In [1]:

```python
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 [3]:

```
train_df=pd.read_csv(r"C:\Users\monim\OneDrive\Desktop\Copy of Data_Train.csv")
train_df
```

Out[3]:

| | 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 [4]:

```
test_df=pd.read_csv(r"C:\Users\monim\OneDrive\Desktop\Copy of Test_set.csv")
test_df
```

Out[4]:

| | 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 [5]:

```
train_df.head()
```

Out[5]:

| | 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 [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]:

```
train_df.tail()
```

Out[7]:

| | 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 [8]:

```
test_df.tail()
```

Out[8]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duratio |
|---|---|---|---|---|---|---|---|---|
| **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 [9]:

```
train_df.describe()
```

Out[9]:

| | 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 [10]:

```python
test_df.describe()
```

Out[10]:

| | 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 [11]:

```python
train_df.shape
```

Out[11]:

```
(10683, 11)
```

In [12]:

```python
test_df.shape
```

Out[12]:

```
(2671, 10)
```

In [13]:

```python
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 [14]:

```python
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 [15]:

```python
train_df.isna().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 [16]:

```python
test_df.isna().sum()
```

Out[16]:

```
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
```
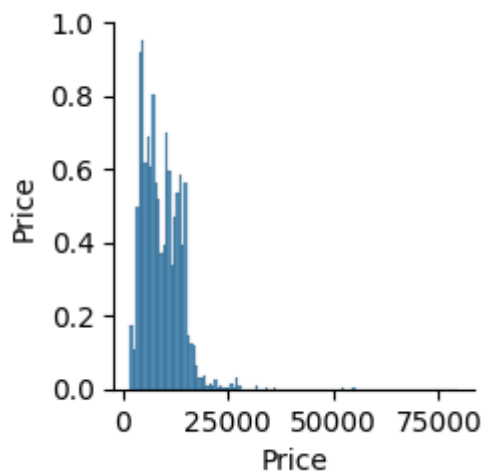
In [17]:

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

In [18]:

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

Out[18]:

```
<seaborn.axisgrid.PairGrid at 0x1aeab461e40>
```



In [20]:

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

Out[20]:

```
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 [22]:

```
A={"Airline":{"Jet Airways":0,"IndiGo":1,"Air India":2,"Multiple carriers":3,"SpiceJet":
 "Multiple carriers Premium economy":7,"Jet Airways Business":8,"Vistara Premium economy
train_df=train_df.replace(A)
train_df
```

Out[22]:

| | 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 | Air Asia | 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 [23]:

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

Out[23]:

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

In [23]:

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

In [24]:

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

Out[24]:

| | 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 | 1! |
| **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** | Air Asia | 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 | : |
| **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 [25]:

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

Out[25]:

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

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

In [27]:

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

Out[27]:

| | 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 | 1! |
| 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 | Air Asia | 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 | : |
| 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 [28]:

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

Out[28]:

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

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

In [29]:

```python
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[29]:

| | 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 5( |
| **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 | 1! |
| **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** | Air Asia | 9/04/2019 | 2 | 2 | CCU → BLR | 19:55 | 22:25 | 2h 3( |
| **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 | : |
| **10681** | 5 | 01/03/2019 | 3 | 4 | BLR → DEL | 11:30 | 14:10 | 2h 4( |
| **10682** | 2 | 9/05/2019 | 1 | 1 | DEL → GOI → BOM → COK | 10:55 | 19:15 | 8h 2( |

10682 rows × 11 columns

In [30]:

```python
sns.displot(train_df['Price'])
```
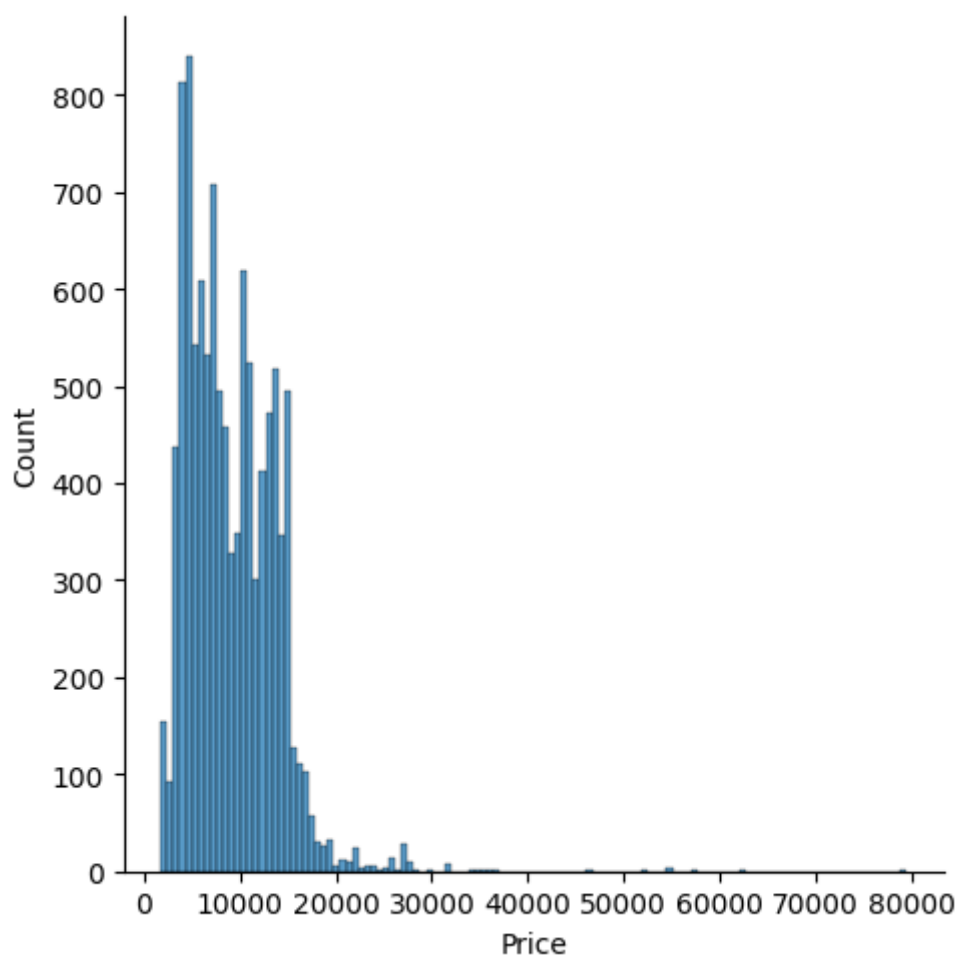
Out[30]:

```
<seaborn.axisgrid.FacetGrid at 0x1aead981330>
```



In [32]:

```python
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
```
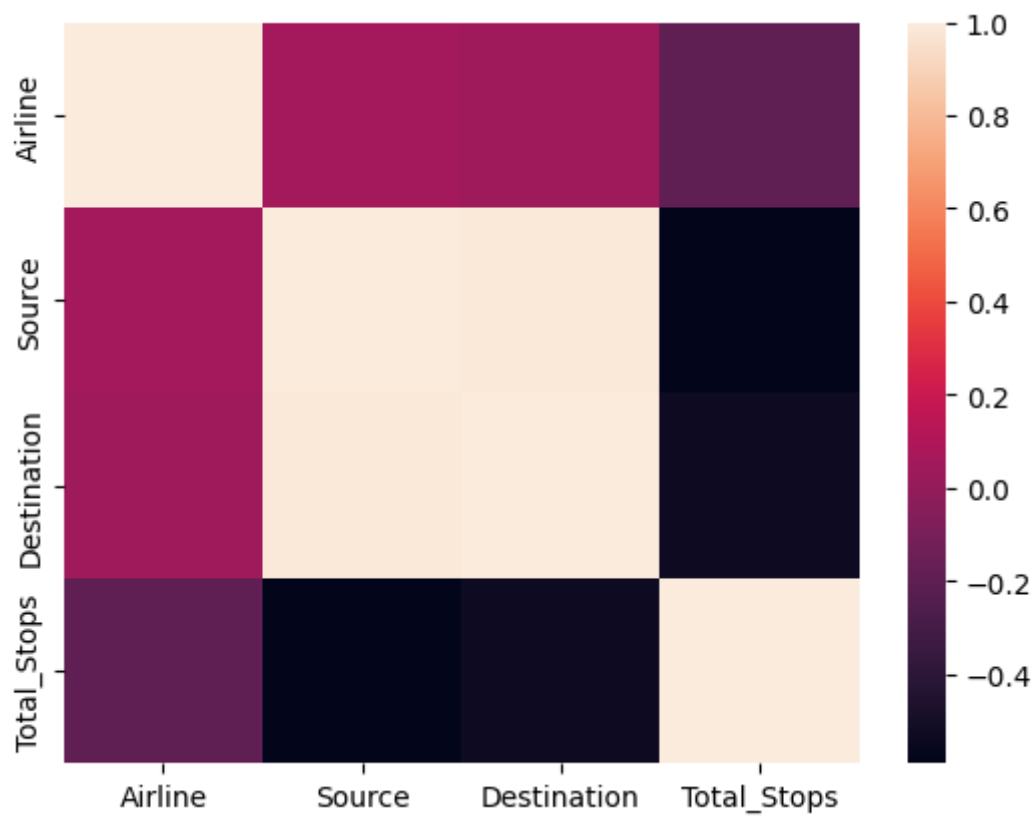
In [33]:

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

In [48]:

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

Out[48]:

```
<Axes: >
```



In [49]:

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

In [50]:

```python
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
```

7821.4711081011665

Out[50]:

|  | coefficient |
|---|---|
| Airline | -326.961757 |
| Source | -3204.781553 |
| Destination | 2451.058322 |
| Total_Stops | 3567.452821 |

In [51]:

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

0.41046053840891994

In [52]:

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

Out[52]:

```
<matplotlib.collections.PathCollection at 0x1aeafe693c0>
```



In [53]:

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

In [54]:

```python
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)
```
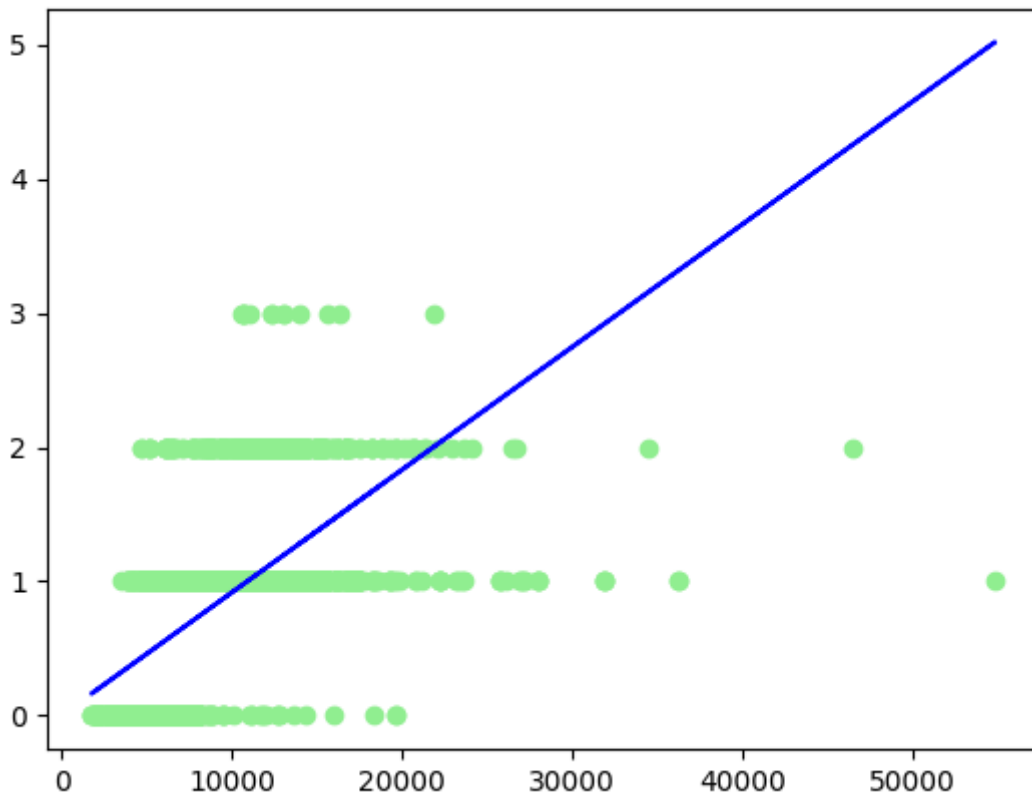
Out[54]:

```
LinearRegression()
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**
**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

In [56]:

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



# LOGISTIC REGRESSION

In [58]:

```python
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 [59]:

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

Out[59]:

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

In [60]:

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

0.431201248049922

# DECISION TREE

In [64]:

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

Out[64]:

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

In [65]:

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

0.921996879875195

# RANDOM FOREST

In [66]:

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

Out[66]:

```
▼ RandomForestClassifier
RandomForestClassifier()
```

In [67]:

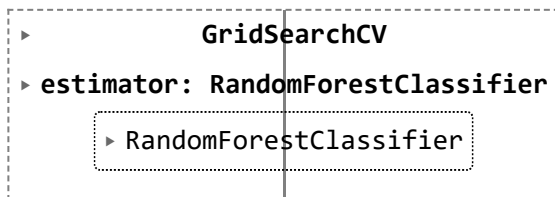```python
rf=RandomForestClassifier()
```

In [68]:

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

In [70]:

```python
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[70]:

```
┌ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┐
│  ▸         GridSearchCV             │
│                                     │
│ ▸ estimator: RandomForestClassifier │
│     ┌ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┐    │
│     ┆ ▸ RandomForestClassifier ┆    │
│     └ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┘    │
└ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┘
```

In [71]:

```python
grid_search.best_score_
```

Out[71]:

```
0.8175743550798769
```

In [72]:

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

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

In [*]:

```python
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 [*]:

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

# 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 [ ]: