

```
jupyter Untitled Last Checkpoint: 13 hours ago (autosaved) Logout
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)
In [102]: X=final_dataset.iloc[:,1:]
          y=final_dataset.iloc[:,0]

In [137]: X['Owner'].unique()
Out[137]: array([0, 1, 3], dtype=int64)

In [103]: X.head()
Out[103]:
   Present_Price  Kms_Driven  Owner  no_year  Fuel_Type_Diesel  Fuel_Type_Petrol  Seller_Type_Individual  Transmission_Manual
0             5.59       27000      0         6                0                1                0                1
1             9.54       43000      0         7                1                0                0                1
2             9.85        6900      0         3                0                1                0                1
3             4.15       5200      0         9                0                1                0                1
4             6.87      42450      0         6                1                0                0                1

In [104]: y.head()
Out[104]:
0    3.35
1    4.75
2    7.25
3    2.85
4    4.60
Name: Selling_Price, dtype: float64

In [105]: ### Feature Importance
          from sklearn.ensemble import ExtraTreesRegressor
          import matplotlib.pyplot as plt
```

```
jupyter Untitled Last Checkpoint: 13 hours ago (autosaved) Logout
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)
Out[105]: ExtraTreesRegressor()

In [106]: print(model.feature_importances_)
[0.37587921 0.03735726 0.00097047 0.08064828 0.22844919 0.00780328
 0.13445871 0.13443361]

In [108]: #Plot graph of feature importances for better visualization
          feat_importances = pd.Series(model.feature_importances_, index=X.columns)
          feat_importances.nlargest(5).plot(kind='barh')
          plt.show()

          no_year
          Transmission_Manual
          Seller_Type_Individual
          Fuel_Type_Diesel
          Present_Price

          0.00 0.05 0.10 0.15 0.20 0.25 0.30 0.35

In [120]: 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)

In [109]: from sklearn.ensemble import RandomForestRegressor

In [111]: regressor=RandomForestRegressor()
```

```
jupyter Untitled Last Checkpoint: 13 hours ago (autosaved) Logout
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)
In [112]: n_estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]
          print(n_estimators)
          [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 1200]

In [113]: from sklearn.model_selection import RandomizedSearchCV

In [115]: #Randomized Search CVsam
          # 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)]
          # max_depth.append(None)
          # 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 [116]: # 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}

          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]}
```

jupyter Untitled Last Checkpoint: 13 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)

```

In [117]: # Use the random grid to search for best hyperparameters
# First create the base model to tune
rf = RandomForestRegressor()

In [124]: # Random search of parameters, using 3 fold cross validation,
# search across 100 different combinations
rf_random = RandomizedSearchCV(estimator = rf, param_distributions = random_grid,scoring='neg_mean_squared_error', n_iter = 10, <
<

In [125]: rf_random.fit(X_train,y_train)
[CV] n_estimators=1100, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, max_depth=15, total= 1.5s
[CV] n_estimators=1100, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, max_depth=15, total= 1.6s
[CV] n_estimators=1100, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, max_depth=15, total= 1.9s
[CV] n_estimators=1100, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, max_depth=15, total= 1.6s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=300, min_samples_split=100, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.4s
[CV] n_estimators=400, min_samples_split=5, min_samples_leaf=5, max_features=auto, max_depth=15, total= 0.6s

In [126]: rf_random.best_params_
Out[126]: {'n_estimators': 1000,

```

Windows Ink Workspace

jupyter Untitled Last Checkpoint: 13 hours ago (autosaved) Logout

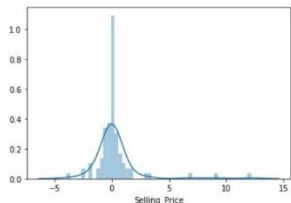
File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)

```

In [128]: predictions=rf_random.predict(X_test)

In [129]: sns.distplot(y_test-predictions)
Out[129]: <matplotlib.axes._subplots.AxesSubplot at 0x2dbf66939c8>


```



```

In [131]: plt.scatter(y_test,predictions)
Out[131]: <matplotlib.collections.PathCollection at 0x2dbf6610708>

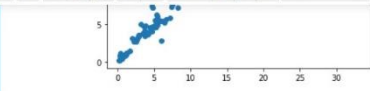
```



Windows Ink Workspace

jupyter Untitled Last Checkpoint: 13 hours ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)



```

In [133]: from sklearn import metrics

In [135]: 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: 0.8849978021977988
MSE: 3.9544237722813156
RMSE: 1.9885733007061408

In [136]: import pickle
# open a file, where you ant to store the data
file = open('random_forest_regression_model.pkl', 'wb')

# dump information to that file
pickle.dump(rf_random, file)

In [ ]:

In [ ]:

```

Windows Ink Workspace