

import libraries

In [145]:

```
1 import pandas as pd
2 import numpy as np
3 import sklearn
4 from sklearn.impute import SimpleImputer
5 from sklearn.compose import ColumnTransformer
6 from sklearn.ensemble import RandomForestRegressor
7 from sklearn.preprocessing import (
8     OneHotEncoder,
9     MinMaxScaler,
10    PowerTransformer,
11    FunctionTransformer,
12    StandardScaler
13 )
14 from sklearn.pipeline import Pipeline, FeatureUnion
15 from feature_engine.encoding import (
16     RareLabelEncoder,
17     MeanEncoder,
18     CountFrequencyEncoder,
19     OrdinalEncoder
20 )
21 from feature_engine.selection import SelectBySingleFeaturePerformance
22 from sklearn.base import BaseEstimator, TransformerMixin
23 from sklearn.metrics.pairwise import rbf_kernel
24 from feature_engine.datetime import DatetimeFeatures
25 from feature_engine.outliers import Winsorizer
26 import warnings
27 import matplotlib.pyplot as plt
28 import os
29 import xgboost as xgb
```

Display Settings

```
In [2]: 1 pd.set_option('display.max_columns',None)
```

```
In [3]: 1 sklearn.set_config(transform_output='pandas')
```

```
In [4]: 1 warnings.filterwarnings('ignore')
```

Read Data

```
In [5]: 1 path = r"D:\Ml projects\car dekho\flights practice\data\train.csv"
```

```
In [6]: 1 train = pd.read_csv(path)
```

```
In [7]: 1 train.head()
```

Out[7]:

	airline	date_of_journey	source	destination	dep_time	arrival_time	duration	total_stops	additional_info	price
0	Jet Airways	2019-03-21	Banglore	New Delhi	08:55:00	19:10:00	615	1.0	In-flight meal not included	7832
1	Jet Airways	2019-03-27	Delhi	Cochin	17:30:00	04:25:00	655	1.0	In-flight meal not included	6540
2	Goair	2019-03-09	Banglore	New Delhi	11:40:00	14:35:00	175	0.0	No Info	7305
3	Air India	2019-06-12	Kolkata	Banglore	09:25:00	18:30:00	545	1.0	No Info	8366
4	Jet Airways	2019-03-12	Banglore	New Delhi	22:55:00	07:40:00	525	1.0	In-flight meal not included	11087

In [8]: 1 train.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6695 entries, 0 to 6694
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   airline          6695 non-null    object  
 1   date_of_journey  6695 non-null    object  
 2   source           6695 non-null    object  
 3   destination      6695 non-null    object  
 4   dep_time         6695 non-null    object  
 5   arrival_time     6695 non-null    object  
 6   duration         6695 non-null    int64  
 7   total_stops      6694 non-null    float64 
 8   additional_info  6695 non-null    object  
 9   price            6695 non-null    int64  
dtypes: float64(1), int64(2), object(7)
memory usage: 523.2+ KB
```

In [9]: 1 train.dropna(inplace=True)

In [10]: 1 train.info()

```
<class 'pandas.core.frame.DataFrame'>
Index: 6694 entries, 0 to 6694
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   airline          6694 non-null    object  
 1   date_of_journey  6694 non-null    object  
 2   source           6694 non-null    object  
 3   destination      6694 non-null    object  
 4   dep_time         6694 non-null    object  
 5   arrival_time     6694 non-null    object  
 6   duration         6694 non-null    int64  
 7   total_stops      6694 non-null    float64 
 8   additional_info  6694 non-null    object  
 9   price            6694 non-null    int64  
dtypes: float64(1), int64(2), object(7)
memory usage: 575.3+ KB
```

```
In [11]: 1 x_train = train.drop(columns='price')
          2 y_train = train.price.copy()
```

Transformation Operation

```
In [12]: 1 x_train.airline
```

```
Out[12]: 0      Jet Airways
          1      Jet Airways
          2          Goair
          3      Air India
          4      Jet Airways
          ...
          6690    Jet Airways
          6691    Air India
          6692    Jet Airways
          6693    Air Asia
          6694    Air India
Name: airline, Length: 6694, dtype: object
```

In [13]:

```
1 air_transformer = Pipeline(steps=[  
2     ('imputer', SimpleImputer(strategy='most_frequent')),  
3     ('grouper', RareLabelEncoder(tol=0.1, replace_with='Other', n_categories=2)),  
4     ('encoder', OneHotEncoder(sparse_output=False, handle_unknown='ignore'))  
5 ])  
6  
7 air_transformer.fit_transform(x_train.loc[:, ['airline']])
```

Out[13]:

	airline_Air India	airline_Indigo	airline_Jet Airways	airline_Multiple Carriers	airline_Other
0	0.0	0.0	1.0	0.0	0.0
1	0.0	0.0	1.0	0.0	0.0
2	0.0	0.0	0.0	0.0	1.0
3	1.0	0.0	0.0	0.0	0.0
4	0.0	0.0	1.0	0.0	0.0
...
6690	0.0	0.0	1.0	0.0	0.0
6691	1.0	0.0	0.0	0.0	0.0
6692	0.0	0.0	1.0	0.0	0.0
6693	0.0	0.0	0.0	0.0	1.0
6694	1.0	0.0	0.0	0.0	0.0

6694 rows × 5 columns

date_of_journey

In [14]:

```
1 feature_to_extract = ['week', 'month', 'day_of_week', 'day_of_year']
```

In [15]:

```
1 doj_transformer = Pipeline(steps=[  
2     ('dt', DatetimeFeatures(features_to_extract=feature_to_extract, yearfirst=True, format='mixed')),  
3     ('scaler', MinMaxScaler())  
4 ])  
5 doj_transformer.fit_transform(x_train.loc[:, ['date_of_journey']])
```

Out[15]:

	date_of_journey_week	date_of_journey_month	date_of_journey_day_of_week	date_of_journey_day_of_year
0	0.176471	0.000000	0.500000	0.169492
1	0.235294	0.000000	0.333333	0.220339
2	0.058824	0.000000	0.833333	0.067797
3	0.882353	1.000000	0.333333	0.872881
4	0.117647	0.000000	0.166667	0.093220
...
6690	0.176471	0.000000	0.500000	0.169492
6691	0.529412	0.666667	0.333333	0.516949
6692	0.764706	1.000000	0.833333	0.779661
6693	1.000000	1.000000	0.000000	0.974576
6694	0.000000	0.000000	0.666667	0.000000

6694 rows × 4 columns

source and destination

```
In [16]: 1 x_train.source
```

```
Out[16]: 0      Banglore
1      Delhi
2      Banglore
3      Kolkata
4      Banglore
...
6690      Delhi
6691      Kolkata
6692      Delhi
6693      Delhi
6694      Banglore
Name: source, Length: 6694, dtype: object
```

```
In [17]: 1 x_train.destination
```

```
Out[17]: 0      New Delhi
1      Cochin
2      New Delhi
3      Banglore
4      New Delhi
...
6690      Cochin
6691      Banglore
6692      Cochin
6693      Cochin
6694      New Delhi
Name: destination, Length: 6694, dtype: object
```

```
In [18]: 1 location_subset = x_train.loc[:,['source','destination']]  
2 location_subset
```

Out[18]:

	source	destination
0	Banglore	New Delhi
1	Delhi	Cochin
2	Banglore	New Delhi
3	Kolkata	Banglore
4	Banglore	New Delhi
...
6690	Delhi	Cochin
6691	Kolkata	Banglore
6692	Delhi	Cochin
6693	Delhi	Cochin
6694	Banglore	New Delhi

6694 rows × 2 columns

In [19]:

```
1 location_pipe1=Pipeline(steps=[  
2     ('grouper',RareLabelEncoder(tol=0.1,replace_with='Other',n_categories=2)),  
3     ('encoder',MeanEncoder()),  
4     ('transformer',PowerTransformer())  
5 ])  
6 location_pipe1.fit_transform(location_subset,y_train)
```

Out[19]:

	source	destination
0	-0.857629	-0.736209
1	1.065619	1.061892
2	-0.857629	-0.736209
3	-0.203923	-0.224330
4	-0.857629	-0.736209
...
6690	1.065619	1.061892
6691	-0.203923	-0.224330
6692	1.065619	1.061892
6693	1.065619	1.061892
6694	-0.857629	-0.736209

6694 rows × 2 columns

In [20]:

```
1 def is_north(x):
2     columns = x.columns.to_list()
3     north_cities=['Delhi', 'Kolkata', 'Mumbai', 'New Delhi']
4     return (
5         x.
6         assign(**{
7             f'{col}_is_north':x.loc[:,col].isin(north_cities).astype(int)
8             for col in columns
9         })
10        .drop(columns=columns)
11    )
12 FunctionTransformer(func=is_north).fit_transform(location_subset)
```

Out[20]:

	source_is_north	destination_is_north
0	0	1
1	1	0
2	0	1
3	1	0
4	0	1
...
6690	1	0
6691	1	0
6692	1	0
6693	1	0
6694	0	1

6694 rows × 2 columns

```
In [21]: 1 location_transformation = FeatureUnion(transformer_list=[  
2     ('part1',location_pipe1),  
3     ('part2',FunctionTransformer(func=is_north))  
4 ])  
5 location_transformation.fit_transform(location_subset,y_train)
```

Out[21]:

	source	destination	source_is_north	destination_is_north
0	-0.857629	-0.736209	0	1
1	1.065619	1.061892	1	0
2	-0.857629	-0.736209	0	1
3	-0.203923	-0.224330	1	0
4	-0.857629	-0.736209	0	1
...
6690	1.065619	1.061892	1	0
6691	-0.203923	-0.224330	1	0
6692	1.065619	1.061892	1	0
6693	1.065619	1.061892	1	0
6694	-0.857629	-0.736209	0	1

6694 rows × 4 columns

dep_time and arrival_time

In [22]:

```
1 time_subset = x_train.loc[:, ['dep_time', 'arrival_time']]  
2 time_subset
```

Out[22]:

	dep_time	arrival_time
0	08:55:00	19:10:00
1	17:30:00	04:25:00
2	11:40:00	14:35:00
3	09:25:00	18:30:00
4	22:55:00	07:40:00
...
6690	10:45:00	18:50:00
6691	09:25:00	18:30:00
6692	14:00:00	19:00:00
6693	07:55:00	13:25:00
6694	11:50:00	08:55:00

6694 rows × 2 columns

In [23]:

```
1 time_pipe1 = Pipeline(steps=[  
2     ('dt', DatetimeFeatures(features_to_extract=['hour', 'minute'])),  
3     ('scaler', MinMaxScaler())  
4 ])  
5 time_pipe1.fit_transform(time_subset)
```

Out[23]:

	dep_time_hour	dep_time_minute	arrival_time_hour	arrival_time_minute
0	0.347826	1.000000	0.826087	0.181818
1	0.739130	0.545455	0.173913	0.454545
2	0.478261	0.727273	0.608696	0.636364
3	0.391304	0.454545	0.782609	0.545455
4	0.956522	1.000000	0.304348	0.727273
...
6690	0.434783	0.818182	0.782609	0.909091
6691	0.391304	0.454545	0.782609	0.545455
6692	0.608696	0.000000	0.826087	0.000000
6693	0.304348	1.000000	0.565217	0.454545
6694	0.478261	0.909091	0.347826	1.000000

6694 rows × 4 columns

In [24]:

```
1 def part_of_day(x, morning=4, noon=12, eve=16, night=20):
2     columns = x.columns.to_list()
3     x_temp = x.assign(**{
4         col: pd.to_datetime(x.loc[:, col]).dt.hour
5         for col in columns
6     })
7     return (
8         x_temp
9         .assign(**{
10             f'{col}_part_of_day': np.select(
11                 [x_temp.loc[:, col].between(morning, noon, inclusive='left'),
12                  x_temp.loc[:, col].between(noon, eve, inclusive='left'),
13                  x_temp.loc[:, col].between(eve, night, inclusive='left')], [
14                      'morning', 'afternoon', 'evening'],
15                      default = 'night'
16                 )
17                 for col in columns
18             })
19             .drop(columns=columns)
20         )
21     FunctionTransformer(func=part_of_day).fit_transform(time_subset)
```

Out[24]:

	dep_time_part_of_day	arrival_time_part_of_day
0	morning	evening
1	evening	morning
2	morning	afternoon
3	morning	evening
4	night	morning
...
6690	morning	evening
6691	morning	evening
6692	afternoon	evening
6693	morning	afternoon
6694	morning	morning

6694 rows × 2 columns

In [25]:

```
1 time_pipe2 = Pipeline(steps=[  
2     ('part1', FunctionTransformer(func=part_of_day)),  
3     ('encoder', CountFrequencyEncoder()),  
4     ('scaler', MinMaxScaler())  
5 ])  
6 time_pipe2.fit_transform(time_subset)
```

Out[25]:

	dep_time_part_of_day	arrival_time_part_of_day
0	1.000000	0.667335
1	0.202773	0.951904
2	1.000000	0.000000
3	1.000000	0.667335
4	0.174177	0.951904
...
6690	1.000000	0.667335
6691	1.000000	0.667335
6692	0.000000	0.667335
6693	1.000000	0.000000
6694	1.000000	0.951904

6694 rows × 2 columns

In [26]:

```
1 time_transformer = FeatureUnion(transformer_list=[  
2     ('part1',time_pipe1),  
3     ('pipe2',time_pipe2)  
4 ])  
5 time_transformer.fit_transform(time_subset)
```

Out[26]:

	dep_time_hour	dep_time_minute	arrival_time_hour	arrival_time_minute	dep_time_part_of_day	arrival_time_part_of_day
0	0.347826	1.000000	0.826087	0.181818	1.000000	0.667335
1	0.739130	0.545455	0.173913	0.454545	0.202773	0.951904
2	0.478261	0.727273	0.608696	0.636364	1.000000	0.000000
3	0.391304	0.454545	0.782609	0.545455	1.000000	0.667335
4	0.956522	1.000000	0.304348	0.727273	0.174177	0.951904
...
6690	0.434783	0.818182	0.782609	0.909091	1.000000	0.667335
6691	0.391304	0.454545	0.782609	0.545455	1.000000	0.667335
6692	0.608696	0.000000	0.826087	0.000000	0.000000	0.667335
6693	0.304348	1.000000	0.565217	0.454545	1.000000	0.000000
6694	0.478261	0.909091	0.347826	1.000000	1.000000	0.951904

6694 rows × 6 columns

duration

In [27]:

```
1 ( 
2     x_train
3     .duration
4     .quantile([0.25,0.5,0.75])
5     .values
6     .reshape(-1,1)
7     # .shape
8 )
```

Out[27]: array([[170.],
 [510.],
 [920.]])

In [28]:

```
1  class RBFpercentileSimilarity(BaseEstimator,TransformerMixin):
2      def __init__(self,variables=None,percentiles=[0.25,0.50,0.75],gamma=0.1):
3          self.variables = variables
4          self.percentiles = percentiles
5          self.gamma = gamma
6
7      def fit(self,x,y=None):
8          if not self.variables:
9              self.variables = x.select_dtypes(include='number').columns.to_list()
10
11         self.reference_values_ = {
12             col:(
13                 x
14                 .loc[:,col]
15                 .quantile(self.percentiles)
16                 .values
17                 .reshape(-1,1)
18             )
19             for col in self.variables
20         }
21
22         return self
23
24     def transform(self,x):
25         objects=[]
26         for col in self.variables:
27             columns = [f'{col}_rbf_int{((percentile*100)}' for percentile in self.percentiles]
28             obj = pd.DataFrame(
29                 data = rbf_kernel(x.loc[:,[col]],self.reference_values_[col],gamma=self.gamma),
30                 columns=columns
31             )
32             objects.append(obj)
33
34         return pd.concat(objects, axis=1)
```

```
In [29]: 1 RBFpercentileSimilarity().fit_transform(x_train.loc[:, ['duration']])
```

Out[29]:

	duration_rbf_int25.0	duration_rbf_int50.0	duration_rbf_int75.0
0	0.000000	0.000000e+00	0.0
1	0.000000	0.000000e+00	0.0
2	0.082085	0.000000e+00	0.0
3	0.000000	6.293989e-54	0.0
4	0.000000	1.691898e-10	0.0
...
6689	0.000000	0.000000e+00	0.0
6690	0.000000	6.293989e-54	0.0
6691	0.000000	0.000000e+00	0.0
6692	0.000000	0.000000e+00	0.0
6693	0.000000	0.000000e+00	0.0

6694 rows × 3 columns

```
In [30]: 1 def duration_category(X, short=180, med=400):  
2     return (  
3         X  
4         .assign(duration_cat=np.select([X.duration.lt(short),  
5                                         X.duration.between(short, med, inclusive="left")],  
6                                         ["short", "medium"],  
7                                         default="long"))  
8         .drop(columns="duration")  
9     )
```

In [31]:

```
1 def is_over(X, value=1000):
2     return (
3         X
4         .assign(**{
5             f"duration_over_{value}": X.duration.ge(value).astype(int)
6         })
7         .drop(columns="duration")
8     )
```

In [32]:

```
1 duration_pipe1 = Pipeline(steps=[  
2     # ('rbf', RBFpercentileSimilarity()),  
3     ('scaler', PowerTransformer())  
4 ])  
5  
6 duration_union = FeatureUnion(transformer_list=[  
7     ('part1', duration_pipe1),  
8     ('part2', StandardScaler())  
9 ])  
10 duration_transformer = Pipeline(steps=[  
11     ('outlier', Winsorizer(capping_method='iqr', fold=1.5)),  
12     ('Imputer', SimpleImputer(strategy='median')),  
13     ('union', duration_union)  
14 ])  
15 duration_transformer.fit_transform(x_train.loc[:, ['duration']])
```

Out[32]:

	part1_duration	part2_duration
0	0.330340	-0.033600
1	0.403254	0.046768
2	-1.020908	-0.917646
3	0.191947	-0.174244
4	0.149506	-0.214428
...
6690	1.733333	2.598445
6691	0.191947	-0.174244
6692	-0.464892	-0.666496
6693	-0.362973	-0.606221
6694	1.196291	1.272376

6694 rows × 2 columns

In [33]:

```
1 ( 
2     x_train
3     .duration
4     .loc[lambda ser: ser > 1000]
5     .to_frame()
6     .quantile([0.25,0.50,0.75])
7 )
```

Out[33]:

	duration
0.25	1220.0
0.50	1420.0
0.75	1590.0

total_stops

In [34]:

```
1 def is_direct(X):
2     return X.assign(is_direct_flight=X.total_stops.eq(0).astype(int))
3
4
5 total_stops_transformer = Pipeline(steps=[
6     ("imputer", SimpleImputer(strategy="most_frequent")),
7     ("", FunctionTransformer(func=is_direct))
8 ])
9
10 total_stops_transformer.fit_transform(x_train.loc[:, ["total_stops"]])
```

Out[34]:

	total_stops	is_direct_flight
0	1.0	0
1	1.0	0
2	0.0	1
3	1.0	0
4	1.0	0
...
6690	2.0	0
6691	1.0	0
6692	1.0	0
6693	1.0	0
6694	1.0	0

6694 rows × 2 columns

additional_info

In [35]:

```
1 info_pipe1 = Pipeline(steps=[  
2     ("group", RareLabelEncoder(tol=0.1, n_categories=2, replace_with="Other")),  
3     ("encoder", OneHotEncoder(handle_unknown="ignore", sparse_output=False))  
4 ])  
5  
6 info_pipe1.fit_transform(x_train.loc[:, ["additional_info"]])
```

Out[35]:

	additional_info_In-flight meal not included	additional_info_No Info	additional_info_Other
0	1.0	0.0	0.0
1	1.0	0.0	0.0
2	0.0	1.0	0.0
3	0.0	1.0	0.0
4	1.0	0.0	0.0
...
6690	0.0	1.0	0.0
6691	0.0	1.0	0.0
6692	1.0	0.0	0.0
6693	0.0	1.0	0.0
6694	0.0	0.0	1.0

6694 rows × 3 columns

In [36]:

```
1 def have_info(X):
2     return X.assign(additional_info=X.additional_info.ne("No Info").astype(int))
3 info_union = FeatureUnion(transformer_list=[
4     ("part1", info_pipe1),
5     ("part2", FunctionTransformer(func=have_info))
6 ])
7 info_transformer = Pipeline(steps=[
8     ("imputer", SimpleImputer(strategy="constant", fill_value="unknown")),
9     ("union", info_union)
10])
11
12 info_transformer.fit_transform(x_train.loc[:, ["additional_info"]])
```

Out[36]:

	additional_info_In-flight meal not included	additional_info_No Info	additional_info_Other	additional_info
0	1.0	0.0	0.0	1
1	1.0	0.0	0.0	1
2	0.0	1.0	0.0	0
3	0.0	1.0	0.0	0
4	1.0	0.0	0.0	1
...
6690	0.0	1.0	0.0	0
6691	0.0	1.0	0.0	0
6692	1.0	0.0	0.0	1
6693	0.0	1.0	0.0	0
6694	0.0	0.0	1.0	1

6694 rows × 4 columns

In []:

1

```
In [37]: 1 x_train.columns.to_list()
```

```
Out[37]: ['airline',  
          'date_of_journey',  
          'source',  
          'destination',  
          'dep_time',  
          'arrival_time',  
          'duration',  
          'total_stops',  
          'additional_info']
```

Column Transformer

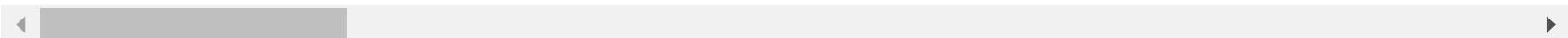
In [38]:

```
1 column_transformer = ColumnTransformer(transformers=[
2     ('air',air_transformer,['airline']),
3     ('doj',doj_transformer,['date_of_journey']),
4     ('location',location_transformation,['source','destination']),
5     ('time',time_transformer,['dep_time','arrival_time']),
6     ('dur',duration_transformer,['duration']),
7     ("stops", total_stops_transformer, ["total_stops"]),
8     ("info", info_transformer, ["additional_info"])
9 ],remainder='passthrough')
10 column_transformer.fit_transform(x_train,y_train)
```

Out[38]:

	air_airline_Air India	air_airline_Indigo	air_airline_Jet Airways	air_airline_Multiple Carriers	air_airline_Other	doj_date_of_journey_week	doj_date_of_journey_year
0	0.0	0.0	1.0	0.0	0.0	0.176471	0.0
1	0.0	0.0	1.0	0.0	0.0	0.235294	0.0
2	0.0	0.0	0.0	0.0	1.0	0.058824	0.0
3	1.0	0.0	0.0	0.0	0.0	0.882353	0.0
4	0.0	0.0	1.0	0.0	0.0	0.117647	0.0
...
6690	0.0	0.0	1.0	0.0	0.0	0.176471	0.0
6691	1.0	0.0	0.0	0.0	0.0	0.529412	0.0
6692	0.0	0.0	1.0	0.0	0.0	0.764706	0.0
6693	0.0	0.0	0.0	0.0	1.0	1.000000	0.0
6694	1.0	0.0	0.0	0.0	0.0	0.000000	0.0

6694 rows × 27 columns



Feature Selection

In [39]:

```

1 estimator = RandomForestRegressor(n_estimators=10, max_depth=3, random_state=42)
2
3 selector = SelectBySingleFeaturePerformance(
4     estimator=estimator,
5     scoring="r2",
6     threshold=0.1
7 )

```

In [40]:

```

1 preprocessor = Pipeline(steps=[
2     ("ct", column_transformer),
3     ("selector", selector)
4 ])
5
6 preprocessor.fit_transform(x_train, y_train)

```

Out[40]:

	air_airline_Indigo	air_airline_Jet Airways	air_airline_Other	doj_date_of_journey_week	doj_date_of_journey_day_of_year	location_sor
0	0.0	1.0	0.0	0.176471	0.169492	-0.857
1	0.0	1.0	0.0	0.235294	0.220339	1.065
2	0.0	0.0	1.0	0.058824	0.067797	-0.857
3	0.0	0.0	0.0	0.882353	0.872881	-0.203
4	0.0	1.0	0.0	0.117647	0.093220	-0.857
...
6690	0.0	1.0	0.0	0.176471	0.169492	1.065
6691	0.0	0.0	0.0	0.529412	0.516949	-0.203
6692	0.0	1.0	0.0	0.764706	0.779661	1.065
6693	0.0	0.0	1.0	1.000000	0.974576	1.065
6694	0.0	0.0	0.0	0.000000	0.000000	-0.857

6694 rows × 11 columns

```
In [41]: 1 feature_performances = preprocessor.named_steps['selector'].feature_performance_
2 feature_performances
```

```
Out[41]: {'air_airline_Air India': 0.00204748830602682,
'air_airline_Indigo': 0.12836215285886088,
'air_airline_Jet Airways': 0.19344947492438477,
'air_airline_Multiple Carriers': 0.019061382721082227,
'air_airline_Other': 0.11846058299683764,
'doj_date_of_journey_week': 0.18545205813867202,
'doj_date_of_journey_month': 0.0891930219053337,
'doj_date_of_journey_day_of_week': 0.005234570009364818,
'doj_date_of_journey_day_of_year': 0.2271022061742051,
'location_source': 0.12699498762521164,
'location_destination': 0.1311104993210662,
'location_source_is_north': 0.02964117275449485,
'location_destination_is_north': 0.02964117275449485,
'time_dep_time_hour': 0.0074911375594401974,
'time_dep_time_minute': 0.03794149357137192,
'time_arrival_time_hour': 0.08068284546718978,
'time_arrival_time_minute': 0.031190920907649595,
'time_dep_time_part_of_day': -0.0011295681198565388,
'time_arrival_time_part_of_day': 0.0314661311801779,
'dur_part1_duration': 0.4247739678003379,
'dur_part2_duration': 0.4248389453042105,
'stops_total_stops': 0.4034522557365127,
'stops_is_direct_flight': 0.37941238701968255,
'info_additional_info_In-flight meal not included': 0.0017193556330226494,
'info_additional_info_No Info': -0.0008149235823306326,
'info_additional_info_Other': 0.01778001687368147,
'info_additional_info': -0.0008149235823306326}
```

```
In [42]: 1 sorted_feat_imp = dict(sorted(feature_performances.items(), key=lambda val: val[1]))
2 sorted_feat_imp
```

```
Out[42]: {'time_dep_time_part_of_day': -0.0011295681198565388,
'info_additional_info_No Info': -0.0008149235823306326,
'info_additional_info': -0.0008149235823306326,
'info_additional_info_In-flight meal not included': 0.0017193556330226494,
'air_airline_Air India': 0.00204748830602682,
'doj_date_of_journey_day_of_week': 0.005234570009364818,
'time_dep_time_hour': 0.0074911375594401974,
'info_additional_info_Other': 0.01778001687368147,
'air_airline_Multiple Carriers': 0.019061382721082227,
'location_source_is_north': 0.02964117275449485,
'location_destination_is_north': 0.02964117275449485,
'time_arrival_time_minute': 0.031190920907649595,
'time_arrival_time_part_of_day': 0.0314661311801779,
'time_dep_time_minute': 0.03794149357137192,
'time_arrival_time_hour': 0.08068284546718978,
'doj_date_of_journey_month': 0.0891930219053337,
'air_airline_Other': 0.11846058299683764,
'location_source': 0.12699498762521164,
'air_airline_Indigo': 0.12836215285886088,
'location_destination': 0.1311104993210662,
'doj_date_of_journey_week': 0.18545205813867202,
'air_airline_Jet Airways': 0.19344947492438477,
'doj_date_of_journey_day_of_year': 0.2271022061742051,
'stops_is_direct_flight': 0.37941238701968255,
'stops_total_stops': 0.4034522557365127,
'dur_part1_duration': 0.4247739678003379,
'dur_part2_duration': 0.4248389453042105}
```


In [43]:

```
1 THRESOLD = 0.1
2 selected_bar = None
3 dropped_bar = None
4 colors = ['red' if score < THRESOLD else 'green' for score in sorted_feat_imp.values()]
5
6 fig,ax = plt.subplots(figsize=(15,4))
7
8 for i,(feature,score) in enumerate(sorted_feat_imp.items()):
9     params = dict(
10         x=i,
11         height = score,
12         edgecolor='black',
13         alpha=0.5
14     )
15     if score<THRESOLD:
16         bar = ax.bar(
17             color = 'red',
18             **params
19         )
20     if not dropped_bar:
21         dropped_bar=bar[0]
22     else:
23         bar = ax.bar(
24             color='green',
25             **params
26         )
27         if not selected_bar:
28             selected_bar = bar[0]
29
30 thresh_line = ax.axhline(
31     y=0.1,
32     color="black",
33     linestyle="--"
34 )
35
36 ax.set_xticks(
37     ticks=range(len(sorted_feat_imp)),
38     labels=list(sorted_feat_imp.keys()),
39     rotation=30,
40     ha="right"
41 )
42
43 ax.set(
```

```

44     xlabel="Feature",
45     ylabel="Score",
46     title="Feature Selection Scores"
47 )
48
49 ax.legend(
50     handles=[selected_bar, dropped_bar, thresh_line],
51     labels=["Selected", "Dropped", "Threshold"],
52     loc="upper left"
53 )
54
55 plt.show()

```



Preprocess Data

```
In [44]: 1 PATH_DIR = r'D:\ML projects\car dekho\flights practice\data'
```

```
In [51]: 1 def get_data(data):  
2     file = f'{data}.csv'  
3     path = os.path.join(PATH_DIR, file)  
4     return pd.read_csv(path)
```

```
In [52]: 1 def export_data(data, name, pre):  
2     # split data into x and y subset  
3     x = data.drop(columns='price')  
4     y = data.price.copy()  
5     x_pre = pre.transform(x)  
6     file_name = f'{name}_pre.csv'  
7     file_path = os.path.join(PATH_DIR, file_name)  
8     x_pre.join(y).to_csv(file_path, index=False)
```

```
In [53]: 1 export_data(train, 'train', preprocessor)
```

```
In [54]: 1 test = get_data('test')  
2 export_data(test, 'test', preprocessor)
```

```
In [55]: 1 val = get_data('val')  
2 export_data(val, 'val', preprocessor)
```

Data

```
In [69]: 1 train_pre = get_data('train_pre')
```

```
In [71]: 1 val_pre = get_data('val_pre')
```

```
In [72]: 1 pre_data = {  
2     'train':train_pre,  
3     'val':val_pre  
4 }
```

Model

Train the Model

```
In [59]: 1 from xgboost import XGBRegressor
```

```
In [75]: 1 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV, train_test_split
2 from sklearn.metrics import mean_squared_error
```


In [146]:

```
1 train_pre = get_data('train_pre')
2 val_pre = get_data('val_pre')
3
4 # Combine train and validation data for cross-validation
5 full_data = pd.concat([train_pre, val_pre], axis=0)
6
7 # Split the data into features and target
8 X = full_data.drop(columns='price')
9 y = full_data['price']
10
11 # Convert the data to DMatrix format
12 dtrain = xgb.DMatrix(X, label=y)
13
14 # Define the hyperparameter grid
15 param_grid = {
16     'eta': [0.05, 0.1, 0.15, 0.2],
17     'max_depth': [3, 4, 5],
18     'alpha': [0, 0.5, 1],
19     'subsample': [0.6, 0.8, 1.0],
20     'colsample_bytree': [0.6, 0.8, 1.0]
21 }
22
23 # Initialize the model
24 model = xgb.XGBRegressor(objective="reg:squarederror", random_state=42)
25
26 # Initialize GridSearchCV or RandomizedSearchCV
27 search = GridSearchCV(
28     estimator=model,
29     param_grid=param_grid,
30     scoring='neg_mean_squared_error',
31     cv=5, # 5-fold cross-validation
32     n_jobs=-1, # Use all available cores
33     verbose=1
34 )
35
36 # Fit the search object
37 search.fit(X, y)
38
39 # Print the best hyperparameters and the corresponding score
40 print("Best hyperparameters: ", search.best_params_)
41 print("Best RMSE: ", (-search.best_score_)**0.5)
42
43 # Evaluate the best model on the validation set
```

```
44 best_model = search.best_estimator_
45 X_val = val_pre.drop(columns='price')
46 y_val = val_pre['price']
47 dval = xgb.DMatrix(X_val, label=y_val)
48 y_pred = best_model.predict(X_val)
49 rmse = mean_squared_error(y_val, y_pred, squared=False)
50 print(f'Validation RMSE: {rmse}')
51
52 # Calculate accuracy (or another metric if needed)
53 accuracy = best_model.score(X_val, y_val)
54 print(f'Validation Accuracy: {accuracy}')
55
```

```
Fitting 5 folds for each of 324 candidates, totalling 1620 fits
Best hyperparameters: {'alpha': 0, 'colsample_bytree': 0.8, 'eta': 0.2, 'max_depth': 4, 'subsample': 1.0}
Best RMSE: 2255.8550082337388
Validation RMSE: 1890.2229800805908
Validation Accuracy: 0.8395137786865234
```

Saving model and preprocessor

In [147]:

```
1 import pickle
2 model_path = 'best_model.pkl'
3 with open(model_path, 'wb') as f:
4     pickle.dump(best_model, f)
5 print(f'Model saved to {model_path}')
```

```
Model saved to best_model.pkl
```

In [148]:

```
1 preprocessor_path = 'preprocess.pkl'
2 with open(preprocessor_path, 'wb') as f:
3     pickle.dump(preprocessor, f)
```

Prediction

In [149]:

```
1 import pandas as pd
2
3 # New data point
4 data = ['Jet Airways', '2019-03-06', 'Banglore', 'New Delhi', '08:00:00', '08:15:00', 1455, 1.0, 'No Info']
5 columns= ['airline', 'date_of_journey', 'source', 'destination', 'dep_time', 'arrival_time', 'duration', 'total_
6
7 # Create a DataFrame
8 new_data = pd.DataFrame([data], columns=columns)
9 new_data
```

Out[149]:

	airline	date_of_journey	source	destination	dep_time	arrival_time	duration	total_stops	additional_info
0	Jet Airways	2019-03-06	Banglore	New Delhi	08:00:00	08:15:00	1455	1.0	No Info

In [150]:

```
1 import pickle
2
3 # Load the model and preprocessor from disk
4 with open('best_model.pkl', 'rb') as f:
5     loaded_model = pickle.load(f)
6
7 with open('preprocess.pkl', 'rb') as f:
8     loaded_preprocessor = pickle.load(f)
9
```

In [151]:

```
1 # Preprocess the new data point
2 new_data_preprocessed = loaded_preprocessor.transform(new_data)
3
4 # Make prediction
5 predicted_price = loaded_model.predict(new_data_preprocessed)
6
7 print(f'Predicted Price: {predicted_price[0]}')
```

Predicted Price: 15125.7548828125

In []: 1

In []: 1