# problem statement: which model is suitable(bestfit) for the given dataset

```
In [3]:  import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         from sklearn import preprocessing,svm
         from sklearn.model_selection import train_test_split
         from sklearn.linear_model import LinearRegression
```

# copy path and read the data frame ¶

In [4]:
```python
train_df=pd.read_csv(r"C:\Users\ubin1\OneDrive\Documents\jupyter\flight_train.
train_df
```

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 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 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 [5]: `test_df=pd.read_csv(r"C:\Users\ubin1\OneDrive\Documents\jupyter\flight_test.cs`
        `test_df`

Out[5]:

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

# Data cleaning and data preprocessing

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

Out[6]:

|   | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | To |
|---|---------|-----------------|--------|-------------|-------|----------|--------------|----------|----|
| 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]: `test_df.head()`

Out[7]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | To |
|---|---|---|---|---|---|---|---|---|---|
| 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 [8]: `train_df.tail()`

Out[8]:

| | 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 |
| 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 [9]: `test_df.tail()`

Out[9]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration |
|---|---|---|---|---|---|---|---|---|
| **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 |
| **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 [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 | Duration |
|---|---|---|---|---|---|---|---|---|
| **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 [12]: `train_df.shape`

Out[12]: `(10683, 11)`

In [13]: `test_df.shape`

Out[13]: `(2671, 10)`

In [14]: `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 [15]: `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 [16]: `train_df.columns`

Out[16]: 
```
Index(['Airline', 'Date_of_Journey', 'Source', 'Destination', 'Route',
       'Dep_Time', 'Arrival_Time', 'Duration', 'Total_Stops',
       'Additional_Info', 'Price'],
      dtype='object')
```

In [17]: `test_df.columns`

Out[17]: 
```
Index(['Airline', 'Date_of_Journey', 'Source', 'Destination', 'Route',
       'Dep_Time', 'Arrival_Time', 'Duration', 'Total_Stops',
       'Additional_Info'],
      dtype='object')
```

# Finding null values and replace them

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

In [21]: 
```python
train_df.isnull().sum()
```

Out[21]: 
```
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 [22]: 
```python
train_df.isnull().sum()
```

Out[22]: 
```
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 [23]: 
```python
train_df.shape
```

Out[23]: (10682, 11)

In [24]: `train_df['Airline'].value_counts()`

Out[24]:
```
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 [25]: `train_df['Source'].value_counts()`

Out[25]:
```
Source
Delhi      4536
Kolkata    2871
Banglore   2197
Mumbai      697
Chennai     381
Name: count, dtype: int64
```

In [26]: `train_df['Destination'].value_counts()`

Out[26]:
```
Destination
Cochin      4536
Banglore    2871
Delhi       1265
New Delhi    932
Hyderabad    697
Kolkata      381
Name: count, dtype: int64
```

In [27]: `train_df['Total_Stops'].value_counts()`

Out[27]:
```
Total_Stops
1 stop      5625
non-stop    3491
2 stops     1520
3 stops       45
4 stops        1
Name: count, dtype: int64
```

# changing strings into values

In [28]:
```python
airline={"Airline":{"Jet Airways":0,"IndiGo":1,"Air India":2,"Multiple carrier
 "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 | Duration |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 24/03/2019 | Banglore | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m |
| 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 |
| 10679 | 2 | 27/04/2019 | Kolkata | Banglore | CCU ? BLR | 20:45 | 23:20 | 2h 35m |
| 10680 | 0 | 27/04/2019 | Banglore | Delhi | BLR ? DEL | 08:20 | 11:20 | 3h |
| 10681 | 5 | 01/03/2019 | Banglore | New Delhi | BLR ? DEL | 11:30 | 14:10 | 2h 40m |
| 10682 | 2 | 9/05/2019 | Delhi | Cochin | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m |

10682 rows × 11 columns

In [29]:
```python
airline={"Airline":{"Jet Airways":0,"IndiGo":1,"Air India":2,"Multiple carrier
 "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[29]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 24/03/2019 | Banglore | New Delhi | BLR ? DEL | 22:20 | 01:10 22 Mar | 2h 50m |
| 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 |
| 10679 | 2 | 27/04/2019 | Kolkata | Banglore | CCU ? BLR | 20:45 | 23:20 | 2h 35m |
| 10680 | 0 | 27/04/2019 | Banglore | Delhi | BLR ? DEL | 08:20 | 11:20 | 3h |
| 10681 | 5 | 01/03/2019 | Banglore | New Delhi | BLR ? DEL | 11:30 | 14:10 | 2h 40m |
| 10682 | 2 | 9/05/2019 | Delhi | Cochin | DEL ? GOI ? BOM ? COK | 10:55 | 19:15 | 8h 20m |

10682 rows × 11 columns

In [30]:
```python
city={"Source":{"Delhi":0,"Kolkata":1,"Banglore":2,
 "Mumbai":3,"Chennai":4}}
train_df=train_df.replace(city)
train_df
```

Out[30]:

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

10682 rows × 11 columns

In [31]:
```python
destination={"Destination":{"Cochin":0,"Banglore":1,"Delhi":2,
 "New Delhi":3,"Hyderabad":4,"Kolkata":5}}
train_df=train_df.replace(destination)
train_df
```

Out[31]:

| | 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 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

In [32]:
```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[32]:

|  | 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 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

In [33]:  `train_df`

Out[33]:

| | 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 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 [34]: `#EDA`
`fdf=train_df[['Airline','Source','Destination','Total_Stops','Price']]`
`sns.heatmap(fdf.corr(),annot=True)`

Out[34]: `<Axes: >`



# Features scaling:To split the data into training data and test data

In [35]: `x=fdf[['Airline','Source','Destination','Total_Stops']]`
`y=fdf['Price']`

In [36]: `#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=`

In [37]: `#Linear Regression`

In [38]:
```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
```

7211.098088897486

Out[38]:

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

In [39]:
```python
#Linear Rgeression
score=regr.score(X_test,y_test)
print(score)
```

0.41083048909283504

In [40]:
```python
predictions=regr.predict(X_test)
```

In [41]:
```python
plt.scatter(y_test,predictions)
```

Out[41]: &lt;matplotlib.collections.PathCollection at 0x1e047428450&gt;



In [42]:
```python
x=np.array(fdf['Price']).reshape(-1,1)
y=np.array(fdf['Total_Stops']).reshape(-1,1)
fdf.dropna(inplace=True)
```

C:\Users\ubinl\AppData\Local\Temp\ipykernel_17928\521034954.py:3: SettingWith
CopyWarning:
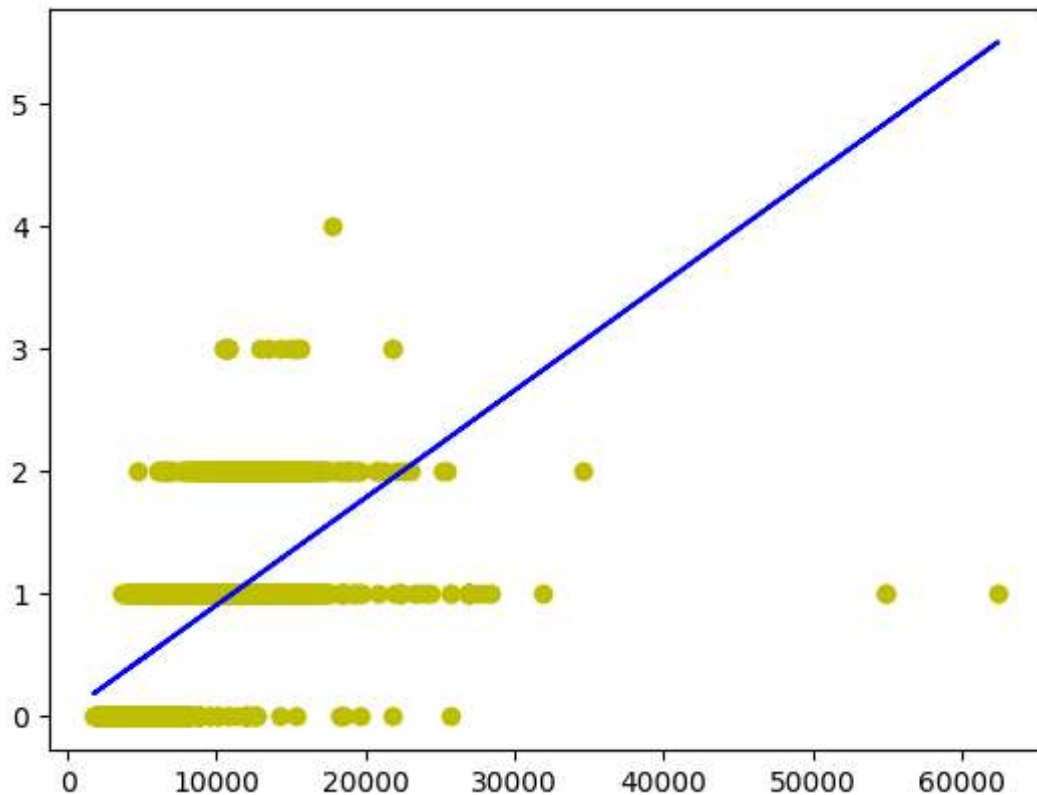A value is trying to be set on a copy of a slice from a DataFrame

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

In [43]:
```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_train,y_train)
```

Out[43]:
```
▼ LinearRegression
LinearRegression()
```

```
In [44]: y_pred=regr.predict(X_test)
         plt.scatter(X_test,y_test,color='y')
         plt.plot(X_test,y_pred,color='b')
         plt.show()
```



## since we disd not get the accuracy for linear Regression we are going to implement Logistic Regression

## Logistic Regression

In [51]: 
```python
#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=
from sklearn.linear_model import LogisticRegression
lr=LogisticRegression(max_iter=10000)
```

C:\Users\ubinl\AppData\Local\Temp\ipykernel_17928\3604832714.py:4: SettingWit
hCopyWarning:
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/s
table/user_guide/indexing.html#returning-a-view-versus-a-copy (https://panda
s.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-ver
sus-a-copy)
    fdf.dropna(inplace=True)

In [52]: 
```python
lr.fit(x_train,y_train)
```

C:\Users\ubinl\AppData\Local\Programs\Python\Python311\Lib\site-packages\skle
arn\utils\validation.py:1143: DataConversionWarning: A column-vector y was pa
ssed when a 1d array was expected. Please change the shape of y to (n_sample
s, ), for example using ravel().
    y = column_or_1d(y, warn=True)

Out[52]: 
```
  ▼          LogisticRegression

LogisticRegression(max_iter=10000)
```
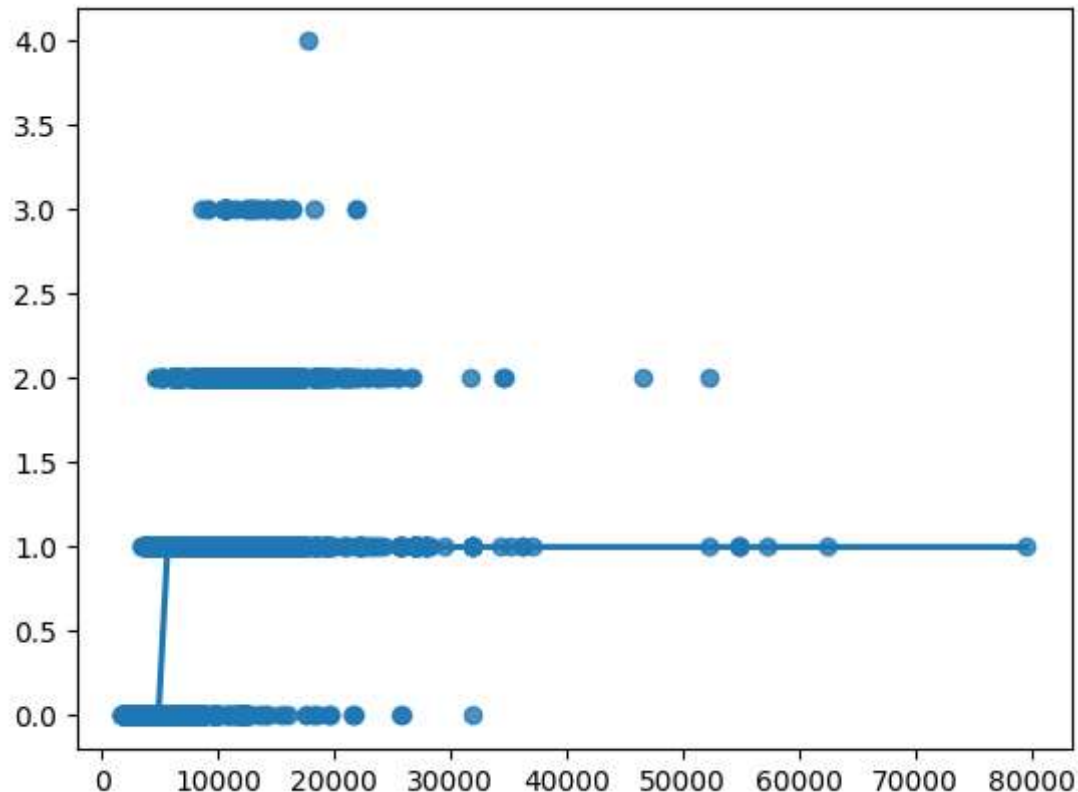
In [53]: 
```python
score=lr.score(x_test,y_test)
print(score)
```

0.7160686427457098

In [54]: 
```python
sns.regplot(x=x,y=y,data=fdf,logistic=True,ci=None)
```

```
C:\Users\ubinl\AppData\Local\Programs\Python\Python311\Lib\site-packages\stat
smodels\genmod\families\links.py:198: RuntimeWarning: overflow encountered in
exp
  t = np.exp(-z)
```

Out[54]: `<Axes: >`



# Since we did not get the accuracy for Logistic Regression we are going to implement Decision Tree and Random Forest and make a comparative study for finding the best model for the dataset

## Decision tree

In [55]: 
```python
from sklearn.tree import DecisionTreeClassifier
clf=DecisionTreeClassifier(random_state=0)
clf.fit(x_train,y_train)
```

Out[55]: 
```
▼          DecisionTreeClassifier

DecisionTreeClassifier(random_state=0)
```

In [56]:
```python
score=clf.score(x_test,y_test)
print(score)
```

0.9369734789391576

# Random Forest

In [49]:
```python
#Random forest classifier
from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
rfc.fit(X_train,y_train)
```

C:\Users\ubinl\AppData\Local\Temp\ipykernel_6292\2470359396.py:4: DataConvers
ionWarning: 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[49]:
```
▾ RandomForestClassifier

RandomForestClassifier()
```

In [50]:
```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]}
from sklearn.model_selection import GridSearchCV
grid_search=GridSearchCV(estimator=rfc,param_grid=params,cv=2,scoring="accurac
```

In [51]:
```python
grid_search.fit(X_train,y_train)
```

C:\Users\ubinl\AppData\Local\Programs\Python\Python311\Lib\site-packages\s
klearn\model_selection\_split.py:700: UserWarning: The least populated cla
ss in y has only 1 members, which is less than n_splits=2.
  warnings.warn(
C:\Users\ubinl\AppData\Local\Programs\Python\Python311\Lib\site-packages\s
klearn\model_selection\_validation.py:686: 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().
  estimator.fit(X_train, y_train, **fit_params)
C:\Users\ubinl\AppData\Local\Programs\Python\Python311\Lib\site-packages\s
klearn\model_selection\_validation.py:686: 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().
  estimator.fit(X_train, y_train, **fit_params)
C:\Users\ubinl\AppData\Local\Programs\Python\Python311\Lib\site-packages\s
klearn\model_selection\_validation.py:686: 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().
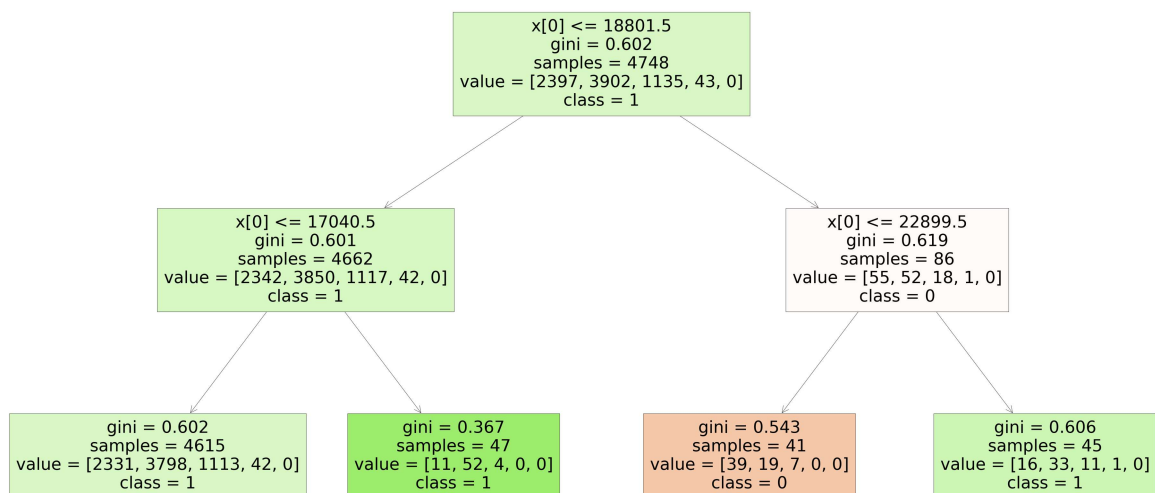  estimator.fit(X_train, y_train, **fit_params)

In [52]:
```python
grid_search.best_score_
```

Out[52]: 0.5244082839178266

In [53]:
```python
rf_best=grid_search.best_estimator_
rf_best
```

Out[53]:
```
                           RandomForestClassifier
RandomForestClassifier(max_depth=2, min_samples_leaf=20, n_estimators=25)
```

In [54]:
```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] <= 18801.5
                           gini = 0.602
                          samples = 4748
                  value = [2397, 3902, 1135, 43, 0]
                            class = 1

        x[0] <= 17040.5                          x[0] <= 22899.5
         gini = 0.601                             gini = 0.619
        samples = 4662                           samples = 86
value = [2342, 3850, 1117, 42, 0]        value = [55, 52, 18, 1, 0]
          class = 1                               class = 0

   gini = 0.602      gini = 0.367       gini = 0.543       gini = 0.606
 samples = 4615    samples = 47       samples = 41       samples = 45
value = [2331,    value = [11, 52,   value = [39, 19,   value = [16, 33,
3798, 1113, 42, 0]   4, 0, 0]          7, 0, 0]          11, 1, 0]
   class = 1          class = 1         class = 0         class = 1
```

In [55]:
```python
score=rfc.score(x_test,y_test)
print(score)
```

0.4836193447737909

# conclusion:Based on the above outcomes we can conclude that the best fit and accurate model is decision tree