

```
In [47]: import pandas as pd  
import numpy as np
```

```
In [2]: df = pd.read_csv('car data.csv')
```

```
In [3]: df.head()
```

```
Out[3]:   Car_Name  Year  Selling_Price  Present_Price  Kms_Driven  Fuel_Type  Seller_Type  Transmission  Owner  
0      ritz  2014          3.35         5.59     27000    Petrol    Dealer    Manual      0  
1      sx4  2013          4.75         9.54     43000   Diesel    Dealer    Manual      0  
2      ciaz  2017          7.25         9.85      6900    Petrol    Dealer    Manual      0  
3    wagon r  2011          2.85         4.15      5200    Petrol    Dealer    Manual      0  
4     swift  2014          4.60         6.87     42450   Diesel    Dealer    Manual      0
```

```
In [5]: df.shape
```

```
Out[5]: (301, 9)
```

```
In [7]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 301 entries, 0 to 300  
Data columns (total 9 columns):  
 #   Column           Non-Null Count  Dtype     
---  --  
 0   Car_Name        301 non-null    object    
 1   Year            301 non-null    int64    
 2   Selling_Price   301 non-null    float64  
 3   Present_Price   301 non-null    float64  
 4   Kms_Driven     301 non-null    int64    
 5   Fuel_Type       301 non-null    object    
 6   Seller_Type     301 non-null    object    
 7   Transmission    301 non-null    object    
 8   Owner           301 non-null    int64    
dtypes: float64(2), int64(3), object(4)  
memory usage: 21.3+ KB
```

```
In [28]: print(df['Seller_Type'].unique())  
print(df['Transmission'].unique())  
print(df['Owner'].unique())  
print(df['Fuel_Type'].unique())
```

```
['Dealer' 'Individual'  
['Manual' 'Automatic'  
[0 1 3]  
['Petrol' 'Diesel' 'CNG']
```

Check Missing or Null Values

```
In [11]: df.isnull().sum()
```

```
Out[11]: Car_Name      0  
Year          0  
Selling_Price  0  
Present_Price  0  
Kms_Driven    0  
Fuel_Type     0  
Seller_Type    0  
Transmission   0  
Owner          0  
dtype: int64
```

```
In [12]: df.describe()
```

```
Out[12]:      Year  Selling_Price  Present_Price  Kms_Driven  Owner  
count  301.000000  301.000000  301.000000  301.000000  301.000000  
mean  2013.627907  4.661296   7.628472  36947.205980  0.043189  
std   2.891554   5.082812   8.644115  38886.883882  0.247915  
min   2003.000000  0.100000  0.320000  500.000000  0.000000  
25%  2012.000000  0.900000  1.200000  15000.000000  0.000000  
50%  2014.000000  3.600000  6.400000  32000.000000  0.000000  
75%  2016.000000  6.000000  9.900000  48767.000000  0.000000  
max  2018.000000 35.000000  92.600000 500000.000000  3.000000
```

```
In [14]: df.columns
```

```
Out[14]: Index(['Car_Name', 'Year', 'Selling_Price', 'Present_Price', 'Kms_Driven',
   'Fuel_Type', 'Seller_Type', 'Transmission', 'Owner'],
   dtype='object')
```

```
In [18]: final_dataset = df[['Year', 'Selling_Price', 'Present_Price', 'Kms_Driven',
   'Fuel_Type', 'Seller_Type', 'Transmission', 'Owner']]
```

```
In [19]: final_dataset.head()
```

```
Out[19]:
```

	Year	Selling_Price	Present_Price	Kms_Driven	Fuel_Type	Seller_Type	Transmission	Owner
0	2014	3.35	5.59	27000	Petrol	Dealer	Manual	0
1	2013	4.75	9.54	43000	Diesel	Dealer	Manual	0
2	2017	7.25	9.85	6900	Petrol	Dealer	Manual	0
3	2011	2.85	4.15	5200	Petrol	Dealer	Manual	0
4	2014	4.60	6.87	42450	Diesel	Dealer	Manual	0

```
In [20]: final_dataset['Current_Year']=2023
```

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

```
Out[21]:
```

	Year	Selling_Price	Present_Price	Kms_Driven	Fuel_Type	Seller_Type	Transmission	Owner	Current_Year
0	2014	3.35	5.59	27000	Petrol	Dealer	Manual	0	2023
1	2013	4.75	9.54	43000	Diesel	Dealer	Manual	0	2023
2	2017	7.25	9.85	6900	Petrol	Dealer	Manual	0	2023
3	2011	2.85	4.15	5200	Petrol	Dealer	Manual	0	2023
4	2014	4.60	6.87	42450	Diesel	Dealer	Manual	0	2023

```
In [22]: final_dataset['no_year'] = final_dataset['Current_Year'] - final_dataset['Year']
```

```
In [23]: final_dataset.head()
```

```
Out[23]:
```

	Year	Selling_Price	Present_Price	Kms_Driven	Fuel_Type	Seller_Type	Transmission	Owner	Current_Year	no_year
0	2014	3.35	5.59	27000	Petrol	Dealer	Manual	0	2023	9
1	2013	4.75	9.54	43000	Diesel	Dealer	Manual	0	2023	10
2	2017	7.25	9.85	6900	Petrol	Dealer	Manual	0	2023	6
3	2011	2.85	4.15	5200	Petrol	Dealer	Manual	0	2023	12
4	2014	4.60	6.87	42450	Diesel	Dealer	Manual	0	2023	9

```
In [24]: final_dataset.drop(['Year'],axis=1,inplace=True)
```

```
In [26]: final_dataset.drop(['Current_Year'],axis=1,inplace=True)
```

```
In [27]: final_dataset.head()
```

```
Out[27]:
```

	Selling_Price	Present_Price	Kms_Driven	Fuel_Type	Seller_Type	Transmission	Owner	no_year
0	3.35	5.59	27000	Petrol	Dealer	Manual	0	9
1	4.75	9.54	43000	Diesel	Dealer	Manual	0	10
2	7.25	9.85	6900	Petrol	Dealer	Manual	0	6
3	2.85	4.15	5200	Petrol	Dealer	Manual	0	12
4	4.60	6.87	42450	Diesel	Dealer	Manual	0	9

```
In [29]: final_dataset = pd.get_dummies(final_dataset,drop_first=True)
```

```
In [30]: final_dataset.head()
```

```
Out[30]:
```

	Selling_Price	Present_Price	Kms_Driven	Owner	no_year	Fuel_Type_Diesel	Fuel_Type_Petrol	Seller_Type_Individual	Transmission_Manual
0	3.35	5.59	27000	0	9	0	1	0	1
1	4.75	9.54	43000	0	10	1	0	0	1
2	7.25	9.85	6900	0	6	0	1	0	1
3	2.85	4.15	5200	0	12	0	1	0	1
4	4.60	6.87	42450	0	9	1	0	0	1

In [31]: `final_dataset.corr()`

Out[31]:

	Selling_Price	Present_Price	Kms_Driven	Owner	no_year	Fuel_Type_Diesel	Fuel_Type_Petrol	Seller_Type_Individual	Transmission_Manual
Selling_Price	1.000000	0.878983	0.029187	-0.088344	-0.236141	0.552339	-0.540571	-0.550724	
Present_Price	0.878983	1.000000	0.203647	0.008057	0.047584	0.473306	-0.465244	-0.512030	
Kms_Driven	0.029187	0.203647	1.000000	0.089216	0.524342	0.172515	-0.172874	-0.101419	
Owner	-0.088344	0.008057	0.089216	1.000000	0.182104	-0.053469	0.055687	0.124269	
no_year	-0.236141	0.047584	0.524342	0.182104	1.000000	-0.064315	0.059959	0.039896	
Fuel_Type_Diesel	0.552339	0.473306	0.172515	-0.053469	-0.064315	1.000000	-0.979648	-0.350467	
Fuel_Type_Petrol	-0.540571	-0.465244	-0.172874	0.055687	0.059959	-0.979648	1.000000	0.358321	
Seller_Type_Individual	-0.550724	-0.512030	-0.101419	0.124269	0.039896	-0.350467	0.358321	1.000000	
Transmission_Manual	-0.367128	-0.348715	-0.162510	-0.050316	-0.000394	-0.098643	0.091013	0.063240	



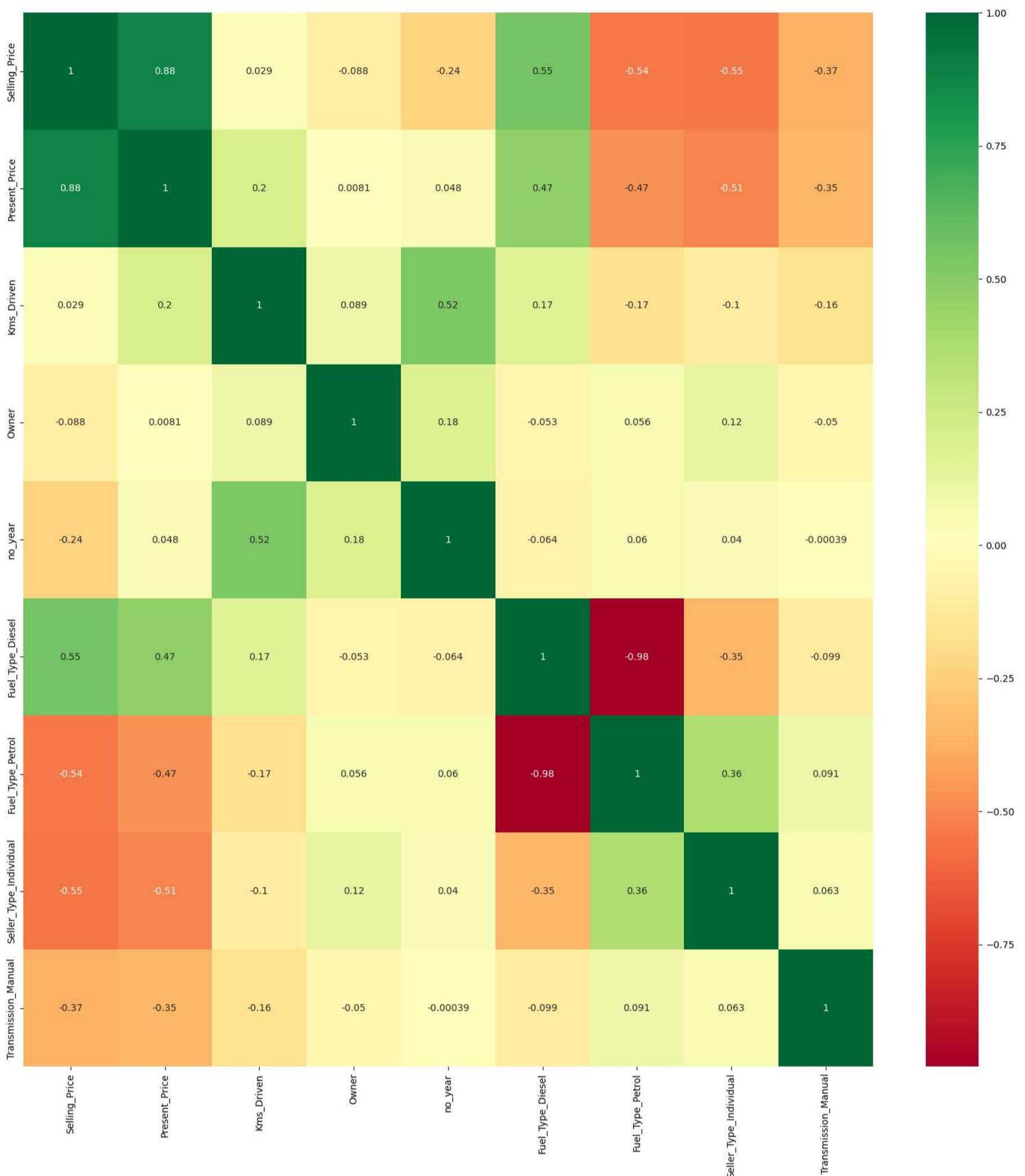
In [35]: `import matplotlib.pyplot as plt`
`%matplotlib inline`

In [36]: `corrmat = final_dataset.corr()`
`top_corr_features = corrmat.index`

```

plt.figure(figsize=(20,20))
g = sns.heatmap(final_dataset[top_corr_features].corr(), annot=True, cmap="RdYlGn")
plt.show()

```



In [37]: `final_dataset.head()`

	Selling_Price	Present_Price	Kms_Driven	Owner	no_year	Fuel_Type_Diesel	Fuel_Type_Petrol	Seller_Type_Individual	Transmission_Manual
0	3.35	5.59	27000	0	9	0	1	0	1
1	4.75	9.54	43000	0	10	1	0	0	1
2	7.25	9.85	6900	0	6	0	1	0	1
3	2.85	4.15	5200	0	12	0	1	0	1
4	4.60	6.87	42450	0	9	1	0	0	1

In [38]: `#independent and dependent features`
`x = final_dataset.iloc[:,1:] #independent feature`
`y = final_dataset.iloc[:,0] #dependent feature(selling price)`

In [39]: `x.head()`

	Present_Price	Kms_Driven	Owner	no_year	Fuel_Type_Diesel	Fuel_Type_Petrol	Seller_Type_Individual	Transmission_Manual
0	5.59	27000	0	9	0	1	0	1
1	9.54	43000	0	10	1	0	0	1
2	9.85	6900	0	6	0	1	0	1
3	4.15	5200	0	12	0	1	0	1
4	6.87	42450	0	9	1	0	0	1

In [40]: `y.head()`

```
Out[40]: 0    3.35
1    4.75
2    7.25
3    2.85
4    4.60
Name: Selling_Price, dtype: float64
```

Feature Importance

```
In [41]: from sklearn.ensemble import ExtraTreesRegressor
model=ExtraTreesRegressor()
model.fit(x,y)
```

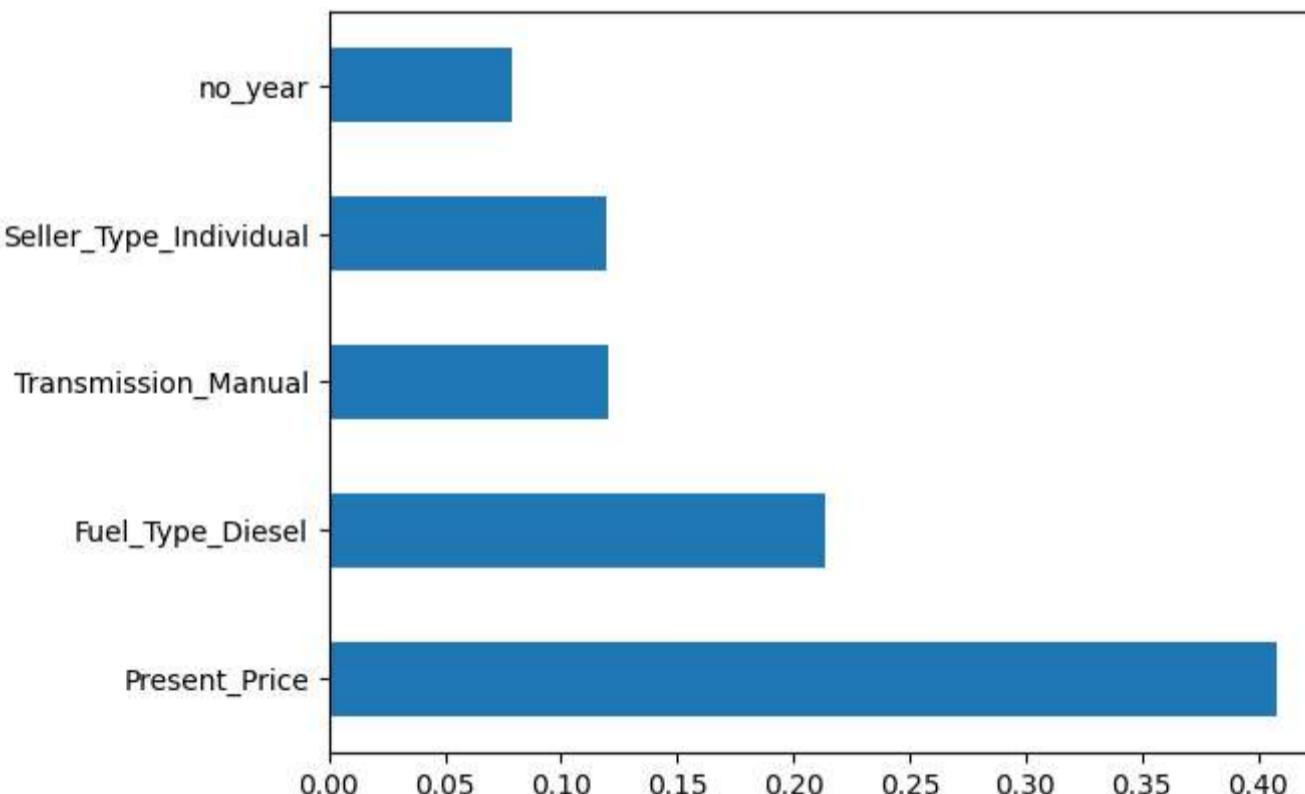
Out[41]:

```
▼ ExtraTreesRegressor
  ExtraTreesRegressor()
```

In [42]: `print(model.feature_importances_)`

```
[0.4081853  0.04053097  0.00042921  0.07918338  0.21367046  0.01849998
 0.11955993  0.11994078]
```

```
In [43]: #plot graph of feature importance for better visualization
feat_importances = pd.Series(model.feature_importances_,index=x.columns)
feat_importances.nlargest(5).plot(kind='barh')
plt.show()
```



```
In [44]: 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)
```

In [45]: `x_train.shape`

```
Out[45]: (240, 8)
```

```
In [46]: from sklearn.ensemble import RandomForestRegressor
rf_random=RandomForestRegressor()
```

```
In [51]: #Hyperparameters
#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 each split
max_features = ['auto','sqrt']
#Maximum number of Levels in tree
max_depth = [int(x) for x in np.linspace(5,30,num=6)]

#maximum number of samples required to split a node
min_samples_split = [2,5,10,15,100]
```

```
#minimum number of samples required at leaf node
min_samples_leaf = [1,2,5,10]
```

```
In [52]: from sklearn.model_selection import RandomizedSearchCV
```

```
In [53]: 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
}
print(random_grid)

{'n_estimators': [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 1200], 'max_features': ['auto', 'sqrt'], 'max_depth': [5, 10, 15, 20, 25, 30], 'min_samples_split': [2, 5, 10, 15, 100], 'min_samples_leaf': [1, 2, 5, 10]}
```

```
In [54]: #Random grid is used in order to get best hyperparameters
#First create the base model to tune
rf = RandomForestRegressor()
```

```
In [55]: rf_random = RandomizedSearchCV(estimator=rf, param_distributions=random_grid, scoring='neg_mean_squared_error', n_iter=10, cv=5)
```

```
In [56]: 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= 0.6s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time= 0.6s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time= 0.6s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time= 0.6s
[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time= 0.6s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total time= 0.8s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total time= 0.8s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total time= 0.8s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total time= 0.8s
[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total time= 1.0s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total time= 0.2s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total time= 0.2s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total time= 0.2s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total time= 0.2s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total time= 0.2s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time= 0.3s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time= 0.3s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time= 0.3s
C:\Users\shubh\AppData\Local\Programs\Python\Python310\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_features=1.0` or remove this parameter as it is also the default value for RandomForestRegressors and ExtraTreesRegressors.
  warn(
[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time= 0.3s
```

Out[56]:

- ▶ **RandomizedSearchCV**
- ▶ **estimator: RandomForestRegressor**
 - ▶ **RandomForestRegressor**

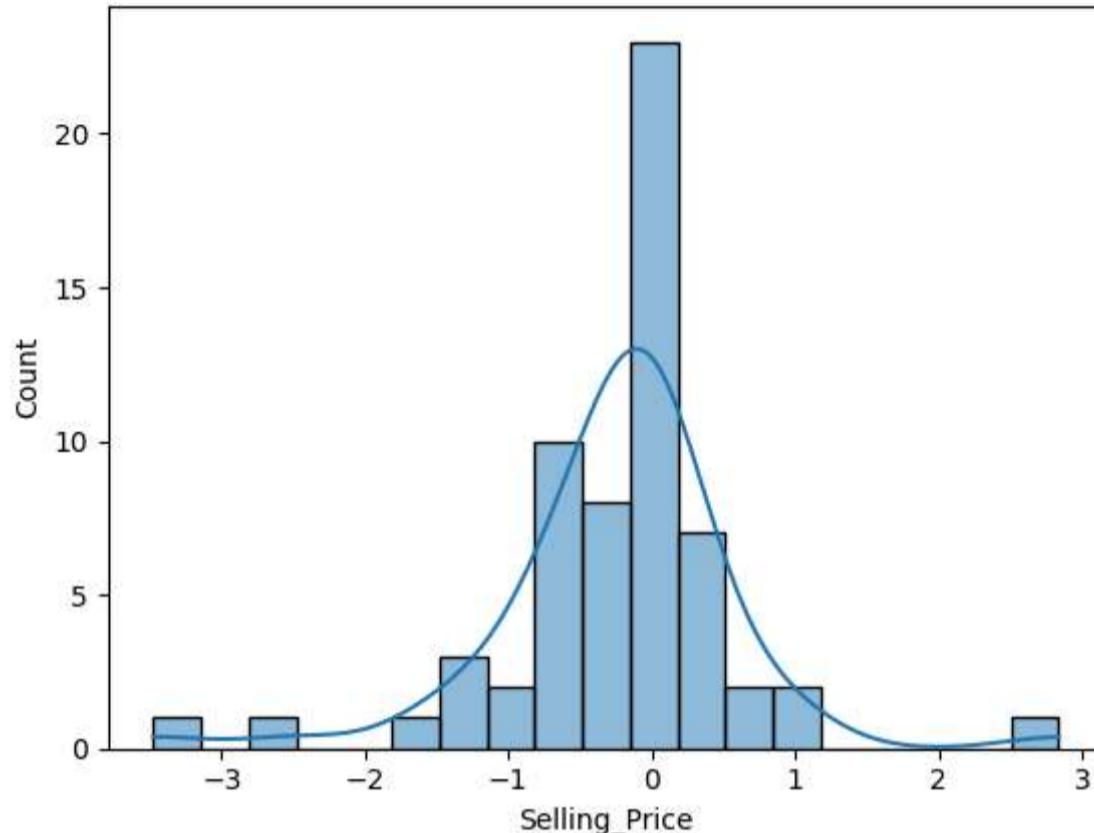
```
In [57]: predictions = rf_random.predict(x_test)
```

```
In [58]: predictions
```

```
Out[58]: array([ 8.52189,  0.50176,  5.39801,  3.37341,  0.78886,  1.39583,
  6.52147,  1.37178,  3.14813,  5.97256,  4.03216,  0.53326,
  0.32355,  10.21373,  4.04336,  5.74945,  0.43036,  3.8751 ,
  8.06685,  4.73819,  5.51812,  5.72391,  0.41593,  7.45853,
  0.67196,  0.58065,  6.99615,  4.59444,  0.96258,  2.66726,
  0.55823,  4.39465,  2.15045,  3.73452,  0.74877,  0.57005,
  0.52684,  1.04826,  3.71459,  3.9197 ,  0.85539,  2.92491,
  4.83557,  0.61142,  5.41368,  1.0781 ,  0.4116 ,  0.33602,
  1.10286,  3.49949,  5.76599,  12.07897,  4.9101 ,  0.25268,
  1.26014,  0.53182,  5.41859,  0.90635,  6.70895,  3.95881,
  0.5658 ])
```

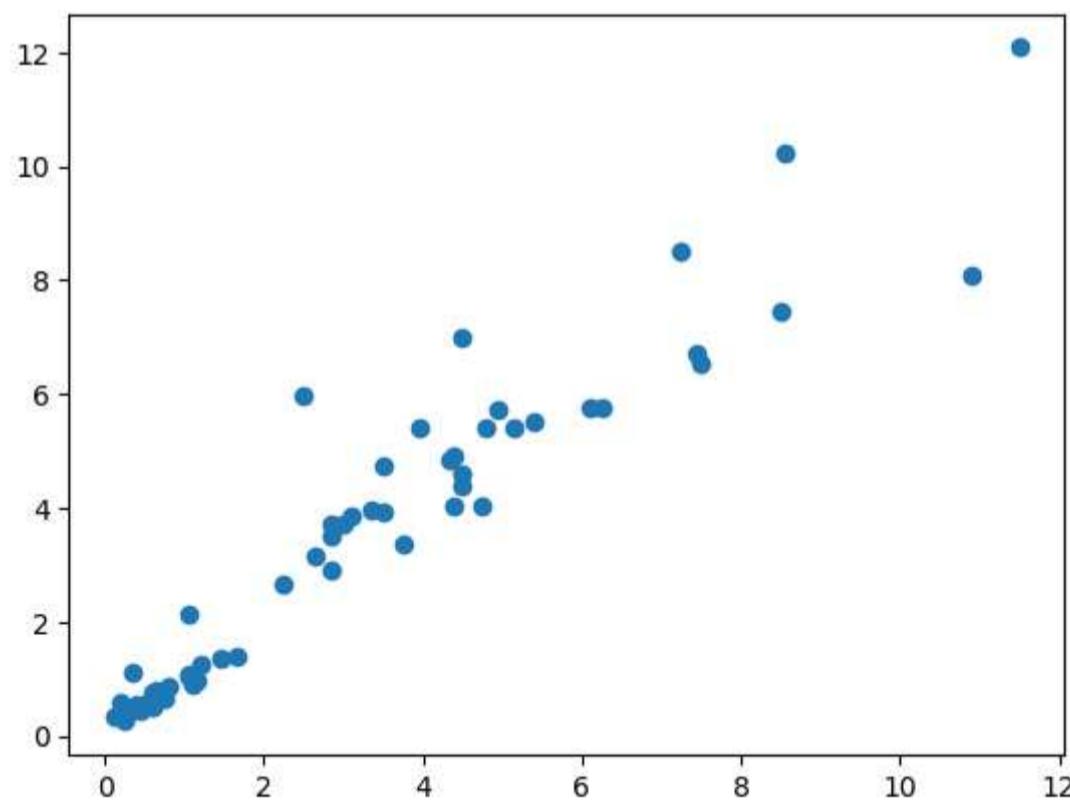
```
In [67]: sns.histplot(y_test-predictions,kde=True)
```

```
Out[67]: <AxesSubplot: xlabel='Selling_Price', ylabel='Count'>
```



```
In [64]: plt.scatter(y_test,predictions)
```

```
Out[64]: <matplotlib.collections.PathCollection at 0x2b4d0c0f9d0>
```



```
In [69]: import pickle
#open file to store data
file = open('random_forest_regression_model.pkl','wb')

#dump info to the file
pickle.dump(rf_random,file)
```

```
In [ ]:
```