import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import numpy as np from sklearn.model selection import train test split from sklearn.linear model import LinearRegression from sklearn.linear model import Lasso from sklearn import metrics import warnings warnings.filterwarnings('ignore') data = pd.read excel('CAR DETAILS FROM CAR DEKHO.csv.xlsx') In [4]: data.head() Out[4]: owner name year selling_price km_driven fuel seller_type transmission 0 Maruti 800 AC 2007 60000 Individual First Owner 70000 Petrol Manual Maruti Wagon R LXI Minor 2007 135000 50000 Petrol Individual Manual First Owner 2 Hyundai Verna 1.6 SX 2012 600000 100000 Diesel Individual Manual First Owner 3 Datsun RediGO T Option 2017 250000 46000 Petrol Individual First Owner Manual Honda Amaze VX i-DTEC 2014 450000 141000 Diesel Individual Manual Second Owner data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 4340 entries, 0 to 4339 Data columns (total 8 columns): Non-Null Count Dtype # Column 0 name 4340 non-null object 4340 non-null int64 1 year 2 selling_price 4340 non-null int64 3 km_driven 4340 non-null int64 4340 non-null object fuel 5 seller_type 4340 non-null object transmission 4340 non-null object 4340 non-null object dtypes: int64(3), object(5) memory usage: 271.4+ KB data.describe() selling_price km_driven year **count** 4340.000000 4.340000e+03 4340.000000 mean 2013.090783 5.041273e+05 66215.777419 4.215344 5.785487e+05 46644.102194 std min 1992.000000 2.000000e+04 1.000000 2011.000000 2.087498e+05 35000.000000 25% 2014.000000 3.500000e+05 60000.000000 2016.000000 6.000000e+05 90000.000000 max 2020.000000 8.900000e+06 806599.000000 sns.pairplot(data) <seaborn.axisgrid.PairGrid at 0x2db09084280> 2020 2015 2010 2005 2000 1995 le6 8 selling price 2 800000 600000 km driven 400000 200000 2000 2010 2020 250000 500000 750000 0.0 7.5 0 km_driven year selling_price In [8]: data.hist(bins=50, figsize=(30,15)) plt.show() data.isnull().sum() Out[9]: name selling_price km driven seller type 0 0 transmission owner dtype: int64 data.nunique() 1491 Out[10]: name 27 year 445 selling_price 770 km driven fuel 5 seller type 3 2 transmission owner dtype: int64 data.shape (4340, 8)data.corr() km_driven year selling_price 1.000000 0.413922 -0.419688 year selling_price 0.413922 1.000000 -0.192289 **km_driven** -0.419688 -0.192289 1.000000 #how old the car is? data['Car_Age'] = 2021 - data.year In [14]: sns.boxplot(x=data['selling price']) Out[14]: <AxesSubplot:xlabel='selling price'> selling_price le6 sns.boxplot(x=data['km driven']) Out[15]: <AxesSubplot:xlabel='km driven'> 0 100000200000300000400000500000600000700000800000 km_driven sns.boxplot(x=data['Car_Age']) Out[16]: <AxesSubplot:xlabel='Car_Age'> 10 15 20 Car_Age def cap data(data2): for col in data2.columns: print("capping the ",col) if (((data2[col].dtype) == 'float64') | ((data2[col].dtype) == 'int64')): Q1 = data2.quantile(0.25)Q3 = data2.quantile(0.75)IQR = Q3 - Q1data2 = data2[~((data2 < (Q1 - 1.5 * IQR)) | (data2 > (Q3 + 1.5 * IQR))).any(axis=1)]else: data2[col]=data2[col] return data2 data=cap_data(data) capping the name capping the year capping the selling_price capping the km driven capping the fuel capping the seller type capping the transmission capping the owner capping the Car Age sns.boxplot(x=data['Car Age']) Out[18]: <AxesSubplot:xlabel='Car Age'> 10 14 16 Car_Age #one hot encoding for categorical data data = pd.get dummies(data, drop first=True) data.head() name_BMW name_Ambassador name_BMW name_BMW name_Audi name_Audi ... fuel_Electric year selling_price km_driven Car_Age Grand 1800 ISZ **5 Series** 7 Series X1 **A6 2.7 TDI A6 2.8 FSI** 730Ld **MPFI PW CL** 530i sDrive20d **0** 2007 60000 70000 0 0 0 0 0 0 14 0 2007 135000 50000 0 0 0 0 0 0 0 14 0 0 0 0 ... **2** 2012 600000 100000 9 0 0 0 **3** 2017 0 0 0 0 250000 46000 0 0 0 0 ... 0 0 0 0 0 0 2014 450000 141000 7 5 rows × 1308 columns data.corr() name_BMW name_Ambassador name_BMW name_Audi name_Audi Grand 1800 ISZ **5 Series** 7 Series year selling_price km_driven Car_Age A6 2.7 TDI **A6 2.8 FSI** MPFI PW CL 530i 730Ld sD1.000000 -0.493062 -1.000000 -0.005392 -0.014060 -0.022727 -0.031395 -0.044405 -0 0.639216 year 0.002585 0.031548 selling_price 0.639216 1.000000 -0.282218 -0.639216 0.017756 0.006033 0.064129 0 -0.005752 **km_driven** -0.493062 -0.007540 0.003188 0.021068 -0.020780 0 -0.282218 1.000000 0.493062 **Car_Age** -1.000000 -0.639216 0.493062 1.000000 0.005392 0.014060 0.022727 0.031395 0.044405 0 name_Ambassador **Grand 1800 ISZ MPFI** -0.005392 0.005392 1.000000 -0.000262 -0.000262 -0.000262 -0.000371 -0 0.002585 -0.005752 **PW CL** transmission_Manual -0.087046 -0.179249 0.087046 -0.063487 -0.063487 -0.089796 -0 0.116311 0.004132 -0.063487 owner_Fourth & -0.002048 -0.002048 -0.131750 -0.108576 0.118843 0.131750 -0.002048 -0.002048 -0.002897 -0 **Above Owner** owner_Second -0.300481 0.027128 -0.013677 0 -0.216583 0.243864 0.300481 -0.009670 -0.009670 0.027128 Owner -0.087382 -0.000871 owner_Test Drive Car 0.092001 0.093820 -0.092001 -0.000871 -0.000871 -0.000871 -0.001232 -0 owner_Third Owner -0.233022 -0.177025 0.196796 0.233022 -0.004263 -0.004263 0.061541 -0.004263 -0.006029 -0 1308 rows × 1308 columns corelation=data.corr() sns.heatmap(corelation,xticklabels=corelation.columns,yticklabels=corelation.columns,annot=True) Out[6]: <AxesSubplot:> - 1.0 - 0.8 0.41 -0.42 1 - 0.6 0.4 0.41 1 -0.19selling price 0.2 0.0 -0.42-0.19driven -0.2- -0.4 Ê selling_price year km driven import numpy as np data['year']=pd.cut(data['year'],bins=[0,2.0,4.0,6.0,np.inf],labels=[1,2,3,4]) data['year'].hist() Out[12]: <AxesSubplot:> 4000 3000 2000 1000 0 3.6 3.8 4.0 4.2 4.4 data=pd.get dummies(data,drop first=True) data.head() name_Audi A4 2.0 TDI name_Audi name_Ambassador name_Audi name_Audi name_Ambassador name_Audi name_Audi km_driven transmission Grand 1800 ISZ 177 Bhp **A4 3.0 TDI A4 30 TFSI** A4 35 TDI .. Classic 2000 Dsz **A4 1.8 TFSI A4 2.0 TDI MPFI PW CL Premium** Quattro **Technology Premium** Plus 0 .. 0 0 70000 0 0 0 0 0 0 0 0 .. 1 50000 0 0 0 0 0 0 0 0 2 100000 0 0 0 0 0 0 0 0 0 .. 3 46000 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 4 141000 5 rows × 1511 columns #independent and dependent feature X=data.iloc[:,1:] Y=data.iloc[:,0] X.head() name_Audi name_Audi name_Ambassador A4 2.0 TDI name_Audi name_Audi name_Audi **A4 35 TDI** name_Ambassador name_Audi name_Audi transmission **Grand 1800 ISZ** 177 Bhp A4 3.0 TDI **A4 30 TFSI A4 35 TDI** A4 2.0 TDI Classic 2000 Dsz **A4 1.8 TFSI** Premium **MPFI PW CL** Premium **Quattro Technology Premium** Plus Plus 0 2 0 0 0 0 0 0 0 0 0 0 3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 4 5 rows × 1510 columns Y.head() 70000 50000 2 100000 3 46000 4 141000 Name: km_driven, dtype: int64 ### importance Feature from sklearn.ensemble import ExtraTreesRegressor model= ExtraTreesRegressor() model.fit(X,Y) ExtraTreesRegressor() feat_features=pd.Series(model.feature_importances_, index= X.columns) feat features.nlargest(5).plot(kind='barh') plt.show Out[92]: <function matplotlib.pyplot.show(close=None, block=None)> owner_Third Owner owner_Second Owner seller_type_1 name_Maruti Swift VXI BSIII fuel 1 0.02 0.04 0.06 0.08 0.00 0.10 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 [94]: X_train.head() Out[94]: name_Audi name_Auc name Audi name_Audi name_Ambassador A4 2.0 TDI name_Audi name_Ambassador A4 35 TC name_Audi name_Audi transmission Grand 1800 ISZ 177 Bhp A4 3.0 TDI **A4 30 TFSI** A4 35 TDI Classic 2000 Dsz **A4 1.8 TFSI** A4 2.0 TDI **Premiur MPFI PW CL Premium Quattro Technology Premium** Plu **Plus** 4250 0 0 0 0 0 0 0 0 0 3995 0 0 0 0 0 0 0 0 643 0 0 0 0 0 0 0 0 0 2537 0 0 0 0 0 0 0 0 59 0 0 0 0 0 0 0 0 0 5 rows × 1510 columns from sklearn.ensemble import RandomForestRegressor re random=RandomForestRegressor() n estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]max_features = ['auto', 'sqrt'] max depth = [int(x) for x in np.linspace(5, 30, num = 6)] $min_samples_split = [2, 5, 10, 15, 100]$ min samples leaf = [1, 2, 5, 10]from sklearn.model_selection import RandomizedSearchCV # 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', 'sqr t'], 'max_depth': [5, 10, 15, 20, 25, 30], 'min_samples_split': [2, 5, 10, 15, 100], 'min_samples_leaf': [1, 2, 5, 10]} rf = RandomForestRegressor() rf_random=RandomizedSearchCV(estimator= rf , param_distributions=random_grid, scoring='neg_mean_squared_error',n_iter= 10, cv= 5, verbose=2,random_state=42,n_jc 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 = 4.3s[CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time [CV] END max_depth=10, max_features=sqrt, min_samples_leaf=5, min_samples_split=5, n_estimators=900; total time = 6.1s [CV] END max depth=10, max features=sqrt, min samples leaf=5, min samples split=5, n estimators=900; total time = 4.1s[CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti 10.9s [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti [CV] END max depth=15, max features=sqrt, min samples leaf=2, min samples split=10, n estimators=1100; total ti [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=1100; total ti me= 7.9s [CV] END max depth=15, max features=sqrt, min samples leaf=2, min samples split=10, n estimators=1100; total ti [CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti me = 17.3s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti me=17.6s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti me=17.7s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti 17.7s [CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=100, n_estimators=300; total ti me=17.6s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time = 23.6s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time = 24.3s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time = 24.6s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time = 24.1s[CV] END max_depth=15, max_features=auto, min_samples_leaf=5, min_samples_split=5, n_estimators=400; total time = 24.9s[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim e = 33.9s[CV] END max depth=20, max features=auto, min samples leaf=10, min samples split=5, n estimators=700; total tim e = 34.6s[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim e = 34.4s[CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim [CV] END max_depth=20, max_features=auto, min_samples_leaf=10, min_samples_split=5, n_estimators=700; total tim e = 34.6s[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim [CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim e = 15.1s[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim e = 15.4s[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim e = 15.2s[CV] END max_depth=25, max_features=sqrt, min_samples_leaf=1, min_samples_split=2, n_estimators=1000; total tim [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti me= 2.2s [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=10, min_samples_split=15, n_estimators=1100; total ti [CV] END max depth=15, max features=sqrt, min samples leaf=1, min samples split=15, n estimators=300; total tim [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim 2.4s [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim 2.3s [CV] END max_depth=15, max_features=sqrt, min_samples_leaf=1, min_samples_split=15, n_estimators=300; total tim [CV] END max depth=15, max features=sqrt, min samples leaf=1, min samples split=15, n estimators=300; total tim e=2.2s[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time = 2.4s[CV] END max depth=5, max features=sqrt, min samples leaf=2, min samples split=10, n estimators=700; total time [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time [CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time = 2.4s[CV] END max_depth=5, max_features=sqrt, min_samples_leaf=2, min_samples_split=10, n_estimators=700; total time = 2.4s[CV] END max depth=20, max features=auto, min samples leaf=1, min samples split=15, n estimators=700; total tim e = 58.5s[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim [CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim e= 1.0min[CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim e= 1.0min [CV] END max_depth=20, max_features=auto, min_samples_leaf=1, min_samples_split=15, n_estimators=700; total tim e= 1.0min Out[101... RandomizedSearchCV(cv=5, estimator=RandomForestRegressor(), n_jobs=1, param_distributions={'max_depth': [5, 10, 15, 20, 25, 30], 'max_features': ['auto', 'sqrt'], 'min_samples_leaf': [1, 2, 5, 10], 'min_samples_split': [2, 5, 10, 15, 100], 'n_estimators': [100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 1100, 1200]}, random_state=42, scoring='neg_mean_squared_error', predictions=rf_random.predict(X_test) In [104... sns.distplot(Y_test-predictions) Out[104... <AxesSubplot:xlabel='km_driven', ylabel='Density'> 1.4 1.2 1.0 0.0 0.0 0.4 0.2 0.0 -200000 100000 -100000200000 300000 km_driven plt.scatter(Y test, predictions) <matplotlib.collections.PathCollection at 0x29f22fe05e0> 200000 150000 100000 50000 50000 100000 150000 200000 250000 300000 350000 from sklearn import metrics print('MAE:', metrics.mean_absolute_error(Y_test, predictions)/10000) print('RMSE:', np.sqrt(metrics.mean_squared_error(Y_test, predictions))/10000) MAE: 2.73215213957661 RMSE: 3.8521705952256498 import pickle file = open('random forest regression model.csv', 'wb') pickle.dump(rf random, file)