

Problem Statement:

To check Which model is best suitable for Flight price Prediction Dataset

In [1]:



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

In [3]:

```
traindf=pd.read_csv(r"C:\Users\DELL\Downloads\Data_Train.csv")
traindf
```

Out[3]:

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

10683 rows × 11 columns

Data cleaning and preprocessing

In [4]:

```
traindf.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]:

```
traindf.tail()
```

Out[5]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration
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 00m
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 [6]:

```
traindf.describe()
```

Out[6]:

	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 [7]:

```
traindf.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 [8]:

```
traindf.shape
```

Out[8]:

```
(10683, 11)
```

In [9]:

```
traindf.duplicated().sum()
```

Out[9]:

```
220
```

In [10]:

```
traindf.columns
```

Out[10]:

```
Index(['Airline', 'Date_of_Journey', 'Source', 'Destination', 'Route',  
      'Dep_Time', 'Arrival_Time', 'Duration', 'Total_Stops',  
      'Additional_Info', 'Price'],  
      dtype='object')
```

In [11]:

```
traindf.isnull().sum()
```

Out[11]:

```
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 [12]:

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

In [13]:

```
traindf.isnull().sum()
```

Out[13]:

```
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 [14]:

```
traindf.shape
```

Out[14]:

```
(10682, 11)
```

In [15]:

```
traindf['Airline'].value_counts()
```

Out[15]:

```
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 [16]:

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

Out[16]:

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

In [17]:

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

Out[17]:

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

In [18]:



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

Out[18]:

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

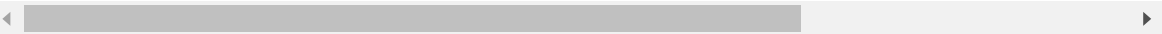
In [19]:

```
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}}
traindf=traindf.replace(airline)
traindf
```

Out[19]:

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

10682 rows × 11 columns



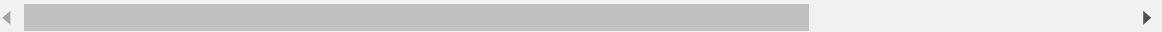
In [20]:

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

Out[20]:

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

10682 rows × 11 columns



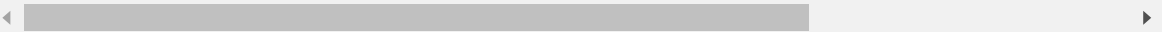
In [21]:

```
destination={"Destination":{"Cochin":0,"Banglore":1,"Delhi":2,
"New Delhi":3,"Hyderabad":4,"Kolkata":5}}
traindf=traindf.replace(destination)
traindf
```

Out[21]:

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

10682 rows × 11 columns



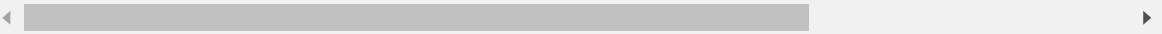
In [22]:

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

Out[22]:

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

10682 rows × 11 columns



In [23]:

```
traindf
```

Out[23]:

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

10682 rows × 11 columns

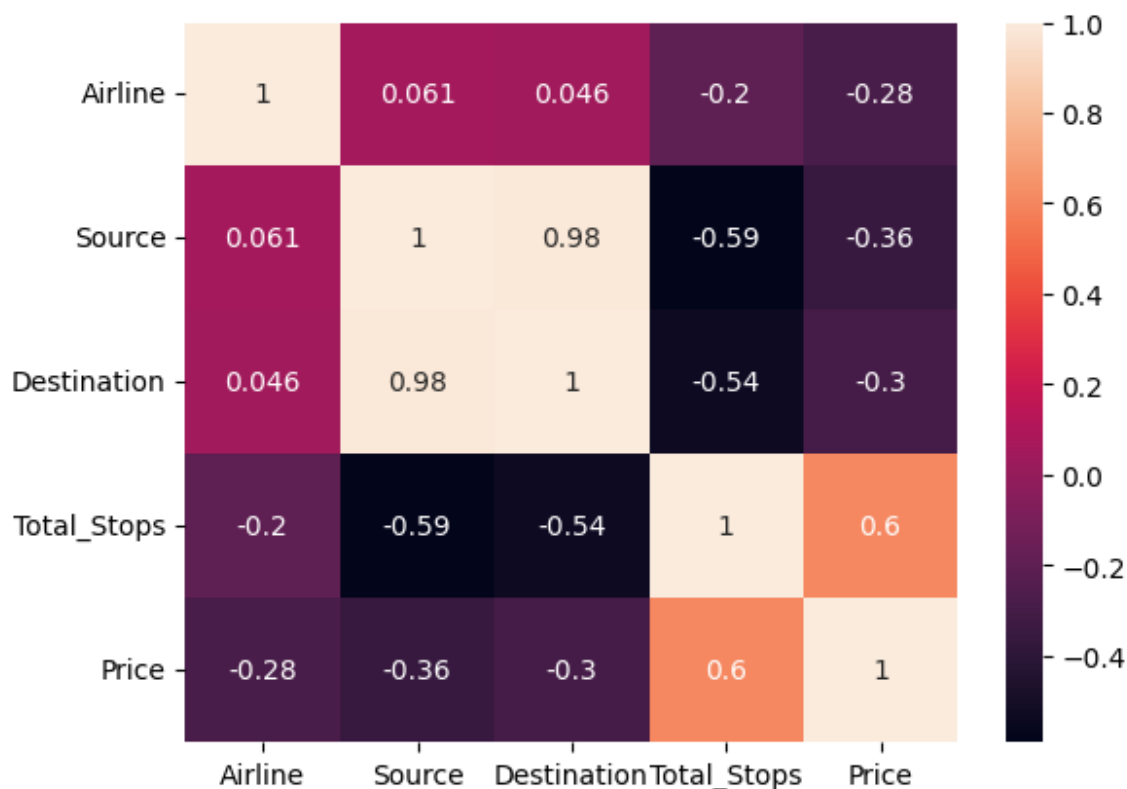
Exploratory Data Analysis

In [24]:

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

Out[24]:

<Axes: >



In [25]:

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

Linear Regression

In [26]:

```
#Linear Regression
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 [27]:



```
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[27]:

	coefficient
Airline	-418.483922
Source	-3275.073380
Destination	2505.480291
Total_Stops	3541.798053

In [28]:



```
#Linear Rgeression
score=regr.score(X_test,y_test)
print(score)
```

0.4108304890928348

In [29]:



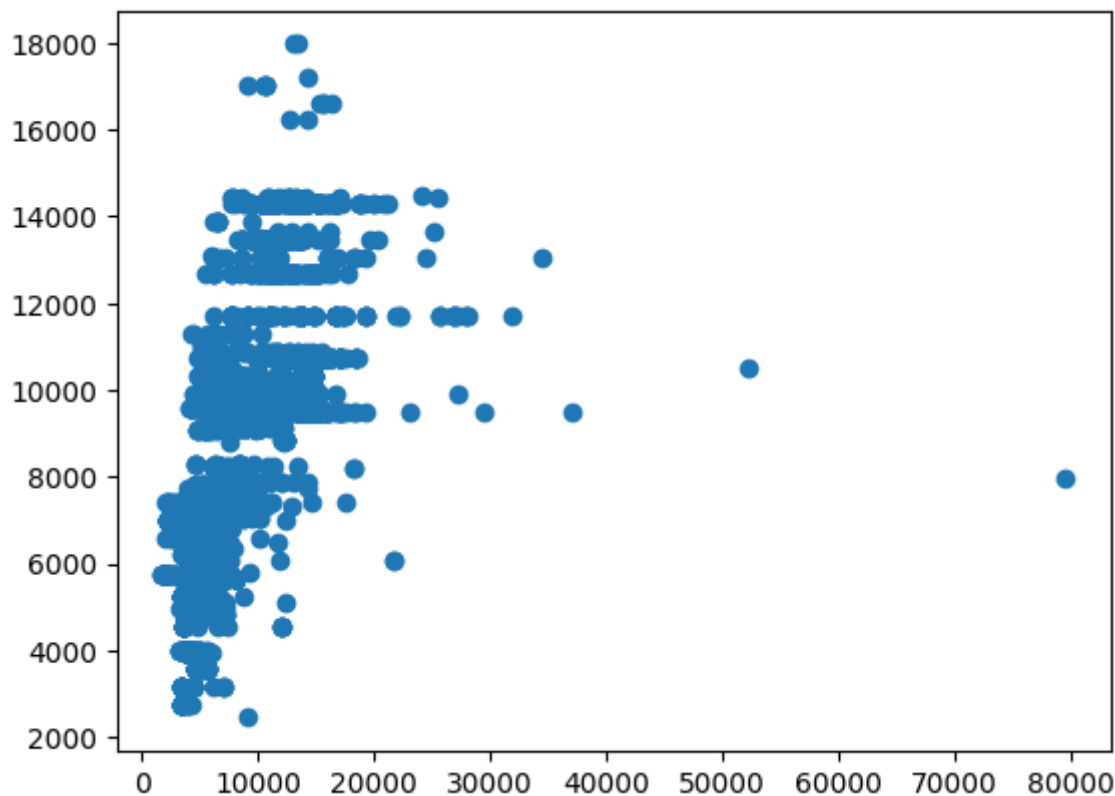
```
predictions=regr.predict(X_test)
```

In [30]:

```
plt.scatter(y_test,predictions)
```

Out[30]:

```
<matplotlib.collections.PathCollection at 0x27b92c8de10>
```



In [31]:

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

C:\Users\DELL\AppData\Local\Temp\ipykernel_9192\521034954.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
fdf.dropna(inplace=True)

In [32]:

```
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)
```

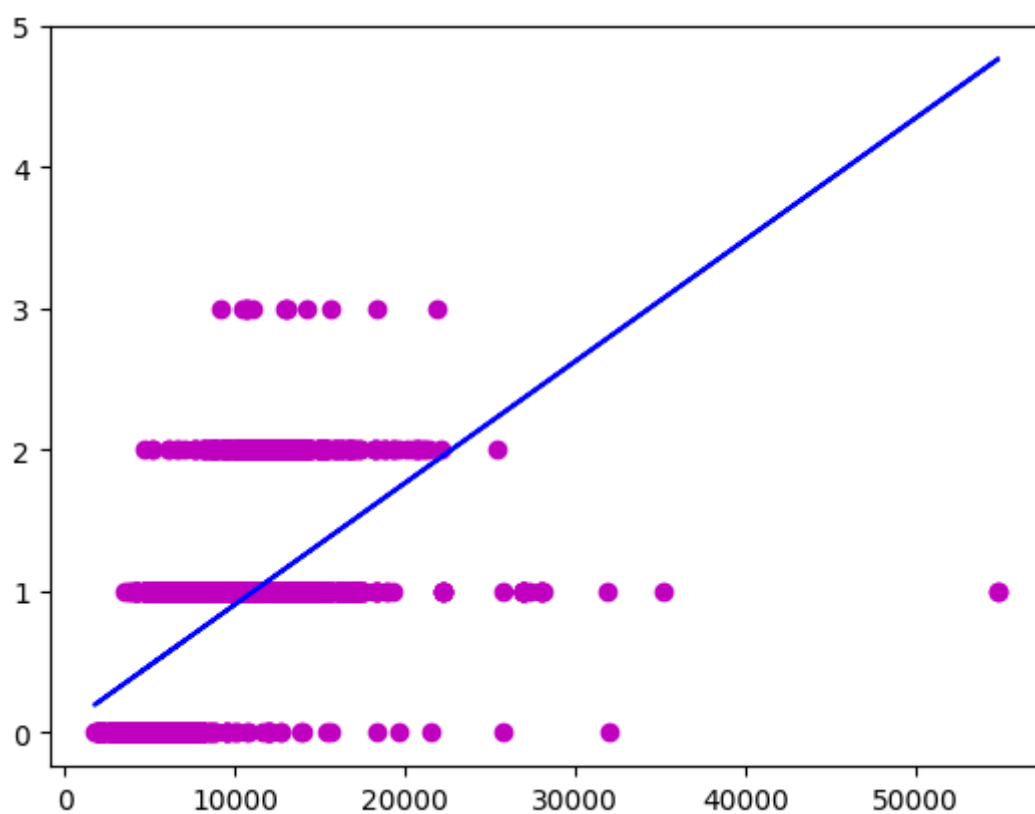
Out[32]:

▼ LinearRegression

LinearRegression()

In [35]:

```
y_pred=regr.predict(X_test)
plt.scatter(X_test,y_test,color='m')
plt.plot(X_test,y_pred,color='b')
plt.show()
```



Logistic Regression

In [36]:

```
#Logistic Regression
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)
```

C:\Users\DELL\AppData\Local\Temp\ipykernel_9192\3604832714.py:4: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

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

In [37]:

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

C:\Users\DELL\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\utils\validation.py:1143: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
y = column_or_1d(y, warn=True)
```

Out[37]:

LogisticRegression
LogisticRegression(max_iter=10000)

In [38]:

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

0.7160686427457098

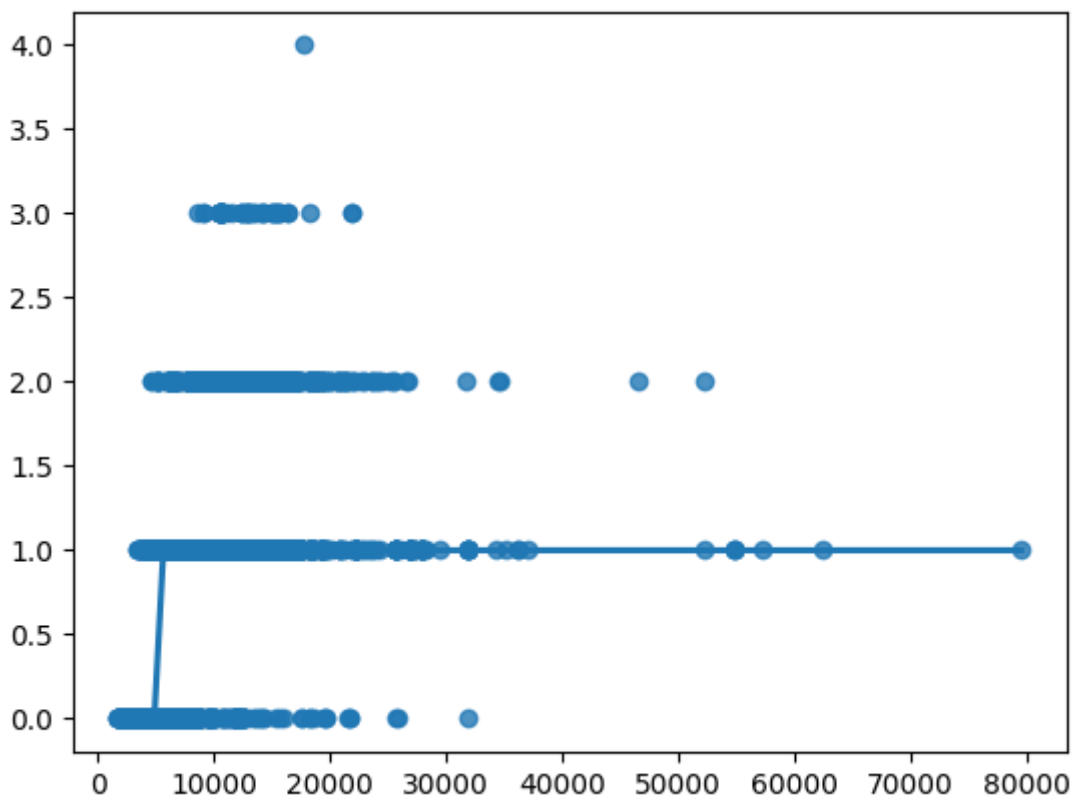
In [39]:

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

```
C:\Users\DELL\AppData\Local\Programs\Python\Python310\lib\site-packages\statsmodels\genmod\family\links.py:198: RuntimeWarning: overflow encountered in exp
  t = np.exp(-z)
```

Out[39]:

<Axes: >



Decision Tree

In [40]:

```
#Decision tree
from sklearn.tree import DecisionTreeClassifier
clf=DecisionTreeClassifier(random_state=0)
clf.fit(x_train,y_train)
```

Out[40]:

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

In [42]:

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

0.9369734789391576

Random Forest

In [43]:

```
#Random forest classifier
from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(X_train,y_train)
```

C:\Users\DELL\AppData\Local\Temp\ipykernel_9192\1232785509.py:4: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
rfc.fit(X_train,y_train)
```

Out[43]:

```
▼ RandomForestClassifier
RandomForestClassifier()
```

In [44]:

```
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 [45]:

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

In [46]:

```
grid_search.fit(X_train,y_train)
```

```
C:\Users\DELL\AppData\Local\Programs\Python\Python310\lib\site-packages
\sklearn\model_selection\_split.py:700: UserWarning: The least populated
class in y has only 1 members, which is less than n_splits=2.
  warnings.warn(
C:\Users\DELL\AppData\Local\Programs\Python\Python310\lib\site-packages
\sklearn\model_selection\_validation.py:686: DataConversionWarning: A co
lumn-vector y was passed when a 1d array was expected. Please change the
shape of y to (n_samples,), for example using ravel().
  estimator.fit(X_train, y_train, **fit_params)
C:\Users\DELL\AppData\Local\Programs\Python\Python310\lib\site-packages
\sklearn\model_selection\_validation.py:686: DataConversionWarning: A co
lumn-vector y was passed when a 1d array was expected. Please change the
shape of y to (n_samples,), for example using ravel().
  estimator.fit(X_train, y_train, **fit_params)
C:\Users\DELL\AppData\Local\Programs\Python\Python310\lib\site-packages
\sklearn\model_selection\_validation.py:686: DataConversionWarning: A co
lumn-vector y was passed when a 1d array was expected. Please change the
shape of y to (n_samples,), for example using ravel().
  estimator.fit(X_train, y_train, **fit_params)
```

In [47]:

```
grid_search.best_score_
```

Out[47]:

```
0.5244080692700013
```

In [48]:

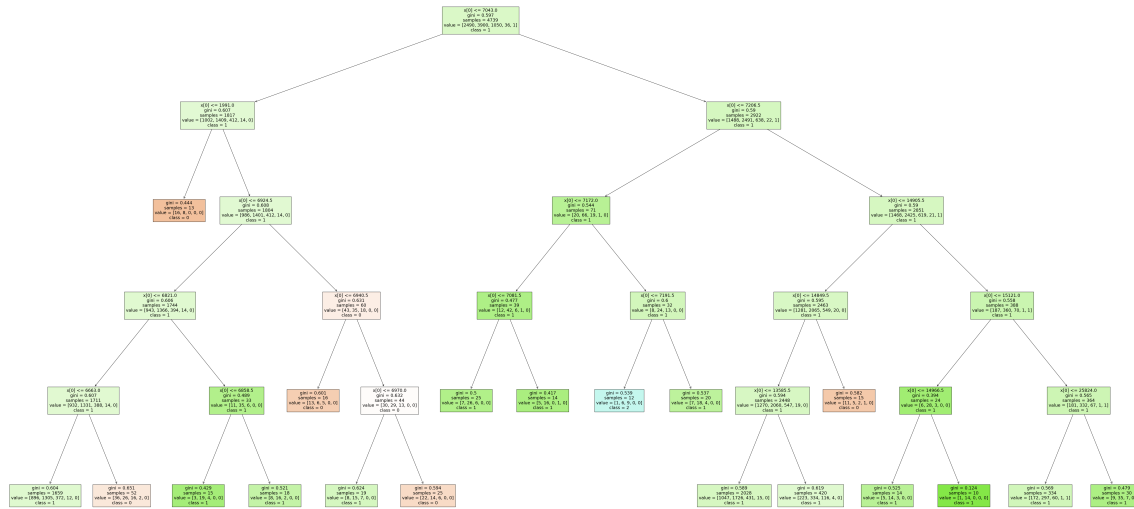
```
rf_best=grid_search.best_estimator_
rf_best
```

Out[48]:

```
RandomForestClassifier
RandomForestClassifier(max_depth=5, min_samples_leaf=10, n_estimators=10)
```

In [49]:

```
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);
```



In [52]:

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

0.45678627145085804

Conclusion:

By observing all the models Decision trees has the best accuracy compared to other models that is 93.

Therefore we can conclude that for flight price prediction Decision tree model is best fit.