

In [1]:

```
#FLIGHT PRICE PREDICTION
```

In []:

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

sns.set()
```

In [2]:

```
train_data = pd.read_csv("Data_Train.csv")
```

In [3]:

```
pd.set_option('display.max_columns', None)
```

In [4]:

```
train_data.head()
```

Out[4]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additio
0	IndiGo	24/03/2019	Banglore	New Delhi	BLR → DEL	22:20	01:10 22 Mar	2h 50m	non-stop	
1	Air India	1/05/2019	Kolkata	Banglore	CCU → IXR → BBI → BLR	05:50	13:15	7h 25m	2 stops	
2	Jet Airways	9/06/2019	Delhi	Cochin	DEL → LKO → BOM → COK	09:25	04:25 10 Jun	19h	2 stops	
3	IndiGo	12/05/2019	Kolkata	Banglore	CCU → NAG → BLR	18:05	23:30	5h 25m	1 stop	
4	IndiGo	01/03/2019	Banglore	New Delhi	BLR → NAG → DEL	16:50	21:35	4h 45m	1 stop	

In [5]:

```
train_data.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 [6]:

```
train_data["Duration"].value_counts()
```

Out[6]:

```
2h 50m      550
1h 30m      386
2h 45m      337
2h 55m      337
2h 35m      329
...
31h 30m       1
30h 25m       1
42h 5m        1
4h 10m        1
47h 40m       1
Name: Duration, Length: 368, dtype: int64
```

In [7]:

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

In [8]:

```
train_data.isnull().sum()
```

Out[8]:

```
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 [9]:

```
train_data["Journey_day"] = pd.to_datetime(train_data.Date_of_Journey, format="%d/%m/%Y").dt.day
```

In [10]:

```
train_data["Journey_month"] = pd.to_datetime(train_data["Date_of_Journey"], format = "%d/%m/%Y").dt.month
```

In [11]:

```
train_data.head()
```

Out[11]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additio
0	IndiGo	24/03/2019	Banglore	New Delhi	BLR → DEL	22:20	01:10 22 Mar	2h 50m	non-stop	
1	Air India	1/05/2019	Kolkata	Banglore	CCU → IXR → BBI → BLR	05:50	13:15	7h 25m	2 stops	
2	Jet Airways	9/06/2019	Delhi	Cochin	DEL → LKO → BOM → COK	09:25	04:25 10 Jun	19h	2 stops	
3	IndiGo	12/05/2019	Kolkata	Banglore	CCU → NAG → BLR	18:05	23:30	5h 25m	1 stop	
4	IndiGo	01/03/2019	Banglore	New Delhi	BLR → NAG → DEL	16:50	21:35	4h 45m	1 stop	

In [12]:

```
train_data.drop(["Date_of_Journey"], axis = 1, inplace = True)
```

In [13]:

```
train_data["Dep_hour"] = pd.to_datetime(train_data["Dep_Time"]).dt.hour
train_data["Dep_min"] = pd.to_datetime(train_data["Dep_Time"]).dt.minute
train_data.drop(["Dep_Time"], axis = 1, inplace = True)
```

In [14]:

```
train_data.head()
```

Out[14]:

	Airline	Source	Destination	Route	Arrival_Time	Duration	Total_Stops	Additional_Info	Price	Journey_day
0	IndiGo	Banglore	New Delhi	BLR → DEL	01:10 22 Mar	2h 50m	non-stop	No info	3897	24
1	Air India	Kolkata	Banglore	CCU → IXR → BBI → BLR	13:15	7h 25m	2 stops	No info	7662	1
2	Jet Airways	Delhi	Cochin	DEL → LKO → BOM → COK	04:25 10 Jun	19h	2 stops	No info	13882	9
3	IndiGo	Kolkata	Banglore	CCU → NAG → BLR	23:30	5h 25m	1 stop	No info	6218	12
4	IndiGo	Banglore	New Delhi	BLR → NAG → DEL	21:35	4h 45m	1 stop	No info	13302	1

In [15]:

```
train_data["Arrival_hour"] = pd.to_datetime(train_data.Arrival_Time).dt.hour
train_data["Arrival_min"] = pd.to_datetime(train_data.Arrival_Time).dt.minute
train_data.drop(["Arrival_Time"], axis = 1, inplace = True)
```

In [16]:

```
train_data.head()
```

Out[16]:

	Airline	Source	Destination	Route	Duration	Total_Stops	Additional_Info	Price	Journey_day	Journey_mor
0	IndiGo	Banglore	New Delhi	BLR → DEL	2h 50m	non-stop	No info	3897	24	
1	Air India	Kolkata	Banglore	CCU → IXR → BBI → BLR	7h 25m	2 stops	No info	7662	1	
2	Jet Airways	Delhi	Cochin	DEL → LKO → BOM → COK	19h	2 stops	No info	13882	9	
3	IndiGo	Kolkata	Banglore	CCU → NAG → BLR	5h 25m	1 stop	No info	6218	12	
4	IndiGo	Banglore	New Delhi	BLR → NAG → DEL	4h 45m	1 stop	No info	13302	1	

In [17]:

```
duration = list(train_data["Duration"])

for i in range(len(duration)):
    if len(duration[i].split()) != 2:
        if "h" in duration[i]:
            duration[i] = duration[i].strip() + " 0m"
        else:
            duration[i] = "0h " + duration[i]

duration_hours = []
duration_mins = []
for i in range(len(duration)):
    duration_hours.append(int(duration[i].split(sep = "h")[0]))
    duration_mins.append(int(duration[i].split(sep = "m")[0].split()[-1]))
```

In [18]:

```
train_data["Duration_hours"] = duration_hours
train_data["Duration_mins"] = duration_mins
```

In [19]:

```
train_data.drop(["Duration"], axis = 1, inplace = True)
```

In [20]:

```
train_data.head()
```

Out[20]:

	Airline	Source	Destination	Route	Total_Stops	Additional_Info	Price	Journey_day	Journey_month	Dep_t
0	IndiGo	Banglore	New Delhi	BLR → DEL	non-stop	No info	3897	24	3	
1	Air India	Kolkata	Banglore	CCU → IXR → BBI → BLR	2 stops	No info	7662	1	5	
2	Jet Airways	Delhi	Cochin	DEL → LKO → BOM → COK	2 stops	No info	13882	9	6	
3	IndiGo	Kolkata	Banglore	CCU → NAG → BLR	1 stop	No info	6218	12	5	
4	IndiGo	Banglore	New Delhi	BLR → NAG → DEL	1 stop	No info	13302	1	3	

In [21]:

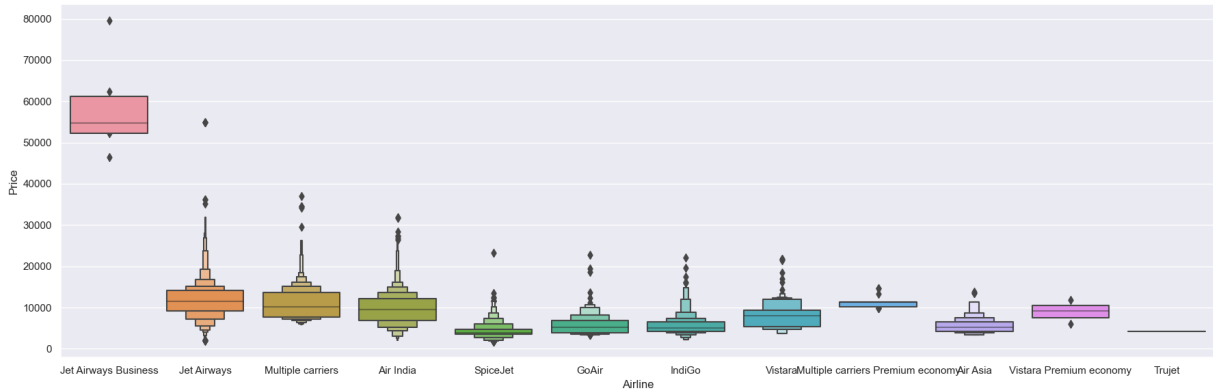
```
train_data["Airline"].value_counts()
```

Out[21]:

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 [22]:

```
sns.catplot(y = "Price", x = "Airline", data = train_data.sort_values("Price", ascending = False), kind="box",
plt.show()
```



In [23]:

```
Airline = train_data[["Airline"]]
Airline = pd.get_dummies(Airline, drop_first= True)
Airline.head()
```

Out[23]:

	Airline_Air India	Airline_GoAir	Airline_IndiGo	Airline_Jet Airways	Airline_Jet Airways Business	Airline_Multiple carriers	Airline_Multiple carriers Premium economy	Airline_SpiceJet
0	0	0	1	0	0	0	0	0
1	1	0	0	0	0	0	0	0
2	0	0	0	1	0	0	0	0
3	0	0	1	0	0	0	0	0
4	0	0	1	0	0	0	0	0

In [24]:

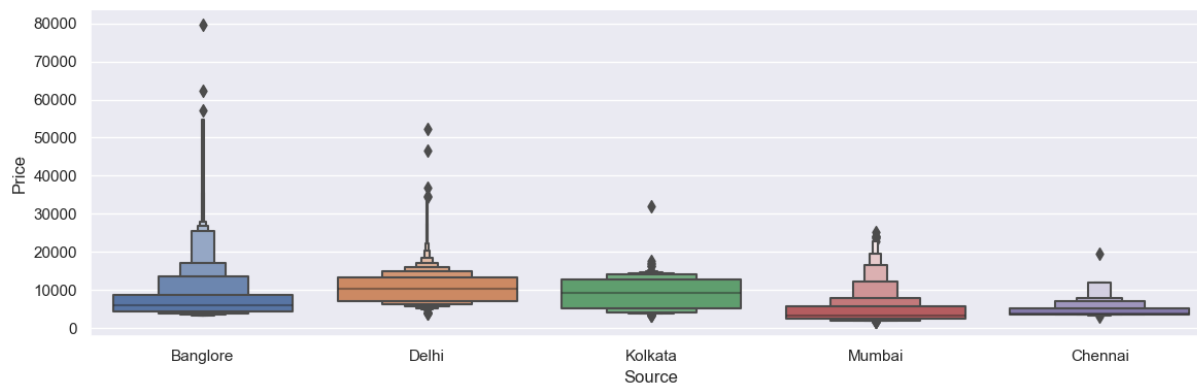
```
train_data["Source"].value_counts()
```

Out[24]:

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

In [25]:

```
sns.catplot(y = "Price", x = "Source", data = train_data.sort_values("Price", ascending = False), kind="box",
plt.show())
```



In [26]:

```
Source = train_data[["Source"]]
Source = pd.get_dummies(Source, drop_first= True)
Source.head()
```

Out[26]:

	Source_Chennai	Source_Delhi	Source_Kolkata	Source_Mumbai
0	0	0	0	0
1	0	0	1	0
2	0	1	0	0
3	0	0	1	0
4	0	0	0	0

In [27]:

```
train_data["Destination"].value_counts()
```

Out[27]:

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


In [28]:

```
Destination = train_data[["Destination"]]

Destination = pd.get_dummies(Destination, drop_first = True)

Destination.head()
```

Out[28]:

	Destination_Cochin	Destination_Delhi	Destination_Hyderabad	Destination_Kolkata	Destination_New Delhi
0	0	0	0	0	1
1	0	0	0	0	0
2	1	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	1

In [29]:

```
train_data["Route"]
```

Out[29]:

```
0          BLR → DEL
1    CCU → IXR → BBI → BLR
2    DEL → LKO → BOM → COK
3          CCU → NAG → BLR
4          BLR → NAG → DEL
...
10678          CCU → BLR
10679          CCU → BLR
10680          BLR → DEL
10681          BLR → DEL
10682    DEL → GOI → BOM → COK
Name: Route, Length: 10682, dtype: object
```

In [30]:

```
train_data.drop(["Route", "Additional_Info"], axis = 1, inplace = True)
```

In [31]:

```
train_data["Total_Stops"].value_counts()
```

Out[31]:

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

In [32]:

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

In [33]:

```
train_data.head()
```

Out[33]:

	Airline	Source	Destination	Total_Stops	Price	Journey_day	Journey_month	Dep_hour	Dep_min	Arrival_h
0	IndiGo	Banglore	New Delhi	0	3897	24	3	22	20	
1	Air India	Kolkata	Banglore	2	7662	1	5	5	50	
2	Jet Airways	Delhi	Cochin	2	13882	9	6	9	25	
3	IndiGo	Kolkata	Banglore	1	6218	12	5	18	5	
4	IndiGo	Banglore	New Delhi	1	13302	1	3	16	50	

In [34]:

```
data_train = pd.concat([train_data, Airline, Source, Destination], axis = 1)
```

In [35]:

```
data_train.head()
```

Out[35]:

	Source	Destination	Total_Stops	Price	Journey_day	Journey_month	Dep_hour	Dep_min	Arrival_hour	Arrival_min
0	Banglore	New Delhi	0	3897	24	3	22	20	1	
1	Kolkata	Banglore	2	7662	1	5	5	50	13	
2	Delhi	Cochin	2	13882	9	6	9	25	4	
3	Kolkata	Banglore	1	6218	12	5	18	5	23	
4	Banglore	New Delhi	1	13302	1	3	16	50	21	

In [36]:

```
data_train.drop(["Airline", "Source", "Destination"], axis = 1, inplace = True)
```

In [37]:

```
data_train.head()
```

Out[37]:

	Total_Stops	Price	Journey_day	Journey_month	Dep_hour	Dep_min	Arrival_hour	Arrival_min	Duration_hou
0	0	3897	24	3	22	20	1	10	
1	2	7662	1	5	5	50	13	15	
2	2	13882	9	6	9	25	4	25	
3	1	6218	12	5	18	5	23	30	
4	1	13302	1	3	16	50	21	35	

In [38]:

```
data_train.shape
```

Out[38]:

(10682, 30)

In []:

```
TEST SET
```

In [39]:

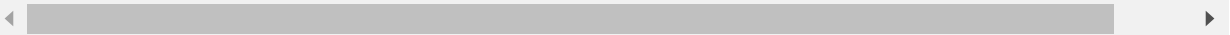
```
test_data = pd.read_csv("Test_set.csv")
```

In [40]:

```
test_data.head()
```

Out[40]:

	Airline	Date_of_Journey	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Additional
0	Jet Airways	6/06/2019	Delhi	Cochin	DEL → BOM → COK	17:30	04:25 07 Jun	10h 55m	1 stop	
1	IndiGo	12/05/2019	Kolkata	Banglore	CCU → MAA → BLR	06:20	10:20	4h	1 stop	
2	Jet Airways	21/05/2019	Delhi	Cochin	DEL → BOM → COK	19:15	19:00 22 May	23h 45m	1 stop	In-flight not
3	Multiple carriers	21/05/2019	Delhi	Cochin	DEL → BOM → COK	08:00	21:00	13h	1 stop	
4	Air Asia	24/06/2019	Banglore	Delhi	BLR → DEL	23:55	02:45 25 Jun	2h 50m	non-stop	



In [41]:

```

print("Test data Info")
print("-"*75)
print(test_data.info())

print()
print()

print("Null values :")
print("-"*75)
test_data.dropna(inplace = True)
print(test_data.isnull().sum())

test_data["Journey_day"] = pd.to_datetime(test_data.Date_of_Journey, format="%d/%m/%Y").dt.day
test_data["Journey_month"] = pd.to_datetime(test_data["Date_of_Journey"], format = "%d/%m/%Y").dt.month
test_data.drop(["Date_of_Journey"], axis = 1, inplace = True)

# Dep_Time
test_data["Dep_hour"] = pd.to_datetime(test_data["Dep_Time"]).dt.hour
test_data["Dep_min"] = pd.to_datetime(test_data["Dep_Time"]).dt.minute
test_data.drop(["Dep_Time"], axis = 1, inplace = True)

test_data["Arrival_hour"] = pd.to_datetime(test_data.Arrival_Time).dt.hour
test_data["Arrival_min"] = pd.to_datetime(test_data.Arrival_Time).dt.minute
test_data.drop(["Arrival_Time"], axis = 1, inplace = True)

duration = list(test_data["Duration"])

for i in range(len(duration)):
    if len(duration[i].split()) != 2:    # Check if duration contains only hour or mins
        if "h" in duration[i]:
            duration[i] = duration[i].strip() + " 0m"    # Adds 0 minute
        else:
            duration[i] = "0h " + duration[i]            # Adds 0 hour

duration_hours = []
duration_mins = []
for i in range(len(duration)):
    duration_hours.append(int(duration[i].split(sep = "h")[0]))    # Extract hours from duration
    duration_mins.append(int(duration[i].split(sep = "m")[0].split()[-1]))    # Extracts only minutes from

# Adding Duration column to test set
test_data["Duration_hours"] = duration_hours
test_data["Duration_mins"] = duration_mins
test_data.drop(["Duration"], axis = 1, inplace = True)

print("Airline")
print("-"*75)
print(test_data["Airline"].value_counts())
Airline = pd.get_dummies(test_data["Airline"], drop_first= True)

print()

print("Source")
print("-"*75)
print(test_data["Source"].value_counts())
Source = pd.get_dummies(test_data["Source"], drop_first= True)

print()

print("Destination")
print("-"*75)
print(test_data["Destination"].value_counts())
Destination = pd.get_dummies(test_data["Destination"], drop_first = True)

test_data.drop(["Route", "Additional_Info"], axis = 1, inplace = True)

test_data.replace({"non-stop": 0, "1 stop": 1, "2 stops": 2, "3 stops": 3, "4 stops": 4}, inplace = True)

```

```
data_test = pd.concat([test_data, Airline, Source, Destination], axis = 1)
data_test.drop(["Airline", "Source", "Destination"], axis = 1, inplace = True)

print()
print()

print("Shape of test data : ", data_test.shape)
```

Test data 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
None

```

Null values :

```

-----
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
Airline

```

```

-----
Jet Airways            897
IndiGo                 511
Air India              440
Multiple carriers      347
SpiceJet               208
Vistara                129
Air Asia               86
GoAir                  46
Multiple carriers Premium economy  3
Vistara Premium economy  2
Jet Airways Business   2
Name: Airline, dtype: int64

```

Source

```

-----
Delhi                1145
Kolkata              710
Bangalore            555
Mumbai               186
Chennai              75
Name: Source, dtype: int64

```

Destination

```

-----
Cochin                1145
Bangalore              710
Delhi                  317
New Delhi             238
Hyderabad             186
Kolkata                75
Name: Destination, dtype: int64

```

Shape of test data : (2671, 28)

In [42]:

```
data_test.head()
```

Out[42]:

	Total_Stops	Journey_day	Journey_month	Dep_hour	Dep_min	Arrival_hour	Arrival_min	Duration_hours	Dur:
0	1	6	6	17	30	4	25	10	
1	1	12	5	6	20	10	20	4	
2	1	21	5	19	15	19	0	23	
3	1	21	5	8	0	21	0	13	
4	0	24	6	23	55	2	45	2	

In []:

```
#FEATURE SELECTION
```

In [43]:

```
data_train.shape
```

Out[43]:

(10682, 30)

In [44]:

```
data_train.columns
```

Out[44]:

```
Index(['Total_Stops', 'Price', 'Journey_day', 'Journey_month', 'Dep_hour',
      'Dep_min', 'Arrival_hour', 'Arrival_min', 'Duration_hours',
      'Duration_mins', 'Airline_Air India', 'Airline_GoAir', 'Airline_IndiGo',
      'Airline_Jet Airways', 'Airline_Jet Airways Business',
      'Airline_Multiple carriers',
      'Airline_Multiple carriers Premium economy', 'Airline_SpiceJet',
      'Airline_Trujet', 'Airline_Vistara', 'Airline_Vistara Premium economy',
      'Source_Chennai', 'Source_Delhi', 'Source_Kolkata', 'Source_Mumbai',
      'Destination_Cochin', 'Destination_Delhi', 'Destination_Hyderabad',
      'Destination_Kolkata', 'Destination_New Delhi'],
      dtype='object')
```


In [45]:

```
X = data_train.loc[:, ['Total_Stops', 'Journey_day', 'Journey_month', 'Dep_hour',
    'Dep_min', 'Arrival_hour', 'Arrival_min', 'Duration_hours',
    'Duration_mins', 'Airline_Air India', 'Airline_GoAir', 'Airline_IndiGo',
    'Airline_Jet Airways', 'Airline_Jet Airways Business',
    'Airline_Multiple carriers',
    'Airline_Multiple carriers Premium economy', 'Airline_SpiceJet',
    'Airline_Trujet', 'Airline_Vistara', 'Airline_Vistara Premium economy',
    'Source_Chennai', 'Source_Delhi', 'Source_Kolkata', 'Source_Mumbai',
    'Destination_Cochin', 'Destination_Delhi', 'Destination_Hyderabad',
    'Destination_Kolkata', 'Destination_New Delhi']]
X.head()
```

Out[45]:

	Journey_day	Journey_month	Dep_hour	Dep_min	Arrival_hour	Arrival_min	Duration_hours	Duration_mins	Ai
0	24	3	22	20	1	10	2	50	
2	1	5	5	50	13	15	7	25	
2	9	6	9	25	4	25	19	0	
1	12	5	18	5	23	30	5	25	
1	1	3	16	50	21	35	4	45	

In [46]:

```
y = data_train.iloc[:, 1]
y.head()
```

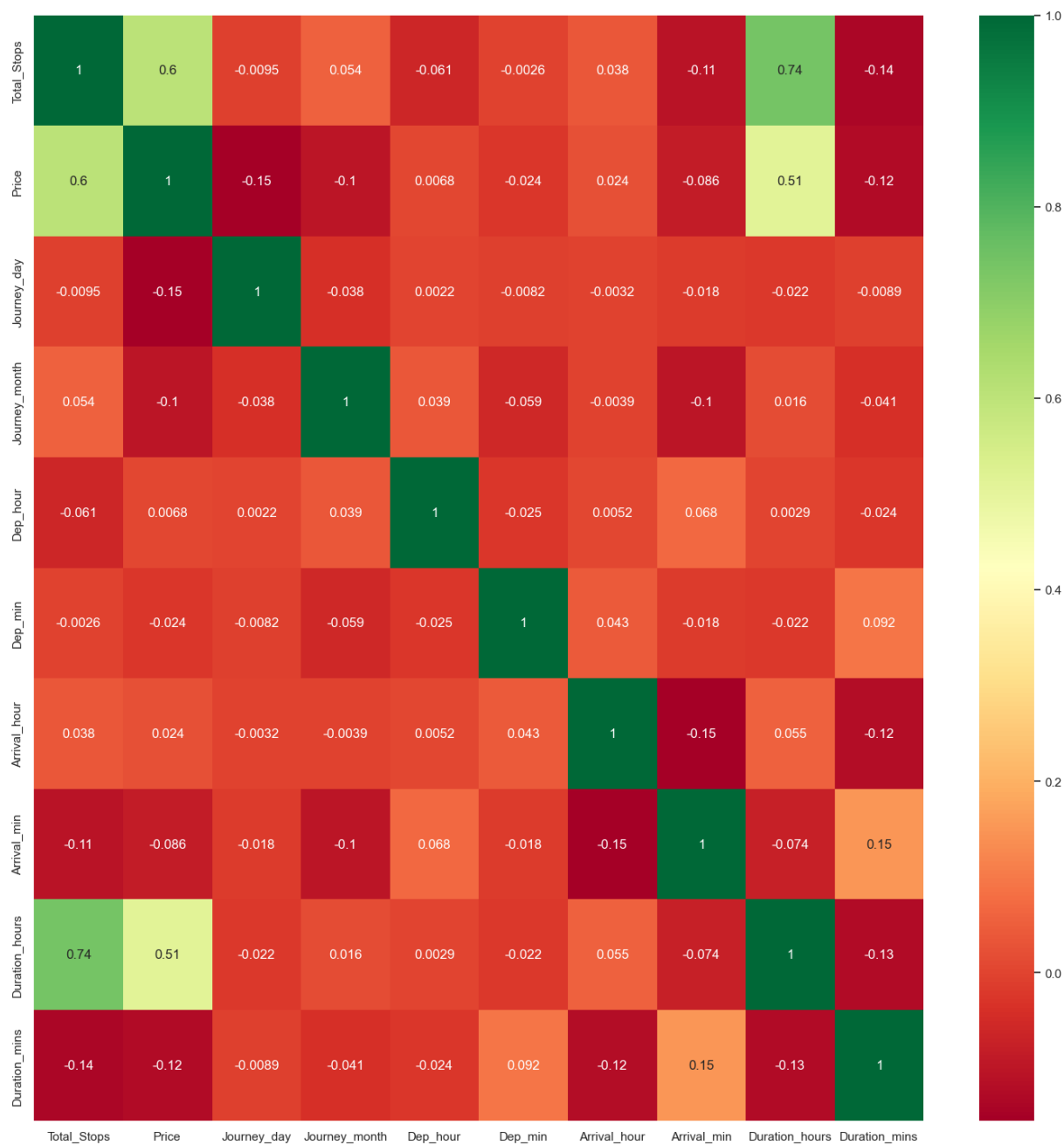
Out[46]:

```
0    3897
1    7662
2   13882
3    6218
4   13302
Name: Price, dtype: int64
```

In [47]:

```
plt.figure(figsize = (18,18))
sns.heatmap(train_data.corr(), annot = True, cmap = "RdYlGn")

plt.show()
```



In [48]:

```
from sklearn.ensemble import ExtraTreesRegressor
selection = ExtraTreesRegressor()
selection.fit(X, y)
```

Out[48]:

ExtraTreesRegressor()

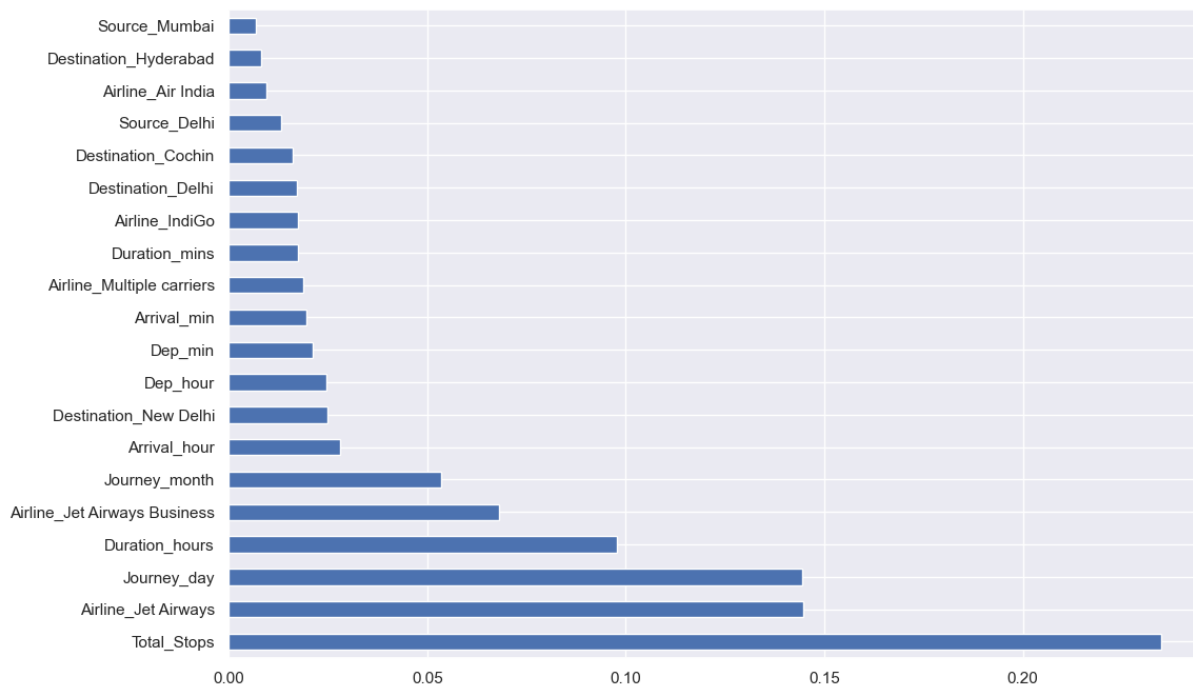
In [49]:

```
print(selection.feature_importances_)
```

```
[2.34765605e-01 1.44308213e-01 5.34872060e-02 2.45388780e-02
 2.12346579e-02 2.79562653e-02 1.94912228e-02 9.77197432e-02
 1.75050671e-02 9.42115992e-03 1.74455454e-03 1.74682473e-02
 1.44687039e-01 6.79946043e-02 1.86904256e-02 8.60664305e-04
 3.03588528e-03 1.25347115e-04 5.04457440e-03 8.41503687e-05
 4.03888347e-04 1.31164257e-02 3.01961064e-03 6.77256821e-03
 1.59923839e-02 1.71426706e-02 8.03409644e-03 5.13522659e-04
 2.48413234e-02]
```

In [50]:

```
plt.figure(figsize = (12,8))
feat_importances = pd.Series(selection.feature_importances_, index=X.columns)
feat_importances.nlargest(20).plot(kind='barh')
plt.show()
```



In []:

```
#FITTING MODEL USING RANDOM FOREST
```

In [51]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
```

In [52]:

```
from sklearn.ensemble import RandomForestRegressor
reg_rf = RandomForestRegressor()
reg_rf.fit(X_train, y_train)
```

Out[52]:

```
RandomForestRegressor()
```

In [53]:

```
y_pred = reg_rf.predict(X_test)
```

In [54]:

```
reg_rf.score(X_train, y_train)
```

Out[54]:

0.9533368127862233

In [55]:

```
reg_rf.score(X_test, y_test)
```

Out[55]:

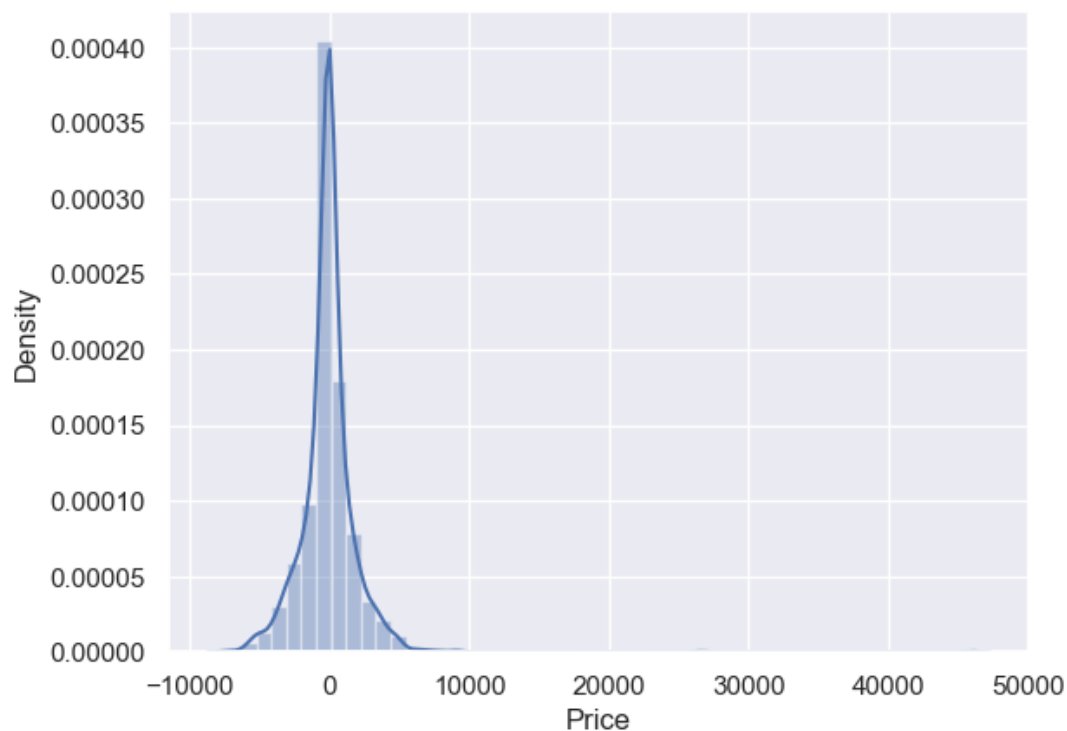
0.7976494001965726

In [56]:

```
sns.distplot(y_test-y_pred)  
plt.show()
```

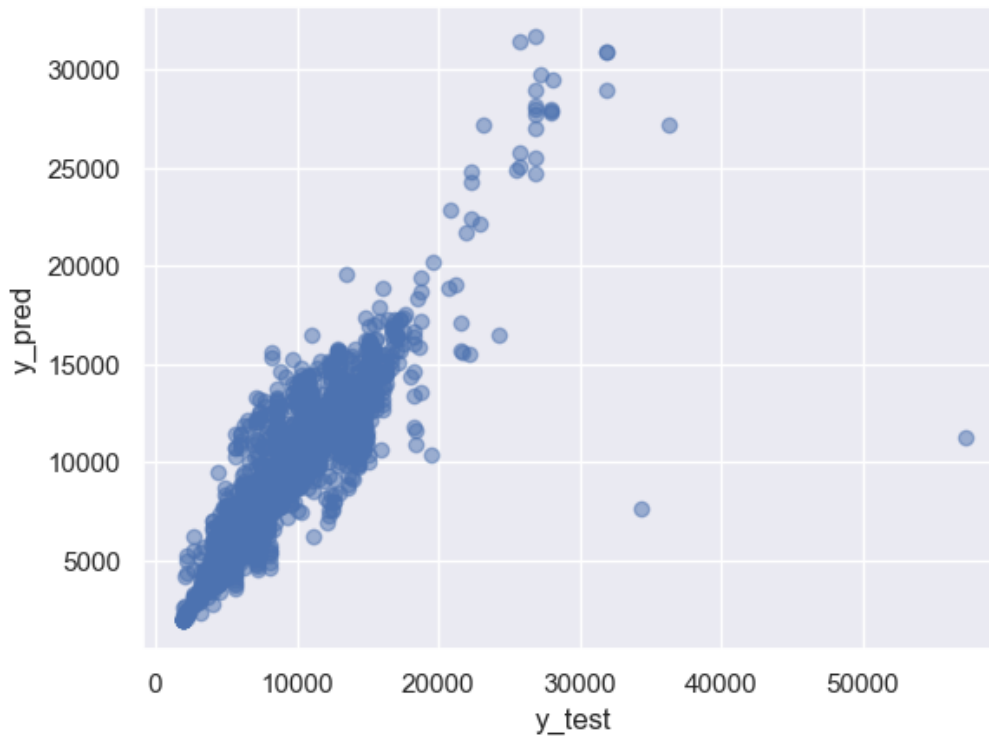
C:\Users\shaik\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)



In [57]:

```
plt.scatter(y_test, y_pred, alpha = 0.5)
plt.xlabel("y_test")
plt.ylabel("y_pred")
plt.show()
```



In [58]:

```
from sklearn import metrics
```

In [59]:

```
print('MAE:', metrics.mean_absolute_error(y_test, y_pred))
print('MSE:', metrics.mean_squared_error(y_test, y_pred))
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)))
```

MAE: 1174.9125929131571

MSE: 4363095.030189421

RMSE: 2088.802295620488

In [60]:

```
2090.5509/(max(y)-min(y))
```

Out[60]:

0.026887077025966846

In [61]:

```
metrics.r2_score(y_test, y_pred)
```

Out[61]:

0.7976494001965726

In []:

```
#Hyperparameter Tuning
```

In [62]:

```
from sklearn.model_selection import RandomizedSearchCV
```

In [63]:

```
n_estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]
max_features = ['auto', 'sqrt']
max_depth = [int(x) for x in np.linspace(5, 30, num = 6)]
min_samples_split = [2, 5, 10, 15, 100]
min_samples_leaf = [1, 2, 5, 10]
```

In [64]:

```
random_grid = {'n_estimators': n_estimators,
               'max_features': max_features,
               'max_depth': max_depth,
               'min_samples_split': min_samples_split,
               'min_samples_leaf': min_samples_leaf}
```

In [65]:

```
rf_random = RandomizedSearchCV(estimator = reg_rf, param_distributions = random_grid, scoring='neg_mean_sq
```

In [66]:

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

24/27


```

ators=1100; total time= 5.8s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total time= 2.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total time= 2.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total time= 2.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total time= 2.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total time= 2.9s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time= 3.7s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time= 3.7s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time= 3.6s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time= 3.7s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time= 3.7s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total time= 24.2s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total time= 24.0s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total time= 23.6s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total time= 23.8s
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total time= 23.9s

```

Out[66]:

```

RandomizedSearchCV(cv=5, estimator=RandomForestRegressor(), n_jobs=1,
                  param_distributions={'max_depth': [5, 10, 15, 20, 25, 30],
                                      'max_features': ['auto', 'sqrt'],
                                      'min_samples_leaf': [1, 2, 5, 10],
                                      'min_samples_split': [2, 5, 10, 15, 100],
                                      'n_estimators': [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 1200]},
                  random_state=42, scoring='neg_mean_squared_error',
                  verbose=2)

```

In [67]:

```
rf_random.best_params_
```

Out[67]:

```

{'n_estimators': 700,
 'min_samples_split': 15,
 'min_samples_leaf': 1,
 'max_features': 'auto',
 'max_depth': 20}

```

In [68]:

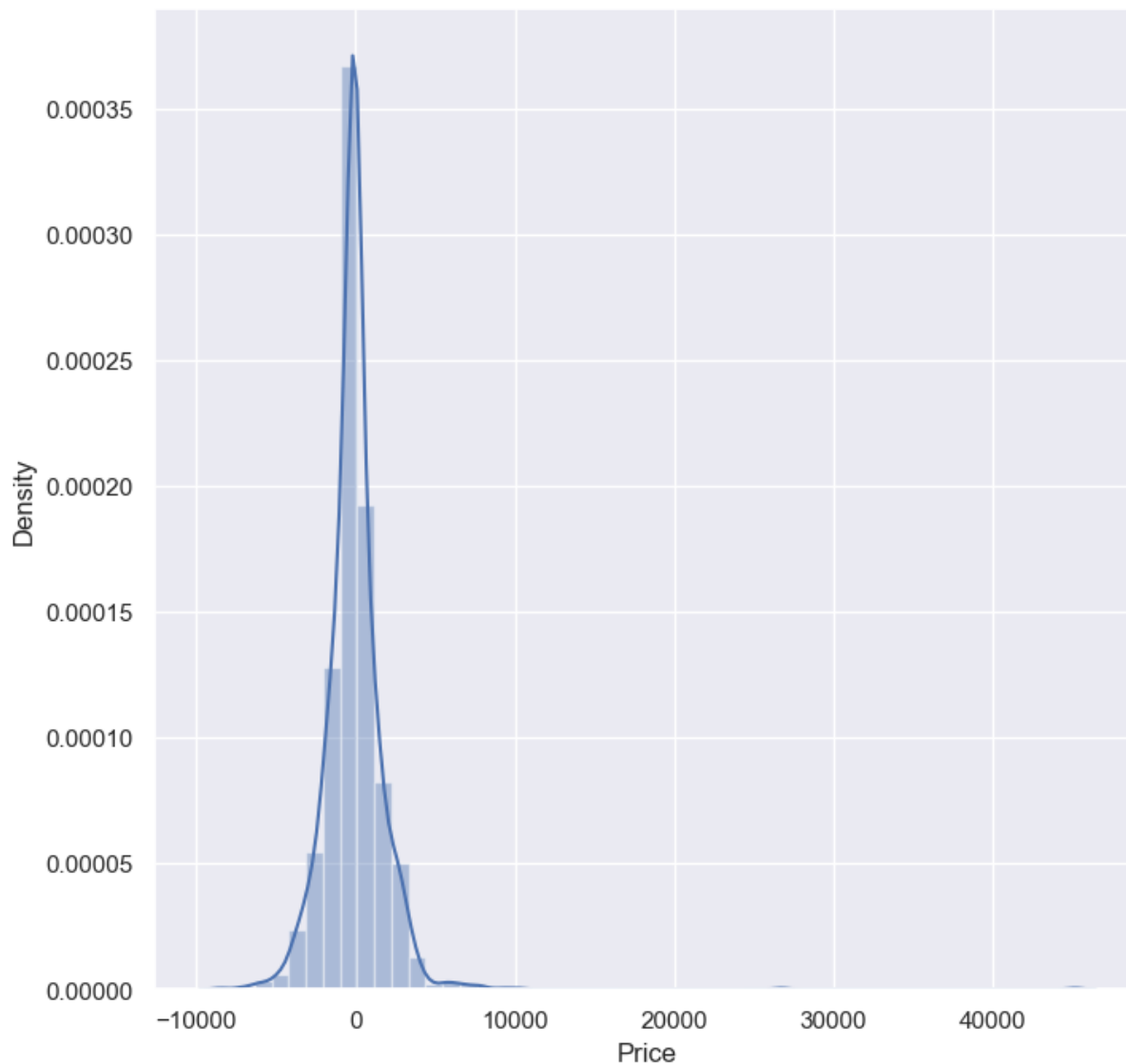
```
prediction = rf_random.predict(X_test)
```

In [69]:

```
plt.figure(figsize = (8,8))  
sns.distplot(y_test-prediction)  
plt.show()
```

C:\Users\shaik\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)



In [70]:

```
plt.figure(figsize = (8,8))
plt.scatter(y_test, prediction, alpha = 0.5)
plt.xlabel("y_test")
plt.ylabel("y_pred")
plt.show()
```



In [71]:

```
print('MAE:', metrics.mean_absolute_error(y_test, prediction))
print('MSE:', metrics.mean_squared_error(y_test, prediction))
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, prediction)))
```

MAE: 1165.6821239351436
MSE: 4052649.1265631523
RMSE: 2013.1192529413534

In []: