	20:45	23:20	2h 3
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 4
10682	2	9/05/2019	Delhi	Cochin	DEL ? GOI ? BOM ? COK	10:55	19:15	8h 2

10682 rows × 11 columns



In [21]:

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

Out[21]:

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

In [22]:

```
S={"Source":{"Delhi":1,"Kolkata":2,"Banglore":3,"Mumbai":4,"Chennai":5}}
train_df=train_df.replace(S)
train_df
```

Out[22]:

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

10682 rows × 11 columns



In [23]:

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

Out[23]:

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

In [24]:

```
D={"Destination":{"Cochin":1,"Banglore":2,"Delhi":3,"New Delhi":4,"Hyderabad":5,"Kolkata":6}
train_df=train_df.replace(D)
train_df
```

Out[24]:

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

10682 rows × 11 columns



In [25]:

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

Out[25]:

Total\_Stops

1 stop            5625

non-stop        3491

2 stops         1520

3 stops          45

4 stops          1

Name: count, dtype: int64



In [26]:

```
T={"Total_Stops":{"1 stop":1,"non-stop":0,"2 stops":2,"3 stops":3,"4 stops":4}}
train_df=train_df.replace(T)
train_df
```

Out[26]:

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

10682 rows × 11 columns

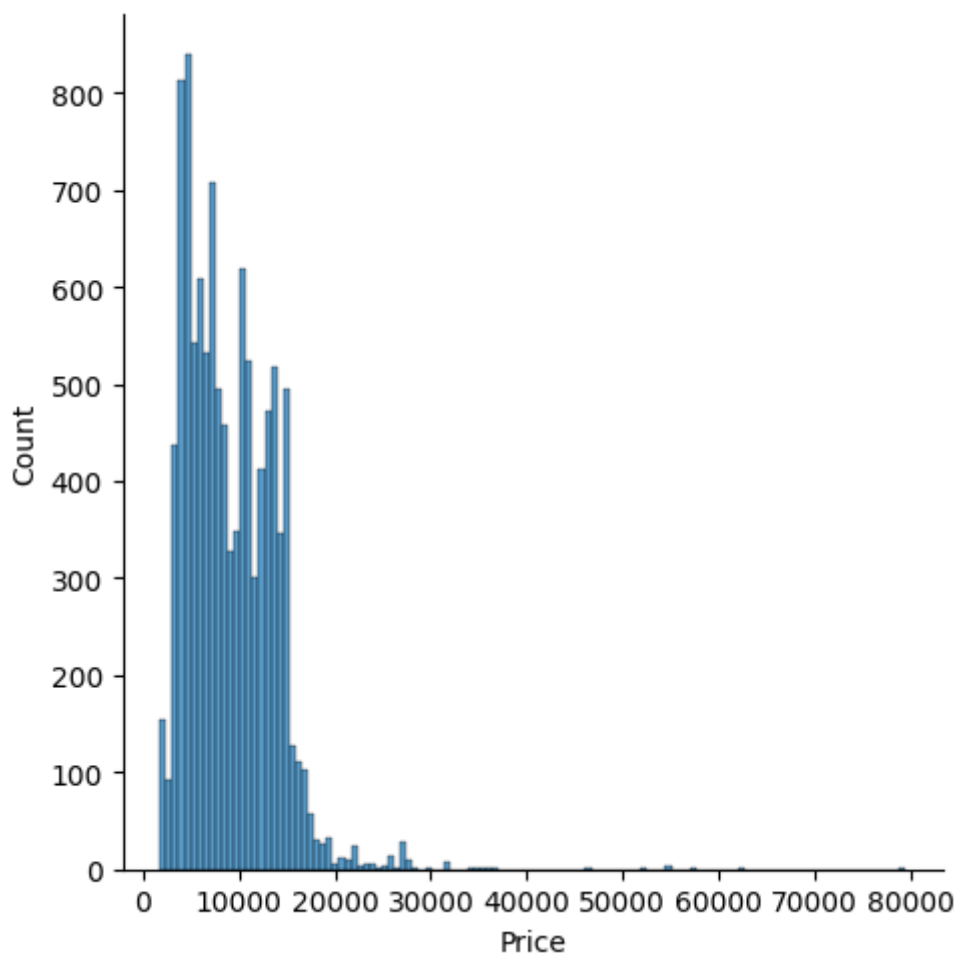


In [27]:

```
sns.displot(train_df['Price'])
```

Out[27]:

<seaborn.axisgrid.FacetGrid at 0x125ee379120>

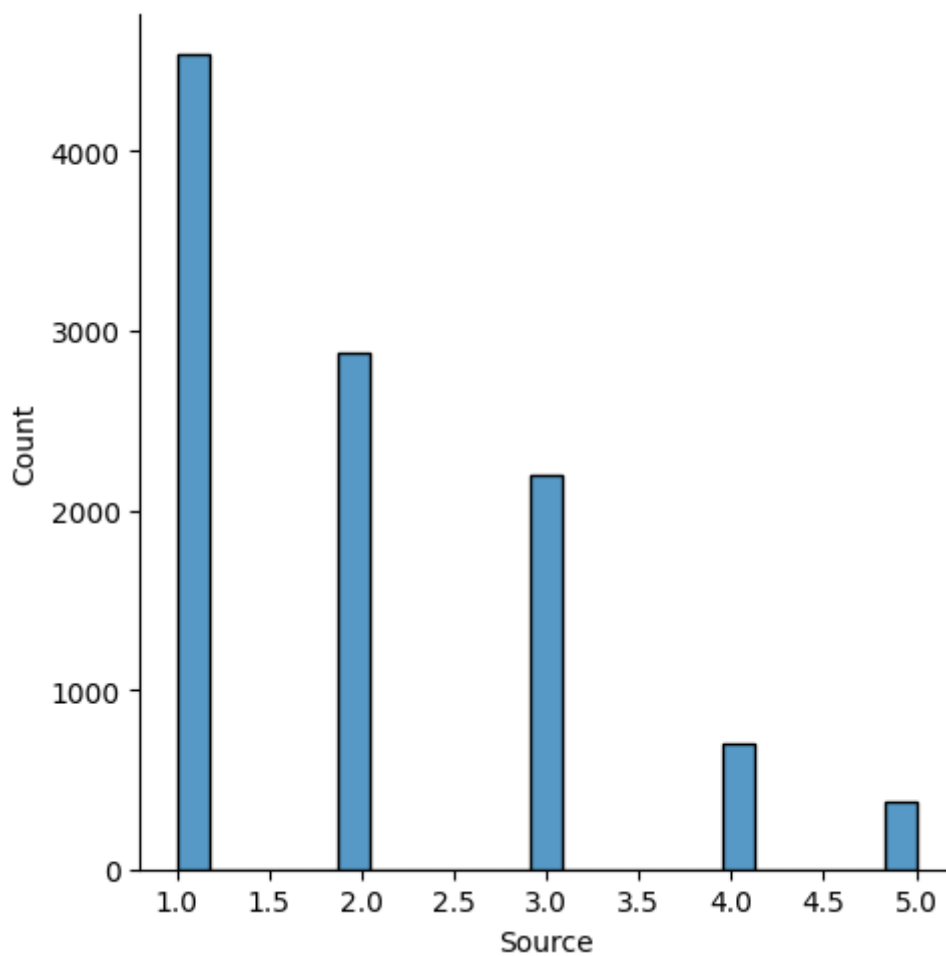


In [28]:

```
sns.displot(train_df['Source'])
```

Out[28]:

<seaborn.axisgrid.FacetGrid at 0x125f1858a90>



## Linear Regression

In [29]:

```
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
```

In [30]:

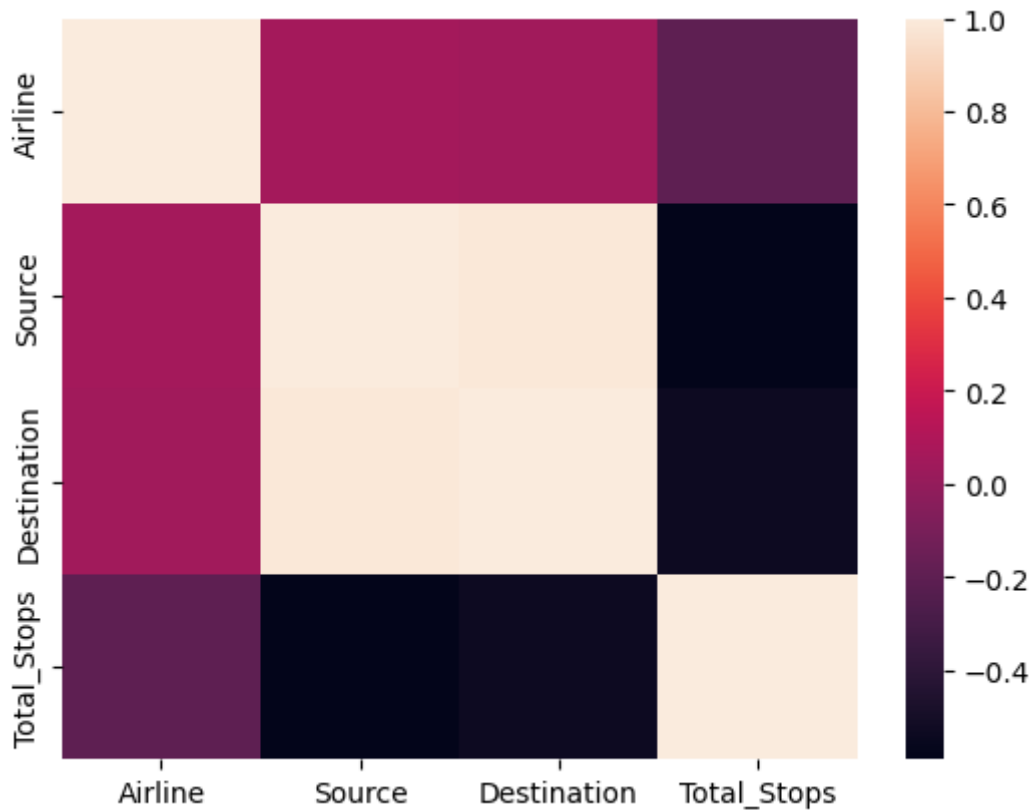
```
x=train_df[['Airline', 'Source', 'Destination', 'Total_Stops']]
y=train_df['Price']
```

In [31]:

```
sns.heatmap(x.corr())
```

Out[31]:

&lt;Axes: &gt;



In [32]:

```
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=100)
```

In [33]:

```
regr = LinearRegression()
regr.fit(x_train,y_train)
print(regr.intercept_)
coeff_train_df=pd.DataFrame(regr.coef_,x.columns,columns=['coefficient'])
coeff_train_df
```

7980.6911780504215

Out[33]:

	coefficient
<b>Airline</b>	-418.483922
<b>Source</b>	-3275.073380
<b>Destination</b>	2505.480291
<b>Total_Stops</b>	3541.798053

In [34]:

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

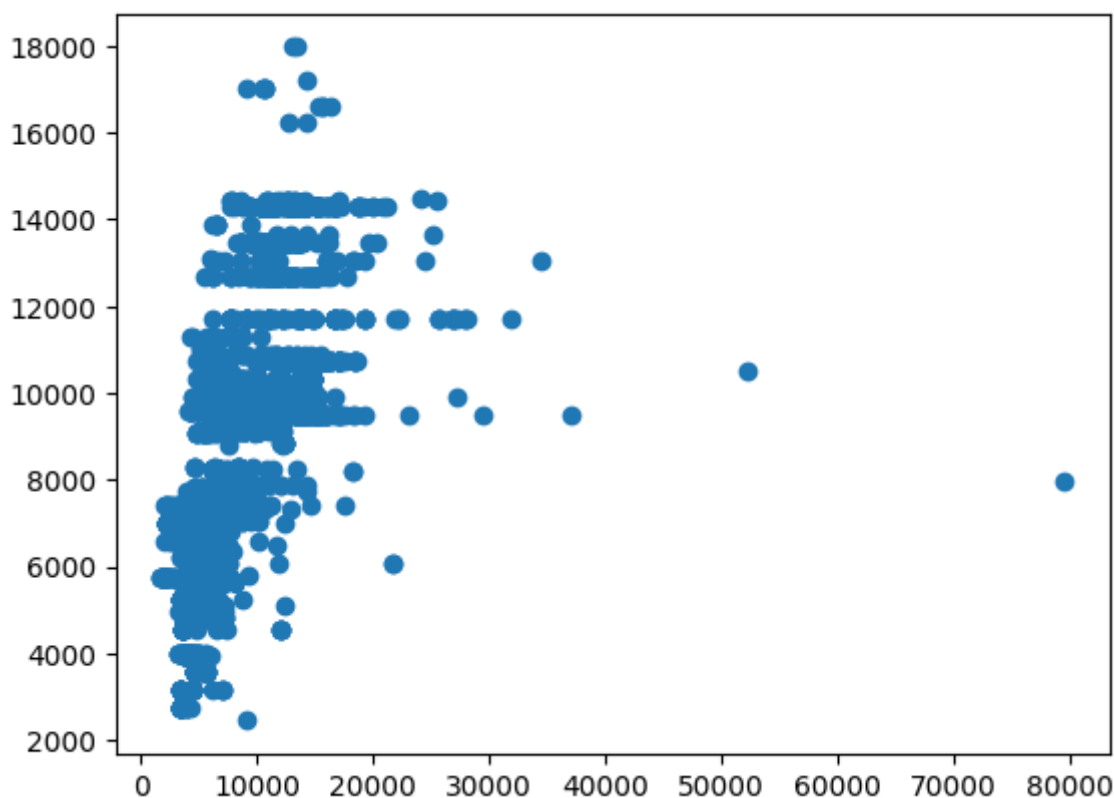
0.41083048909283437

In [35]:

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

Out[35]:

&lt;matplotlib.collections.PathCollection at 0x125f1b75c90&gt;



In [36]:

```
x=np.array(train_df['Price']).reshape(-1,1)
y=np.array(train_df['Total_Stops']).reshape(-1,1)
```

In [37]:

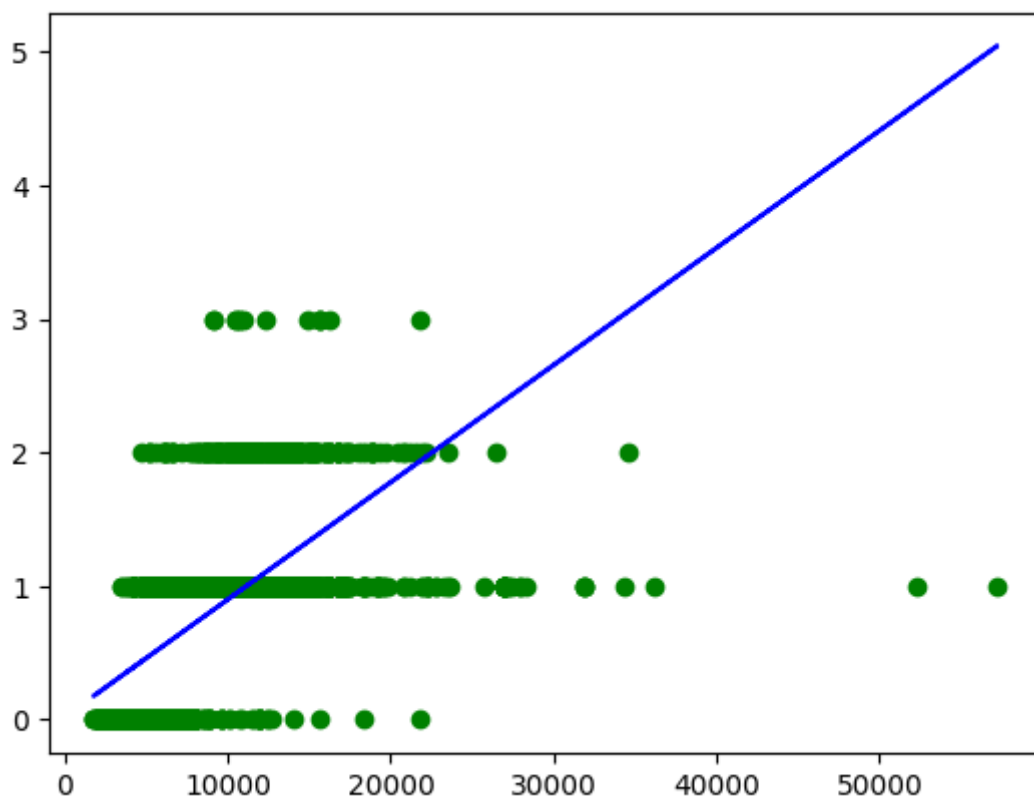
```
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
regr.fit(x_train,y_train)
regr.fit(x_test,y_test)
```

Out[37]:

```
▼ LinearRegression
LinearRegression()
```

In [38]:

```
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 [39]:

```
from sklearn.linear_model import LogisticRegression
x=np.array(train_df['Price']).reshape(-1,1)
y=np.array(train_df['Destination']).reshape(-1,1)
train_df.dropna(inplace=True)
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=1)
lr=LogisticRegression(max_iter=100000)
import warnings
warnings.simplefilter(action='ignore')
```

In [40]:

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

Out[40]:

```
LogisticRegression
LogisticRegression(max_iter=100000)
```

In [41]:

```
score=lr.score(x_test,y_test)
print(score)
```

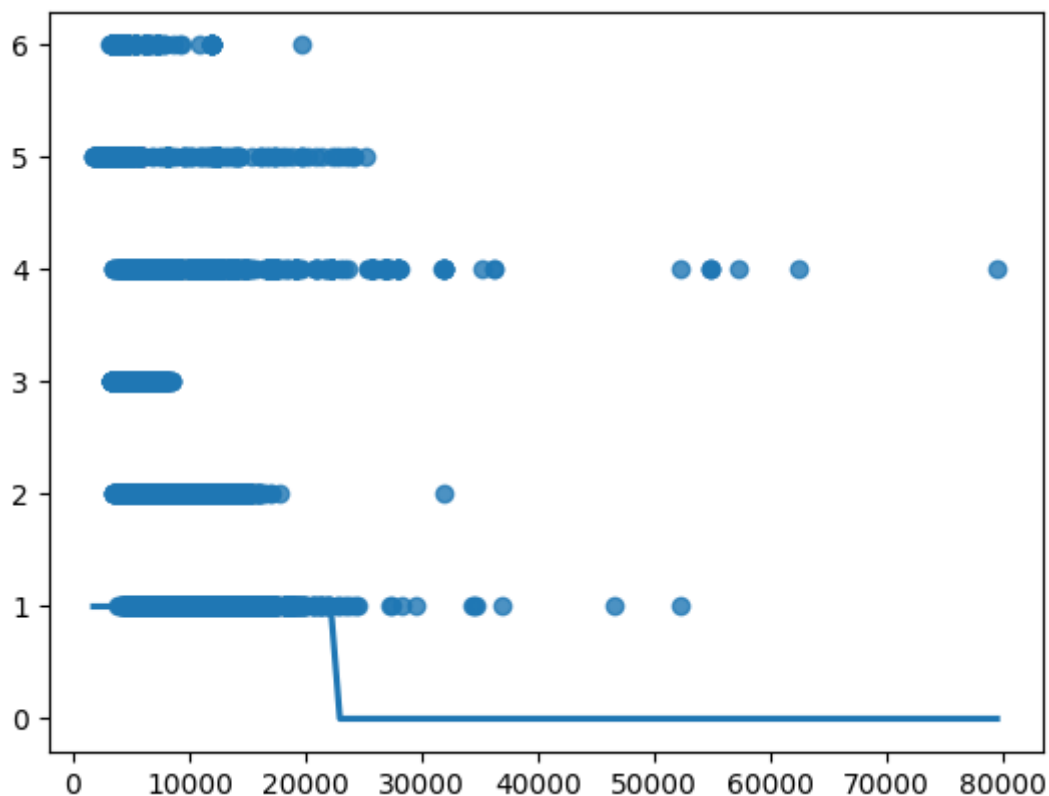
0.431201248049922

In [42]:

```
sns.regplot(x=x,y=y,data=train_df,logistic=True,ci=None)
```

Out[42]:

&lt;Axes: &gt;



## 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)
```

0.921996879875195

## Random Forest

In [45]:

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

Out[45]:

```
▼ RandomForestClassifier
RandomForestClassifier()
```

In [52]:

```
rf=RandomForestClassifier()
```

In [53]:

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

In [54]:

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

Out[54]:

```
► GridSearchCV
► estimator: RandomForestClassifier
  ► RandomForestClassifier
```

In [55]:

```
grid_search.best_score_
```

Out[55]:

0.8217206355693483



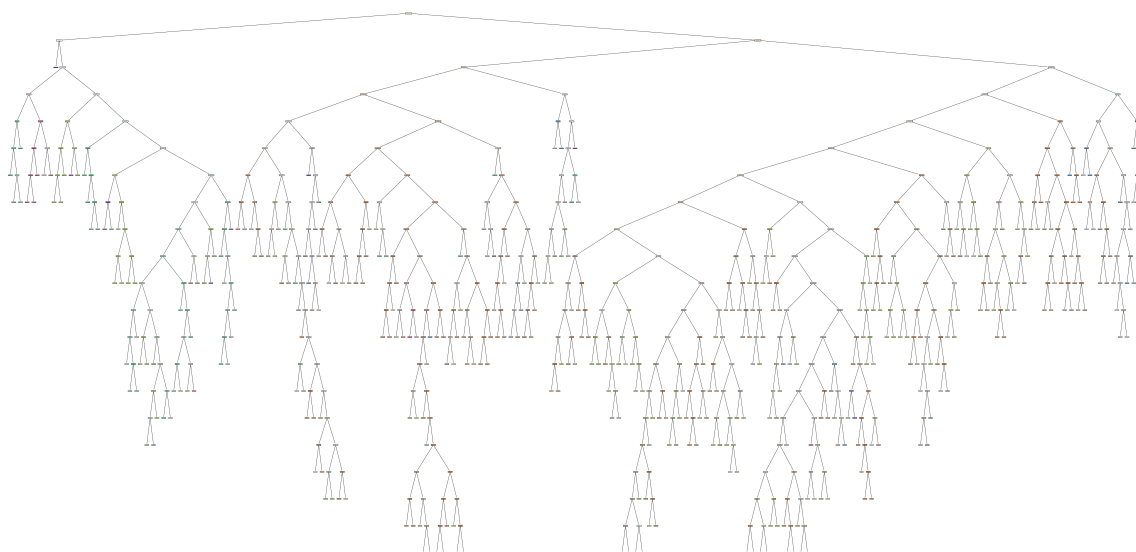
In [56]:

```
rf_best=grid_search.best_estimator_  
print(rf_best)
```

```
RandomForestClassifier(max_depth=20, min_samples_leaf=5, n_estimators=30)
```

In [66]:

```
from sklearn.tree import plot_tree  
from sklearn.tree import DecisionTreeClassifier  
import matplotlib.pyplot as plt  
plt.figure(figsize=(80,40))  
plot_tree(rf_best.estimators_[5],filled=True);
```



In [67]:

```
score=rfc.score(x_test,y_test)  
print(score)
```

```
0.9213728549141965
```

## Conclusion

From The Given Flight Price Dataset,we have performed on different models like Linear Regression,Logistic Regression,Random Forest,Decision Tree.By observing the score or model prediction in this models, In Decision Tree model got the best score and best accuracy.

In [ ]: