```
In [1]:  import numpy as np
         import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns

         sns.set()
```

# Importing dataset

Since data is in form of excel file we have to use pandas read_excel to load the data After loading it is important to check the complete information of data as it can indication many of the hidden infomation such as null values in a column or a row Check whether any null values are there or not. if it is present then following can be done, Imputing data using Imputation method in sklearn Filling NaN values with mean, median and mode using fillna() method Describe data --> which can give statistical analysis

```
In [3]:  train_data = pd.read_excel("Data_Train.xlsx")
```

```
In [4]:  pd.set_option('display.max_columns', None)
```

```
In [5]:  train_data.head()
```

Out[5]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_Stops | Additional_Info | Price |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | 24/03/2019 | Banglore | New Delhi | BLR → DEL | 22:20 | 01:10 22 Mar | 2h 50m | non-stop | No info | 3897 |
| 1 | Air India | 1/05/2019 | Kolkata | Banglore | CCU → IXR → BBI → BLR | 05:50 | 13:15 | 7h 25m | 2 stops | No info | 7662 |
| 2 | Jet Airways | 9/06/2019 | Delhi | Cochin | DEL → LKO → BOM → COK | 09:25 | 04:25 10 Jun | 19h | 2 stops | No info | 13882 |
| 3 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU → NAG → BLR | 18:05 | 23:30 | 5h 25m | 1 stop | No info | 6218 |
| 4 | IndiGo | 01/03/2019 | Banglore | New Delhi | BLR → NAG → DEL | 16:50 | 21:35 | 4h 45m | 1 stop | No info | 13302 |

```
In [6]:  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 [7]:  train_data["Duration"].value_counts()
```

```
Out[7]:  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 [8]:  train_data.dropna(inplace = True)
```

```
In [9]:  train_data.isnull().sum()
```

```
Out[9]:  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
```

# EDA

From description we can see that Date_of_Journey is a object data type, Therefore, we have to convert this datatype into timestamp so as to use this column properly for prediction

For this we require pandas to_datetime to convert object data type to datetime dtype.

.dt.day method will extract only day of that date .dt.month method will extract only month of that date

```
In [10]: train_data["Journey_day"] = pd.to_datetime(train_data.Date_of_Journey, format="%d/%m/%Y").dt.day
```

```
In [11]: train_data["Journey_month"] = pd.to_datetime(train_data["Date_of_Journey"], format = "%d/%m/%Y").dt.month
```

```
In [12]: train_data.head()
```

Out[12]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_Stops | Additional_Info | Price | Journey_da |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | 24/03/2019 | Banglore | New Delhi | BLR → DEL | 22:20 | 01:10 22 Mar | 2h 50m | non-stop | No info | 3897 | 2 |
| 1 | Air India | 1/05/2019 | Kolkata | Banglore | CCU → IXR → BBI → BLR | 05:50 | 13:15 | 7h 25m | 2 stops | No info | 7662 | |
| 2 | Jet Airways | 9/06/2019 | Delhi | Cochin | DEL → LKO → BOM → COK | 09:25 | 04:25 10 Jun | 19h | 2 stops | No info | 13882 | |
| 3 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU → NAG → BLR | 18:05 | 23:30 | 5h 25m | 1 stop | No info | 6218 | 1 |
| 4 | IndiGo | 01/03/2019 | Banglore | New Delhi | BLR → NAG → DEL | 16:50 | 21:35 | 4h 45m | 1 stop | No info | 13302 | |

```
In [13]: # Since we have converted Date_of_Journey column into integers, Now we can drop as it is of no use.

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

```
In [14]: # Departure time is when a plane leaves the gate.
         # Similar to Date_of_Journey we can extract values from Dep_Time

         # Extracting Hours
         train_data["Dep_hour"] = pd.to_datetime(train_data["Dep_Time"]).dt.hour

         # Extracting Minutes
         train_data["Dep_min"] = pd.to_datetime(train_data["Dep_Time"]).dt.minute

         # Now we can drop Dep_Time as it is of no use
         train_data.drop(["Dep_Time"], axis = 1, inplace = True)
```

```
In [15]: train_data.head()
```

| | Airline | Source | Destination | Route | Arrival_Time | Duration | Total_Stops | Additional_Info | Price | Journey_day | Journey_month | Dep_hour |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | Banglore | New Delhi | BLR → DEL | 01:10 22 Mar | 2h 50m | non-stop | No info | 3897 | 24 | 3 | 22 |
| 1 | Air India | Kolkata | Banglore | CCU → IXR → BBI → BLR | 13:15 | 7h 25m | 2 stops | No info | 7662 | 1 | 5 | 5 |
| 2 | Jet Airways | Delhi | Cochin | DEL → LKO → BOM → COK | 04:25 10 Jun | 19h | 2 stops | No info | 13882 | 9 | 6 | 9 |
| 3 | IndiGo | Kolkata | Banglore | CCU → NAG → BLR | 23:30 | 5h 25m | 1 stop | No info | 6218 | 12 | 5 | 18 |
| 4 | IndiGo | Banglore | New Delhi | BLR → NAG → DEL | 21:35 | 4h 45m | 1 stop | No info | 13302 | 1 | 3 | 16 |

In [16]:
```python
# Arrival time is when the plane pulls up to the gate.
# Similar to Date_of_Journey we can extract values from Arrival_Time

# Extracting Hours
train_data["Arrival_hour"] = pd.to_datetime(train_data.Arrival_Time).dt.hour

# Extracting Minutes
train_data["Arrival_min"] = pd.to_datetime(train_data.Arrival_Time).dt.minute

# Now we can drop Arrival_Time as it is of no use
train_data.drop(["Arrival_Time"], axis = 1, inplace = True)
```

In [17]:
```python
train_data.head()
```

Out[17]:

| | Airline | Source | Destination | Route | Duration | Total_Stops | Additional_Info | Price | Journey_day | Journey_month | Dep_hour | Dep_min | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | Banglore | New Delhi | BLR → DEL | 2h 50m | non-stop | No info | 3897 | 24 | 3 | 22 | 20 | |
| 1 | Air India | Kolkata | Banglore | CCU → IXR → BBI → BLR | 7h 25m | 2 stops | No info | 7662 | 1 | 5 | 5 | 50 | |
| 2 | Jet Airways | Delhi | Cochin | DEL → LKO → BOM → COK | 19h | 2 stops | No info | 13882 | 9 | 6 | 9 | 25 | |
| 3 | IndiGo | Kolkata | Banglore | CCU → NAG → BLR | 5h 25m | 1 stop | No info | 6218 | 12 | 5 | 18 | 5 | |
| 4 | IndiGo | Banglore | New Delhi | BLR → NAG → DEL | 4h 45m | 1 stop | No info | 13302 | 1 | 3 | 16 | 50 | |

In [18]:
```python
# Time taken by plane to reach destination is called Duration
# It is the differnce betwwen Departure Time and Arrival time


# Assigning and converting Duration column into list
duration = list(train_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 durat
```

In [19]:
```
# Adding duration_hours and duration_mins list to train_data dataframe

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

In [20]:
```
train_data.drop(["Duration"], axis = 1, inplace = True)
```

In [21]:
```
train_data.head()
```

Out[21]:

| | Airline | Source | Destination | Route | Total_Stops | Additional_Info | Price | Journey_day | Journey_month | Dep_hour | Dep_min | Arrival_hour |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | Banglore | New Delhi | BLR → DEL | non-stop | No info | 3897 | 24 | 3 | 22 | 20 | 1 |
| 1 | Air India | Kolkata | Banglore | CCU → IXR → BBI → BLR | 2 stops | No info | 7662 | 1 | 5 | 5 | 50 | 13 |
| 2 | Jet Airways | Delhi | Cochin | DEL → LKO → BOM → COK | 2 stops | No info | 13882 | 9 | 6 | 9 | 25 | 4 |
| 3 | IndiGo | Kolkata | Banglore | CCU → NAG → BLR | 1 stop | No info | 6218 | 12 | 5 | 18 | 5 | 23 |
| 4 | IndiGo | Banglore | New Delhi | BLR → NAG → DEL | 1 stop | No info | 13302 | 1 | 3 | 16 | 50 | 21 |

# Handling Categorical Data

One can find many ways to handle categorical data. Some of them categorical data are,

Nominal data --> data are not in any order --> OneHotEncoder is used in this case Ordinal data --> data are in order --> LabelEncoder is used in this case

In [22]:
```
train_data["Airline"].value_counts()
```

Out[22]:
```
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
```
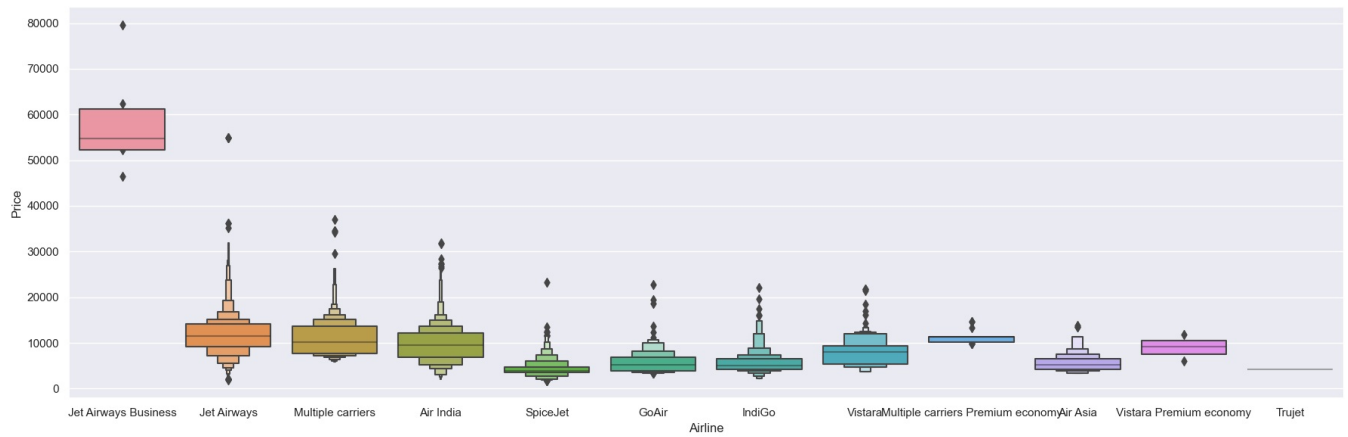
In [23]:
```
# From graph we can see that Jet Airways Business have the highest Price.
# Apart from the first Airline almost all are having similar median

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

```
In [24]:   # As Airline is Nominal Categorical data we will perform OneHotEncoding

           Airline = train_data[["Airline"]]

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

           Airline.head()
```

Out[24]:

| | 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_V |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | |
| 3 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |

```
In [25]:   train_data["Source"].value_counts()
```
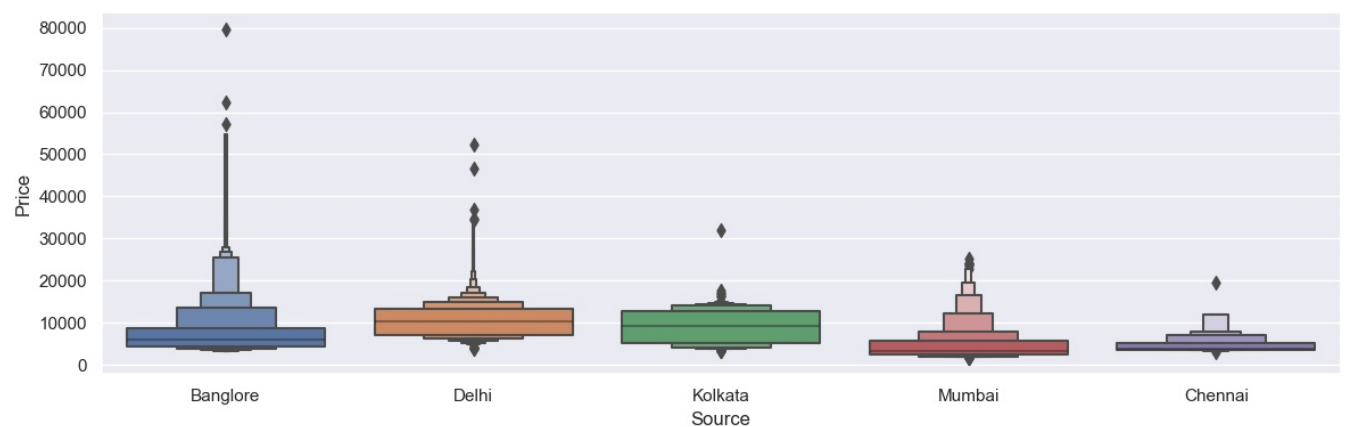
```
Out[25]:   Delhi      4536
           Kolkata    2871
           Banglore   2197
           Mumbai      697
           Chennai     381
           Name: Source, dtype: int64
```

```
In [26]:   # Source vs Price

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



```
In [27]:   # As Source is Nominal Categorical data we will perform OneHotEncoding

           Source = train_data[["Source"]]

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

           Source.head()
```

Out[27]:

| | 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 [28]: train_data["Destination"].value_counts()
```

```
Out[28]: Cochin      4536
         Banglore    2871
         Delhi       1265
         New Delhi    932
         Hyderabad    697
         Kolkata      381
         Name: Destination, dtype: int64
```

```
In [29]: # As Destination is Nominal Categorical data we will perform OneHotEncoding

         Destination = train_data[["Destination"]]

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

         Destination.head()
```

Out[29]:

| | 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 [30]: train_data["Route"]
```

```
Out[30]: 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 [31]: # Additional_Info contains almost 80% no_info
         # Route and Total_Stops are related to each other

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

```
In [32]: train_data["Total_Stops"].value_counts()
```

```
Out[32]: 1 stop      5625
         non-stop    3491
         2 stops     1520
         3 stops       45
         4 stops        1
         Name: Total_Stops, dtype: int64
```

```
In [33]: # As this is case of Ordinal Categorical type we perform LabelEncoder
         # Here Values are assigned with corresponding keys

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

```
In [34]: train_data.head()
```

Out[34]:

| | Airline | Source | Destination | Total_Stops | Price | Journey_day | Journey_month | Dep_hour | Dep_min | Arrival_hour | Arrival_min | Duration_ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | Banglore | New Delhi | 0 | 3897 | 24 | 3 | 22 | 20 | 1 | 10 | |
| 1 | Air India | Kolkata | Banglore | 2 | 7662 | 1 | 5 | 5 | 50 | 13 | 15 | |
| 2 | Jet Airways | Delhi | Cochin | 2 | 13882 | 9 | 6 | 9 | 25 | 4 | 25 | |
| 3 | IndiGo | Kolkata | Banglore | 1 | 6218 | 12 | 5 | 18 | 5 | 23 | 30 | |
| 4 | IndiGo | Banglore | New Delhi | 1 | 13302 | 1 | 3 | 16 | 50 | 21 | 35 | |

```
In [35]:   # Concatenate dataframe --> train_data + Airline + Source + Destination

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

```
In [36]:   data_train.head()
```

Out[36]:

| | Airline | Source | Destination | Total_Stops | Price | Journey_day | Journey_month | Dep_hour | Dep_min | Arrival_hour | Arrival_min | Duration_ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | IndiGo | Banglore | New Delhi | 0 | 3897 | 24 | 3 | 22 | 20 | 1 | 10 | |
| 1 | Air India | Kolkata | Banglore | 2 | 7662 | 1 | 5 | 5 | 50 | 13 | 15 | |
| 2 | Jet Airways | Delhi | Cochin | 2 | 13882 | 9 | 6 | 9 | 25 | 4 | 25 | |
| 3 | IndiGo | Kolkata | Banglore | 1 | 6218 | 12 | 5 | 18 | 5 | 23 | 30 | |
| 4 | IndiGo | Banglore | New Delhi | 1 | 13302 | 1 | 3 | 16 | 50 | 21 | 35 | |

```
In [37]:   data_train.drop(["Airline", "Source", "Destination"], axis = 1, inplace = True)
```

```
In [38]:   data_train.head()
```

Out[38]:

| | Total_Stops | Price | Journey_day | Journey_month | Dep_hour | Dep_min | Arrival_hour | Arrival_min | Duration_hours | Duration_mins | Airline_A Inc |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 3897 | 24 | 3 | 22 | 20 | 1 | 10 | 2 | 50 | |
| 1 | 2 | 7662 | 1 | 5 | 5 | 50 | 13 | 15 | 7 | 25 | |
| 2 | 2 | 13882 | 9 | 6 | 9 | 25 | 4 | 25 | 19 | 0 | |
| 3 | 1 | 6218 | 12 | 5 | 18 | 5 | 23 | 30 | 5 | 25 | |
| 4 | 1 | 13302 | 1 | 3 | 16 | 50 | 21 | 35 | 4 | 45 | |

```
In [39]:   data_train.shape
```

Out[39]:   (10682, 30)

## Test set

```
In [41]:   test_data = pd.read_excel("Test_set.xlsx")
```

```
In [42]:   test_data.head()
```

Out[42]:

| | Airline | Date_of_Journey | Source | Destination | Route | Dep_Time | Arrival_Time | Duration | Total_Stops | Additional_Info |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Jet Airways | 6/06/2019 | Delhi | Cochin | DEL → BOM → COK | 17:30 | 04:25 07 Jun | 10h 55m | 1 stop | No info |
| 1 | IndiGo | 12/05/2019 | Kolkata | Banglore | CCU → MAA → BLR | 06:20 | 10:20 | 4h | 1 stop | No info |
| 2 | Jet Airways | 21/05/2019 | Delhi | Cochin | DEL → BOM → COK | 19:15 | 19:00 22 May | 23h 45m | 1 stop | In-flight meal not included |
| 3 | Multiple carriers | 21/05/2019 | Delhi | Cochin | DEL → BOM → COK | 08:00 | 21:00 | 13h | 1 stop | No info |
| 4 | Air Asia | 24/06/2019 | Banglore | Delhi | BLR → DEL | 23:55 | 02:45 25 Jun | 2h 50m | non-stop | No info |

```
In [43]:   # Preprocessing

           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())

           # EDA

           # Date_of_Journey
```

```python
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)

# Arrival_Time
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
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 durat

# 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)


# Categorical data

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)

# Additional_Info contains almost 80% no_info
# Route and Total_Stops are related to each other
test_data.drop(["Route", "Additional_Info"], axis = 1, inplace = True)

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

# Concatenate dataframe --> test_data + Airline + Source + Destination
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
Banglore    555
Mumbai      186
Chennai      75
Name: Source, dtype: int64

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


Shape of test data :  (2671, 28)
```

In [44]: `data_test.head()`

Out[44]:

| | Total_Stops | Journey_day | Journey_month | Dep_hour | Dep_min | Arrival_hour | Arrival_min | Duration_hours | Duration_mins | Air India | GoAir | Ir |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 6 | 6 | 17 | 30 | 4 | 25 | 10 | 55 | 0 | 0 | |
| 1 | 1 | 12 | 5 | 6 | 20 | 10 | 20 | 4 | 0 | 0 | 0 | |
| 2 | 1 | 21 | 5 | 19 | 15 | 19 | 0 | 23 | 45 | 0 | 0 | |
| 3 | 1 | 21 | 5 | 8 | 0 | 21 | 0 | 13 | 0 | 0 | 0 | |
| 4 | 0 | 24 | 6 | 23 | 55 | 2 | 45 | 2 | 50 | 0 | 0 | |

# Feature Selection

Finding out the best feature which will contribute and have good relation with target variable. Following are some of the feature selection methods,

heatmap feature*importance* SelectKBest

```
In [45]: data_train.shape
```

```
Out[45]: (10682, 30)
```

```
In [46]: data_train.columns
```

```
Out[46]: 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 [48]: 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[48]:

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

```
In [49]: y = data_train.iloc[:, 1]
         y.head()
```

```
Out[49]: 0     3897
         1     7662
         2    13882
         3     6218
         4    13302
         Name: Price, dtype: int64
```

```
In [50]: # Finds correlation between Independent and dependent attributes

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

         plt.show()
```

```
C:\Users\HP\AppData\Local\Temp\ipykernel_16092\3228867913.py:4: FutureWarning: The default value of numeric_onl
y in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or
specify the value of numeric_only to silence this warning.
  sns.heatmap(train_data.corr(), annot = True, cmap = "RdYlGn")
```

In [51]: `# Important feature using ExtraTreesRegressor`

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

Out[51]: ▾ ExtraTreesRegressor

ExtraTreesRegressor()

In [52]: 
```python
print(selection.feature_importances_)
```
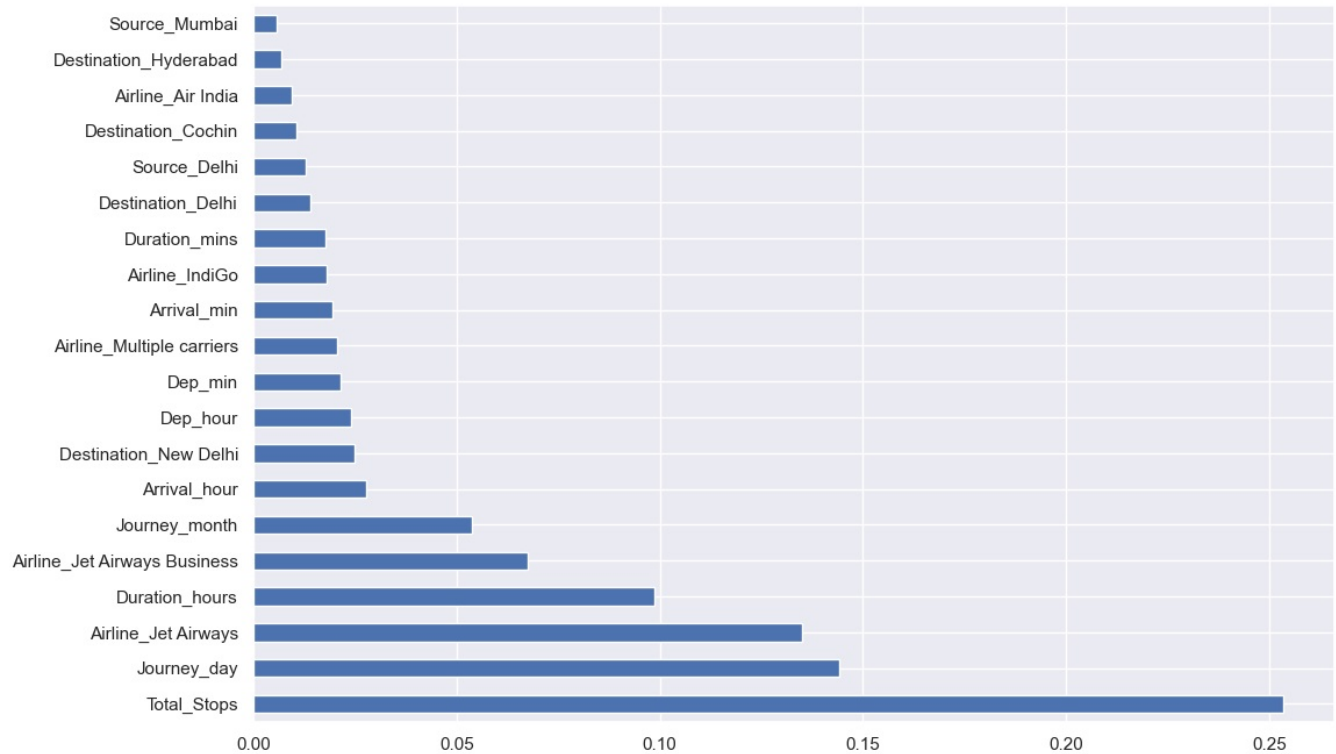
```
[2.53470796e-01 1.44094068e-01 5.36331525e-02 2.39955361e-02
 2.13409191e-02 2.77962739e-02 1.94331763e-02 9.85932605e-02
 1.77649302e-02 9.46535510e-03 1.83641379e-03 1.80382821e-02
 1.35071199e-01 6.75795840e-02 2.05349566e-02 8.26382052e-04
 3.07806504e-03 1.11001738e-04 4.98111153e-03 8.03756456e-05
 4.50273490e-04 1.27076352e-02 3.13285406e-03 5.54249137e-03
 1.04585774e-02 1.39380875e-02 6.77032032e-03 4.24529024e-04
 2.48503928e-02]
```

In [53]:
```python
#plot graph of feature importances for better visualization

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



# Fitting model using Random Forest

Split dataset into train and test set in order to prediction w.r.t X_test If needed do scaling of data Scaling is not done in Random forest Import model Fit the data Predict w.r.t X_test In regression check RSME Score Plot graph

In [54]:
```python
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 [55]:
```python
from sklearn.ensemble import RandomForestRegressor
reg_rf = RandomForestRegressor()
reg_rf.fit(X_train, y_train)
```

Out[55]:
```
▾ RandomForestRegressor

RandomForestRegressor()
```

In [56]:
```python
y_pred = reg_rf.predict(X_test)
```

In [57]:
```python
reg_rf.score(X_train, y_train)
```

Out[57]: 0.9532348142789775

In [58]:
```python
reg_rf.score(X_test, y_test)
```

Out[58]: 0.797010031600254

In [59]:
```python
sns.distplot(y_test-y_pred)
plt.show()
```

```
In [60]: plt.scatter(y_test, y_pred, alpha = 0.5)
         plt.xlabel("y_test")
         plt.ylabel("y_pred")
         plt.show()
```



```
In [61]: from sklearn import metrics
```

```
In [62]: 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: 1173.3705151589168
MSE: 4376881.131875144
RMSE: 2092.0996945354073
```

```
In [63]: # RMSE/(max(DV)-min(DV))

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

```
Out[63]: 0.026887077025966846
```

```
In [64]: metrics.r2_score(y_test, y_pred)

Out[64]: 0.797010031600254
```

# Hyperparameter Tuning

Choose following method for hyperparameter tuning RandomizedSearchCV --> Fast GridSearchCV Assign hyperparameters in form of dictionery Fit the model Check best paramters and best score

```
In [65]: from sklearn.model_selection import RandomizedSearchCV
```

```
In [66]: #Randomized Search CV

         # Number of trees in random forest
         n_estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]
         # Number of features to consider at every split
         max_features = ['auto', 'sqrt']
         # Maximum number of levels in tree
         max_depth = [int(x) for x in np.linspace(5, 30, num = 6)]
         # Minimum number of samples required to split a node
         min_samples_split = [2, 5, 10, 15, 100]
         # Minimum number of samples required at each leaf node
         min_samples_leaf = [1, 2, 5, 10]
```

```
In [67]: # Create the random grid

         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 [69]: # Random search of parameters, using 5 fold cross validation,
         # search across 100 different combinations
         rf_random = RandomizedSearchCV(estimator = reg_rf, param_distributions = random_grid,scoring='neg_mean_squared_
```

```
In [70]: rf_random.fit(X_train,y_train)

         Fitting 5 folds for each of 10 candidates, totalling 50 fits
         [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time
         =   9.8s
         [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time
         =   9.7s
         [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time
         =   9.8s
         [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time
         =   9.4s
         [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time
         =   9.3s
         [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti
         me=   15.3s
         [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti
         me=   15.5s
         [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti
         me=   15.3s
         [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti
         me=   15.6s
         [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti
         me=   15.3s
         C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
         ` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
         es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
         ssors.
           warn(
         [CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti
         me=    8.4s
         C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
         ` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
         es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
         ssors.
           warn(
         [CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti
         me=    8.4s
         C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
         ` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
         es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
         ssors.
           warn(
         [CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti
         me=    8.2s
```

```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti
me=   8.2s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti
me=   7.9s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time
=   14.7s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time
=   16.1s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time
=   15.9s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time
=   15.6s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time
=   16.0s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim
e=   24.1s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim
e=   24.4s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim
e=   24.0s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim
e=   23.8s
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```
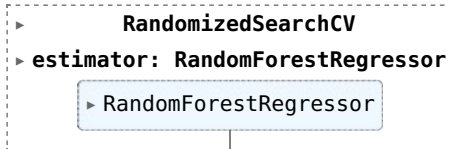
```
[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim
e=  24.3s
[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim
e=  27.4s
[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim
e=  27.4s
[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim
e=  26.7s
[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim
e=  26.7s
[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim
e=  27.0s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti
me=   7.2s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti
me=   7.6s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti
me=   7.3s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti
me=   7.1s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti
me=   7.1s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim
e=   4.0s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim
e=   3.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim
e=   3.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim
e=   3.9s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim
e=   4.0s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time
=   5.5s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time
=   4.5s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time
=   4.6s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time
=   4.6s
[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time
=   4.6s
```

```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```
```
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim
e=  30.9s
```
```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```
```
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim
e=  30.6s
```
```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```
```
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim
e=  30.5s
```
```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```
```
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim
e=  31.1s
```
```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```
```
[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim
e=  31.5s
```
```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\ensemble\_forest.py:413: FutureWarning: `max_features='auto'
` has been deprecated in 1.1 and will be removed in 1.3. To keep the past behaviour, explicitly set `max_featur
es=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegre
ssors.
  warn(
```

```
          ▸          RandomizedSearchCV
        ▸ estimator: RandomForestRegressor
               ▸ RandomForestRegressor
```

In [71]: `rf_random.best_params_`

Out[71]:
```
{'n_estimators': 700,
 'min_samples_split': 15,
 'min_samples_leaf': 1,
 'max_features': 'auto',
 'max_depth': 20}
```

In [72]: `prediction = rf_random.predict(X_test)`

In [73]:
```python
plt.figure(figsize = (8,8))
sns.distplot(y_test-prediction)
plt.show()
```
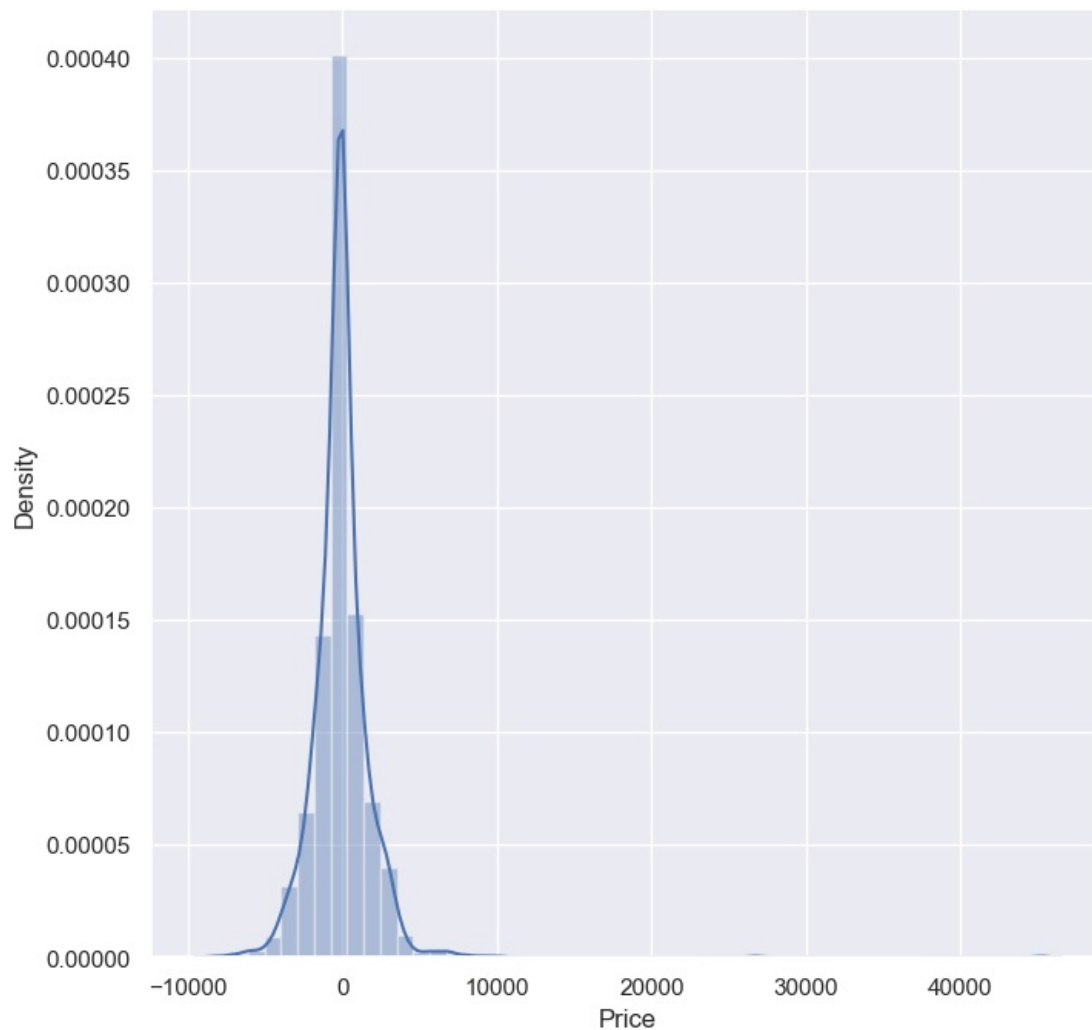
```
C:\Users\HP\AppData\Local\Temp\ipykernel_16092\1574001921.py:2: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

  sns.distplot(y_test-prediction)
```
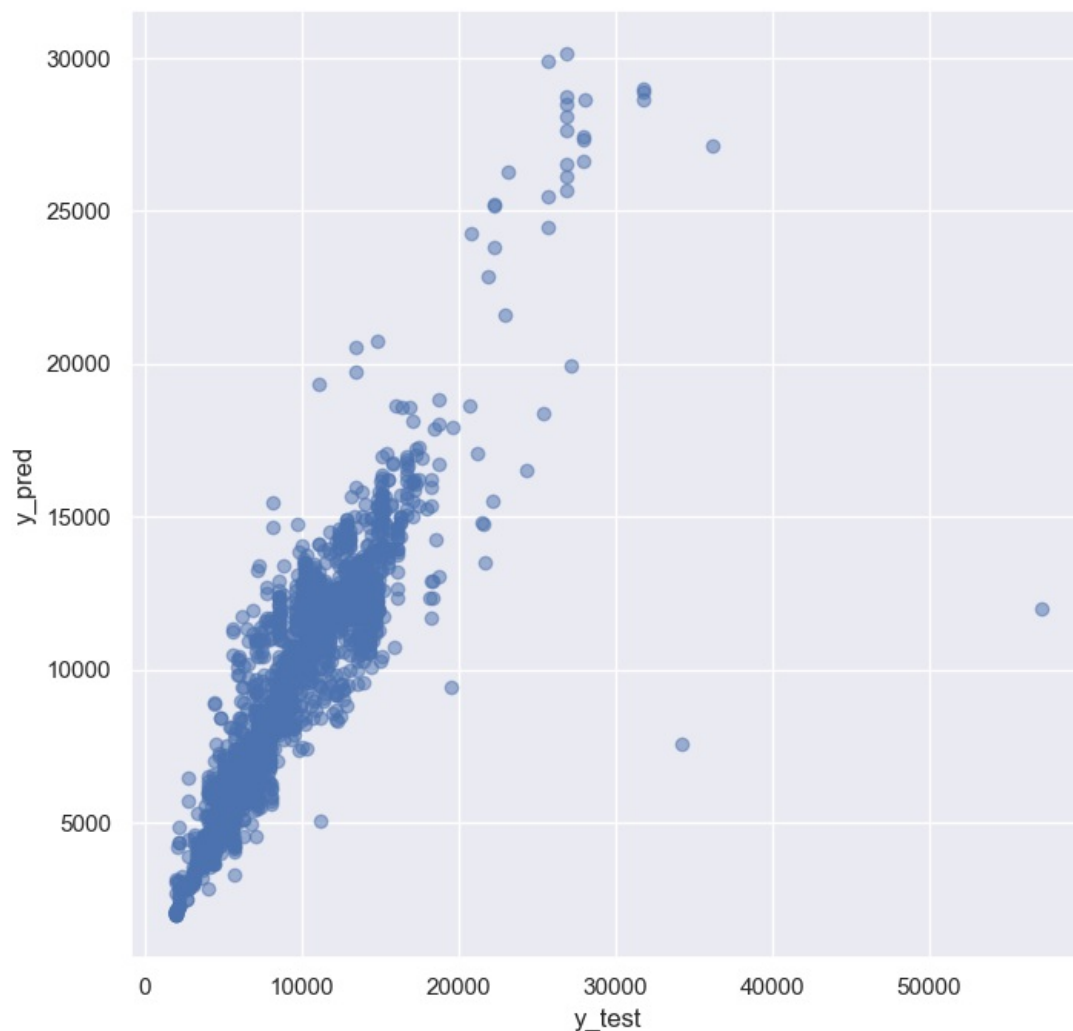


In [74]:
```python
plt.figure(figsize = (8,8))
plt.scatter(y_test, prediction, alpha = 0.5)
plt.xlabel("y_test")
plt.ylabel("y_pred")
plt.show()
```

In [75]:
```python
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: 1164.475741487747
MSE: 4051298.9496157276
RMSE: 2012.7838805037484
```

## Save the model to reuse it again

In [88]:
```python
import pickle
# open a file, where you ant to store the data
file = open('flight_rf.pkl', 'wb')

# dump information to that file
pickle.dump(reg_rf, file)
```

In [92]:
```python
model = open('flight_rf.pkl','rb')
forest = pickle.load(model)
```

In [93]:
```python
y_prediction = forest.predict(X_test)
```

In [94]:
```python
metrics.r2_score(y_test, y_prediction)
```

Out[94]:
```
0.797010031600254
```

In [ ]: