

```
In [1]:
```

```
In [2]:
```

```
import library
```

```
In [3]:
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
```

Data Pre Processing

```
In [4]:
```

```
#Load data
df=pd.read_csv('C:/Users/Prachi/Desktop/ready/archive (1)/car_data.csv')
```

```
In [5]:
```

```
#Printing first five row
df.head()
```

```
Out[5]:
```

	Car_Name	Year	Selling Price	Present Price	Kms_Driven	Fuel_Type	Seller_Type	Transmission	Owner
0	Vol 2014	2014	5.75	5.54	27000	Petrol	Dealer	Manual	1
1	Vol 2013	2013	4.75	5.54	40000	Diesel	Dealer	Manual	0
2	Vol 2017	2017	7.25	5.95	4000	Petrol	Dealer	Manual	1
3	Volvo 2011	2011	3.85	4.15	3000	Petrol	Dealer	Manual	0
4	Vol 2014	2014	4.05	5.57	4000	Diesel	Dealer	Manual	1

```
In [6]:
```

```
#df.shape: number of row and column present in a dataset
```

```
Out[6]:
```

```
(281, 10)
```

```
In [7]:
```

```
df.columns: Returning name of all the columns
```

```
Out[7]:
```

```
Index(['Car_Name', 'Year', 'Selling Price', 'Present Price', 'Kms_Driven',  
       'Fuel_Type', 'Seller_Type', 'Transmission', 'Owner'],  
      dtype='object')
```

```
In [8]:
```

```
#renaming the car_name column  
#Name of company won't affect car's price, price depends upon how many year it's been used, fuel type etc.,  
#So I will drop the column car_name from original dataframe  
df.drop('Car_Name',axis=1,inplace=True)
```

```
In [9]:
```

```
df.isnull().sum() -> Are there any null value present
```

```
Out[9]:
```

```
Year          0  
Selling Price 0  
Present Price 0  
Kms_Driven    0  
Fuel_Type     0  
Seller_Type   0  
Transmission  0  
Owner         0  
dtype: int64
```

```
In [10]:
```

```
#Check the number of rows,columns and data type  
df.info()
```

```

columns: 'year', 'selling_price', 'present_price', 'kms_driven', 'fuel_type', 'seller_type', 'transmission', 'owner'
dtypes: float64(2), int64(2), object(2)
memory usage: 28.9+ KB

```

As we see there are some categorical features are present as we have to store them in one column

```

print(df['seller_type'].unique())
print(df['fuel_type'].unique())
print(df['transmission'].unique())
print(df['owner'].unique())

['Dealer', 'Individual']
['Petrol', 'Diesel', 'CNG']
['Manual', 'Automatic']
[0 1 2]

```

```
df.head()
```

	Year	Selling Price	Present Price	Kms Driven	Fuel Type	Seller Type	Transmission	Owner
0	2014	1.35	1.58	27000	Petrol	Dealer	Manual	0
1	2015	4.75	3.54	40000	Diesel	Dealer	Manual	1
2	2017	7.25	6.95	8000	Petrol	Dealer	Manual	0
3	2011	3.85	4.15	9000	Petrol	Dealer	Manual	0
4	2014	4.65	6.67	42000	Diesel	Dealer	Manual	0

year represent the year in which car have been purchased.
 And now we can calculate the number of year car has been used ?

```
df['Current_Year']=2020
```

```

df['No. of years']=df['Current_Year']-df['Year']
df=df.drop(['Current_Year','Year'],axis=1)
df.head()

```

	Selling Price	Present Price	Kms Driven	Fuel Type	Seller Type	Transmission	Owner	No. of years
0	1.35	1.58	27000	Petrol	Dealer	Manual	0	6
1	4.75	3.54	40000	Diesel	Dealer	Manual	1	7
2	7.25	6.95	8000	Petrol	Dealer	Manual	0	1
3	3.85	4.15	9000	Petrol	Dealer	Manual	0	6
4	4.65	6.67	42000	Diesel	Dealer	Manual	0	6

```

df=df.get_dummies(df.drop(['Year'],axis=1))
df.head()

```

	Selling Price	Present Price	Kms Driven	Owner	No. of years	Fuel_Type_Diesel	Fuel_Type_Petrol	Seller_Type_Individual	Transmission_Man
0	1.35	1.58	27000	0	6	0	1	0	1

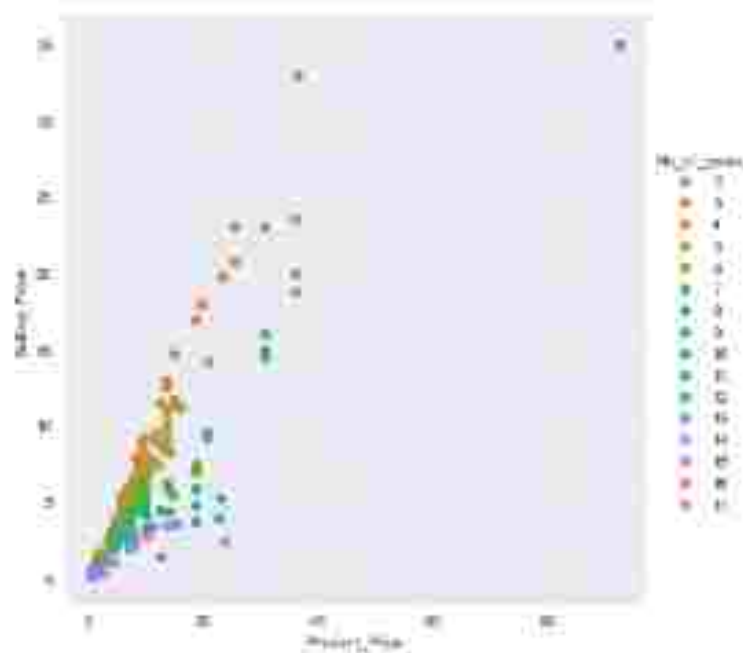
1	475	8.34	40000	0	7	1	0	0
2	720	8.88	40000	0	3	0	1	0
3	2.85	4.15	50000	0	8	0	1	0
4	8.80	8.57	40000	0	9	1	0	0

```
df.corr()
```

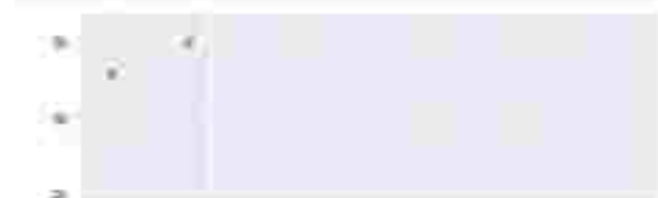
	Selling_Price	Present_Price	Kms_Driven	Owner	No_of_years	Fuel_Type_Diesel	Fuel_Type_Petrol	Seller_Type_Individual
Selling_Price	1.000000	0.978882	0.026187	-0.088849	-0.286441	0.553218	-0.540071	-0.33
Present_Price	0.978882	1.000000	0.026847	-0.090057	-0.247554	0.477038	-0.480244	-0.31
Kms_Driven	0.026187	0.026847	1.000000	0.038214	0.324340	0.172713	-0.172074	-0.19
Owner	-0.088849	-0.090057	0.038214	1.000000	0.132406	0.023489	-0.099047	0.12
No_of_years	-0.286441	-0.247554	0.324340	0.132406	1.000000	0.064013	-0.064013	0.00
Fuel_Type_Diesel	0.553218	0.477038	0.172713	0.023489	0.064013	1.000000	-0.070648	-0.33
Fuel_Type_Petrol	-0.540071	-0.480244	-0.172074	-0.099047	-0.064013	-0.070648	1.000000	0.33
Seller_Type_Individual	-0.330000	-0.310000	-0.190000	0.120000	0.000000	-0.330000	0.330000	1.00
Transaction_Mount	-0.377128	-0.340715	-0.182087	0.088016	0.000000	-0.377128	0.377128	0.00

Data Visualization

```
plt.figure(figsize=(10,8))
plt.scatter(df['No_of_years'], df['Present_Price'], c=df['Selling_Price'], s=100, alpha=0.5)
plt.xlabel('No_of_years')
```



```
plt.figure(figsize=(10,8))
plt.scatter(df['Present_Price'], df['Selling_Price'], c=df['No_of_years'], s=100, alpha=0.5)
plt.xlabel('Present_Price')
```



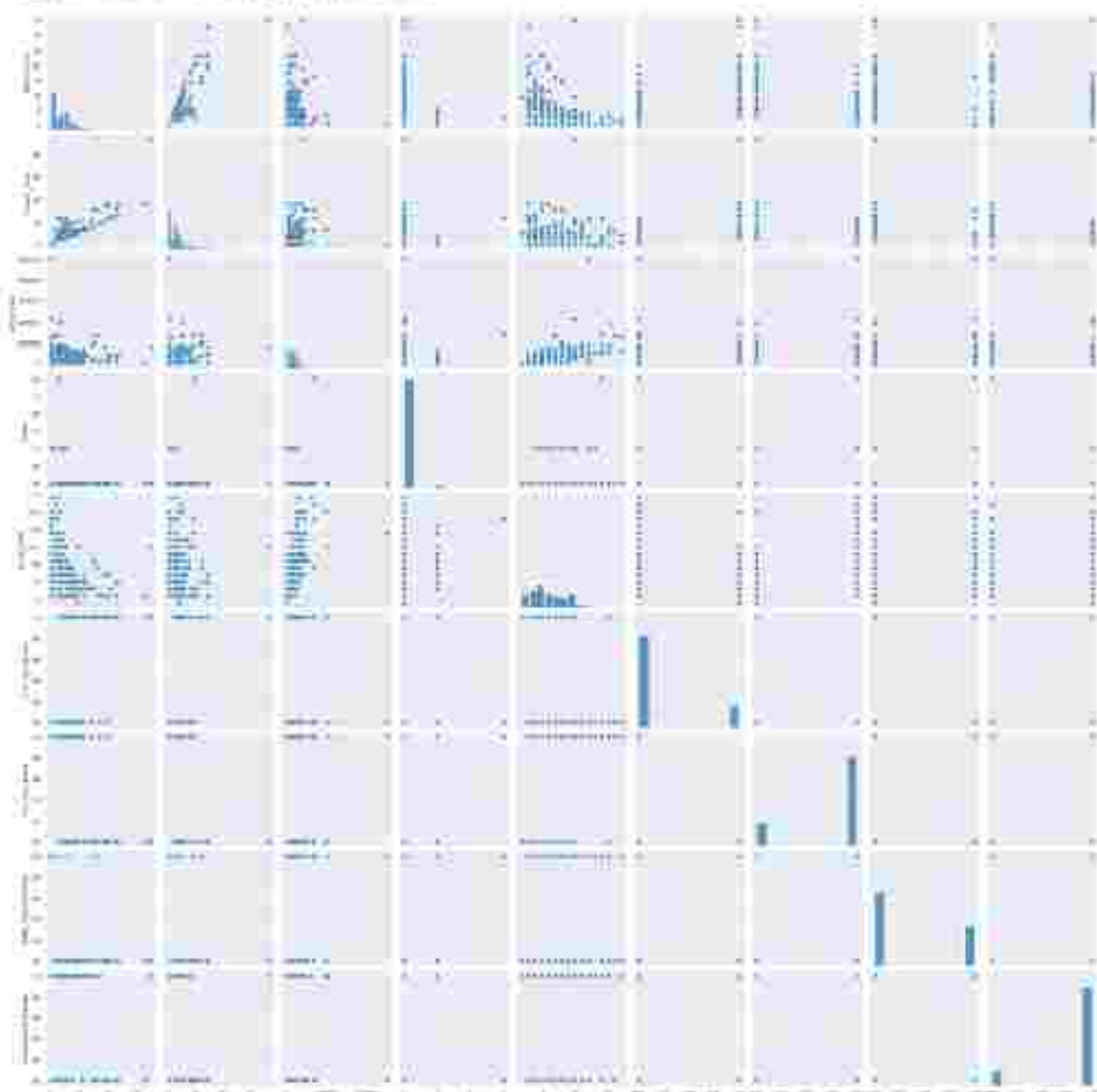


As the car would be driven higher will be
 With cost as we see the graph at 300 distance $x = 340000$ kilometers the car's cost is near to zero as we can

1.0.0.0

www.mind.org.uk (011-1) 1111

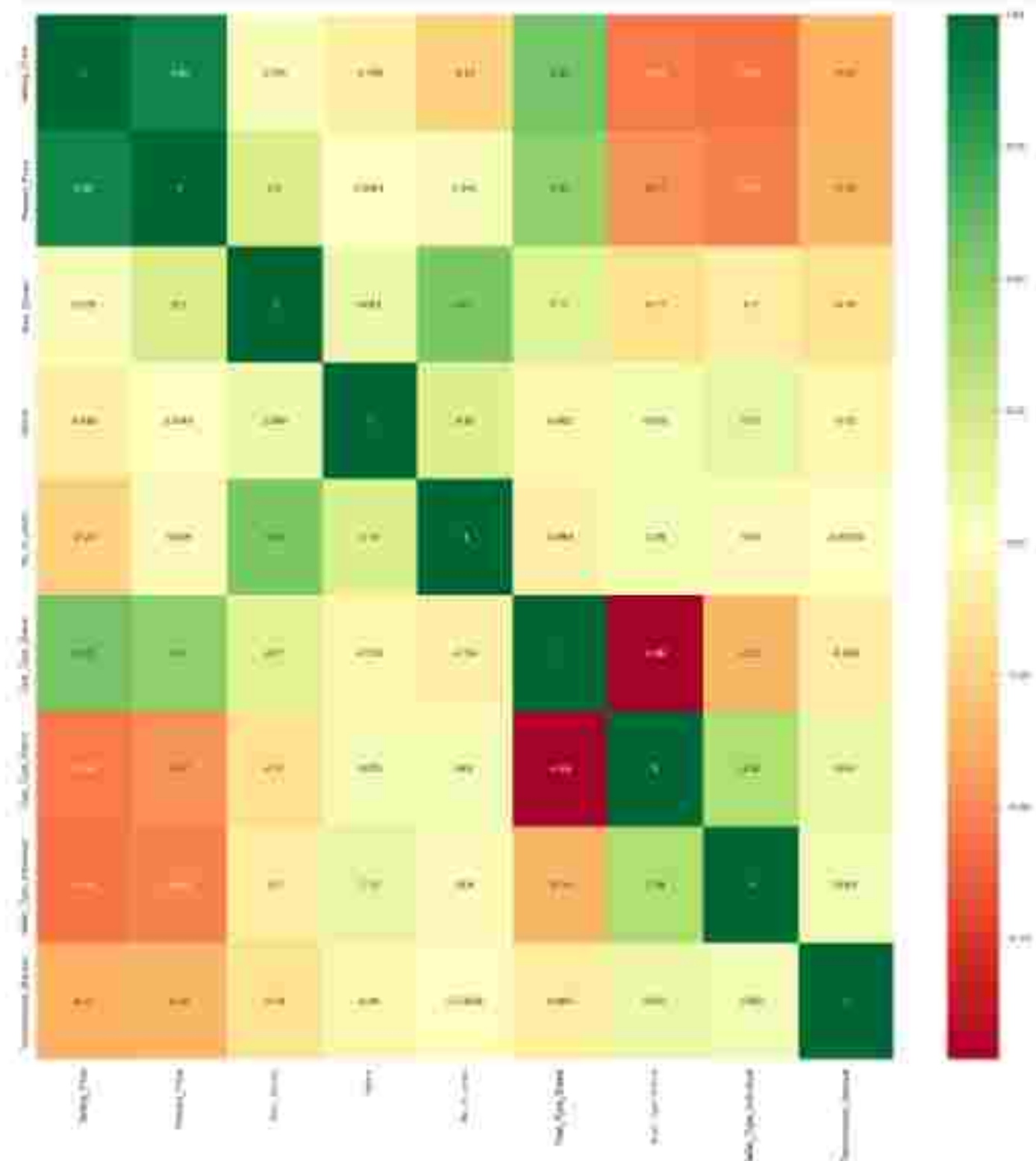
WILEY



10 | 101 | All we see there are very high overlapping in dataset as well as we cannot use for linear regression, too bad!

11 | 111 |

```
import pandas as pd
top_corr_features = extract.index
plt.figure(figsize=(20,20))
plt.imshow(
    pd.DataFrame(top_corr_features).corr().annot=True, cmap=plt.cm.
```



12 | 121 | *PCA analysis* : some analysis involve single variable most prominently it is used to find the pattern

13 | 131 |

```
pc1 = pca.components_[0]
pc2 = pca.components_[1]
```

```

19 [10]: X['Owner'].fillna(0)

20 [20]: array([0, 1, 0], dtype=int64)

```

```

21 [21]: A.head()

```

```

22 [22]:

```

	Present_Price	Kms_Driven	Owner	No_of_years	Fuel_Type_Diesel	Fuel_Type_Petrol	Seat_Type_Individual	Transmission_Manual
0	2.55	27000	1	6	1	1	1	
1	9.54	40000	1	7	1	0	1	
2	6.90	5000	0	3	0	1	1	
3	4.16	10000	0	9	0	1	0	
4	5.07	42400	0	3	1	0	1	

```

23 [23]: y.head()

```

```

24 [24]:
0    2.25
1    4.75
2    7.25
3    2.35
4    4.88
Name: Selling_Price, dtype: float64

```

```

25 [25]:
from sklearn.metrics import ExtraTreeRegressor
import matplotlib.pyplot as plt
model = ExtraTreeRegressor()
model.fit(X,y)

```

```

26 [26]: ExtraTreeRegressor()

```

```

27 [27]: print(model.feature_importances_)

```

```

[0.48963854 0.34265175 0.38942411 0.30036432 0.21644101 0.61429507
 0.22817204 0.1348836 ]

```

```

28 [28]:
#plot graph of feature importances for better visualization:
feat_importances = plt.Series(model.feature_importances_, index=X.columns)
feat_importances.nlargest(5).plot(kind='bar')
plt.show()

```



```

29 [29]:
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)

```

```

30 [30]:
from sklearn.metrics import mean_squared_error

```



```

144 regressor=RandomForestRegressor()

145 # Create the random grid
146 n_estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]
147 print(n_estimators)

148 [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 1200]

149 # Create the random grid
150 # Number of trees in random forest
151 n_estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]
152 # Number of features to consider at every split
153 max_features = ['auto', 'sqrt']
154 # Maximum number of levels in tree
155 max_depth = [int(x) for x in np.linspace(5, 20, num = 6)]
156 # min_depth.append(None)
157 # Minimum number of samples required to split a node
158 min_samples_split = [2, 5, 10, 15, 100]
159 # Minimum number of samples required at each leaf node
160 min_samples_leaf = [1, 2, 3, 20]

161 # Create the random grid
162 random_grid = {'n_estimators': n_estimators,
163                'max_features': max_features,
164                'max_depth': max_depth,
165                'min_samples_split': min_samples_split,
166                'min_samples_leaf': min_samples_leaf}

167 print(random_grid)

168 {'n_estimators': [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 1200], 'max_features': ['auto', 'sqrt'],
169  'max_depth': [5, 10, 15, 20, 25, 30], 'min_samples_split': [2, 5, 10, 15, 100], 'min_samples_leaf': [1, 2, 3, 20]}

169 # Use the random grid to search for best hyperparameters
170 # First create the base model to tune
171 rf = RandomForestRegressor()

172 # Random search of parameters, using 5 fold cross validation
173 # Search across 100 different combinations
174 rf_random = RandomizedSearchCV(estimator = rf, param_distributions = random_grid, scoring = 'neg_mean_squared_error',
175                               cv = 5, n_iter = 100)

176 rf_random.fit(X_train, y_train)

177 Fitting 5 folds for each of 100 candidates, totalling 500 fits
178 [CV] 200 max_depth=10, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=200, total time=
179 1.8s
179 [CV] 200 max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=200, total time=
180 1.8s
180 [CV] 200 max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=200, total time=
181 0.9s
181 [CV] 200 max_depth=10, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=200, total time=
182 0.9s
182 [CV] 200 max_depth=10, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=200, total time=
183 0.9s
183 [CV] 200 max_depth=15, max_features=sqrt, min_samples_leaf=5, min_samples_split=10, n_estimators=2100, total time=
184 1.1s
184 [CV] 200 max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=2100, total time=
185 1.1s
185 [CV] 200 max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=2100, total time=
186 1.1s
186 [CV] 200 max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=2100, total time=
187 1.1s
187 [CV] 200 max_depth=15, max_features=auto, min_samples_leaf=3, min_samples_split=100, n_estimators=280, total time=
188 0.2s
188 [CV] 200 max_depth=15, max_features=auto, min_samples_leaf=2, min_samples_split=100, n_estimators=280, total time=
189 0.2s
189 [CV] 200 max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300, total time=
190 0.2s

```



```
In [57]: rf_random_test_patho_
```

```
In [62]: {'n_estimators': 1000,  
         'min_sample_split': 2,  
         'min_sample_leaf': 1,  
         'max_features': 'sqrt',  
         'max_depth': 25}
```

```
In [63]: rf_random_test_score_
```

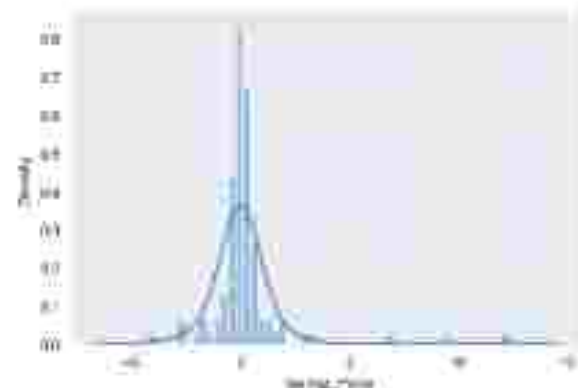
```
Out[62]: 0.86296606471117
```

```
In [64]: gradient_boost_random_predict(X_test)
```

```
In [66]: plt.display(y_test_predictions)
```

Warning: plt.show() is deprecated and will be removed in a future version. Please adapt your code to use either 'display' (a plt-level function) with similar flexibility or 'show' (an axes-level function for histograms).

```
Out[66]: AxesSubplot(= <matplotlib.figure.Figure>)
```



```
In [68]: plt.scatter(y_test_predictions)
```

```
Out[68]: matplotlib.collections.PathCollection at 0x101bce97000
```



```
In [69]: from sklearn import metrics
```

```
In [69]: print MAE: metrics.mean_absolute_error(y_test, predictions)  
print MSE: metrics.mean_squared_error(y_test, predictions)  
print RMSE: np.sqrt(metrics.mean_squared_error(y_test, predictions))  
  
MAE: 0.1881278752720218  
MSE: 4.133937830837505  
RMSE: 2.033211773338484
```

```
In [101]: import pickle
# open a file, where you want to store the data
file = open('random_forest_regression_model.pkl', 'wb')

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

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

In []:

Learning Python Data Science Tools and Libraries