#import libraries import pandas as pd import numpy as np import os import matplotlib.pyplot as plt from sklearn.preprocessing import LabelEncoder, StandardScaler %matplotlib inline import warnings warnings.filterwarnings('ignore') os.chdir('C:\\Users\\utpala mohapatra\\Documents\\python folder\\python datasets') #load the dataset train df = pd.read csv('mercedes train.csv') train df.head() y X0 X1 X2 X3 X4 X5 X6 X8 ... X375 X376 X377 X378 X379 X380 X382 X383 X384 X385 0 130.81 0 0 1 0 0 0 0 0 0 0 at d 0 88.53 0 0 0 0 0 0 6 0 0 0 av d 0 0 0 0 0 0 0 0 7 76.26 0 1 az W 0 0 0 0 80.62 az 0 0 0 0 0 0 0 0 0 0 0 0 d 0 13 78.02 d az 5 rows × 378 columns # Select the independent and Dependent Variables X_train = train_df.drop(['ID','y'],axis=1) X train.head(10) X380 X382 X383 X384 X385 X0 X1 X2 X3 X4 X5 X6 X8 X10 X11 ... X375 X376 X377 X378 X379 0 ... 0 0 0 0 0 0 0 0 0 0 d 1 at u 0 0 ... 0 0 0 0 0 0 0 0 0 av d 0 ... 0 0 0 0 0 0 0 0 2 0 0 az d 1 0 ... 0 0 0 0 0 0 0 0 0 0 0 az n d е d h d 0 0 ... 0 0 0 0 0 0 0 0 0 0 az 0 ... 0 0 0 0 h 0 0 0 0 0 b d 6 f 0 0 ... 0 0 0 0 0 0 0 0 0 0 d h al 0 ... 0 0 0 0 0 as 0 0 0 0 0 0 0 0 ... 0 0 0 0 0 0 0 1 0 0 d h f 0 ... 0 0 0 0 0 b aq а 10 rows × 376 columns In [4]: y_train = train_df['y'] y_train 130.81 Out[4]: 0 88.53 1 76.26 3 80.62 78.02 4 . . . 4204 107.39 4205 108.77 4206 109.22 4207 87.48 4208 110.85 Name: y, Length: 4209, dtype: float64 Task-1> Check for null and unique values for test and train sets. # check for null values X train.isna().sum().sum() # no null values Out[5]: 0 #check for number of unique values in each column X train.nunique(axis=0) Out[6]: X0 47 27 X1 X2 ХЗ 7 X4 4 X380 X382 X383 X384 X385 Length: 376, dtype: int64 Task-2> If for any column(s), the variance is equal to zero, then you need to remove those variable(s). # Drop the columns with zero variance # first We have to remove the categorical variables for applying VarianceThreshold X_train_num = X_train.iloc[:,8:] X_train_num.shape Out[7]: (4209, 368) from sklearn.feature_selection import VarianceThreshold var thres = VarianceThreshold(threshold=0) var_thres.fit(X_train_num) Out[8]: VarianceThreshold(threshold=0) In [9]: const_var = [col for col in X_train_num.columns if col not in X_train_num.columns[var_thres.get_support()]] print(len(const_var)) X train num = X train num.drop(const var,axis=1) X train num X10 X12 X13 X14 X15 X16 X17 X18 X19 X20 ... X375 X376 X377 X378 X379 X380 X382 X383 X384 X385 0 0 0 0 0 0 0 0 0 0 0 0 4 0 0 0 0 0 0 0 0 0 0 4204 0 0 0 0 0 4205 0 4206 0 4207 0 4208 0 0 0 0 4209 rows × 356 columns # now select the categirical data and check for constant value X train cat = X train.iloc[:,:8] X train cat X0 X1 X2 X3 X4 X5 X6 X8 at а d u j av az d d az az 4204 as d aa q 4205 d aa 4206 d aa е 4207 4208 d aa ae 4209 rows × 8 columns # check whether variable X4 is constant by checking unique values X train cat['X4'].unique() Out[12]: array(['d', 'b', 'c', 'a'], dtype=object) Task-3> Apply label encoder. #Convert the categorical variable to numeric by Label encoder le = LabelEncoder() for i in X_train_cat.columns: X_train_cat[i]=le.fit_transform(X_train_cat[i]) X_train_cat.head() X0 X1 X2 X3 X4 X5 X6 X8 **0** 32 23 17 0 3 24 9 14 21 19 3 28 32 11 14 20 24 34 2 3 27 9 23 20 21 3 27 34 **4** 20 23 34 5 3 12 3 13 In [14]: X_train_clean = pd.concat([X_train_cat, X_train_num], axis=1) X_train_clean.head() X0 X1 X2 X3 X4 X5 X6 X8 X10 X12 **0** 32 23 17 0 3 24 9 14 0 ... 0 1 0 0 0 0 0 0 0 **1** 32 21 19 3 28 11 14 **2** 20 24 34 2 3 27 9 23 0 0 0 0 0 0 0 1 0 0 0 20 21 34 3 27 11 0 0 0 0 20 23 34 5 3 12 3 13 0 ... 0 0 0 0 0 0 5 rows × 364 columns # scaling the data from sklearn.preprocessing import StandardScaler scaler = StandardScaler() X_train_scaled = scaler.fit_transform(X_train_clean) Task-4> Perform dimensionality reduction. # Transform the Data to Principal Components with 95% variance explained from sklearn.decomposition import PCA pca = PCA(n components = 0.95)X train pca = pca.fit transform(X train scaled) print(X_train_scaled.shape, X_train_pca.shape) (4209, 364) (4209, 148) pca.components Out[18]: array([[-5.35850147e-02, 3.02768146e-02, 7.67089079e-02, ..., 2.54830761e-03, -1.03542621e-03, -3.83323362e-03], [-4.40293464e-02, 7.64543717e-02, 3.03300345e-02, ..., -6.26757593e-03, -1.08513170e-04, 2.65683265e-03], [-7.54999549e-02, 1.10395307e-01, -1.18728567e-02, ..., -6.28244169e-03, 3.31178413e-03, 1.49452025e-02], [8.57814681e-02, 7.47065097e-02, 9.46676765e-02, ..., 6.59797907e-02, 4.70565367e-02, 1.02763315e-02], [-8.60453052e-02, 3.49035492e-02, -2.52742047e-02, ..., 9.07406575e-02, 3.11910364e-02, -1.25180042e-02], [-4.69196642e-02, -9.46502808e-03, -2.44435088e-03, ..., -3.48238073e-04, -7.65026584e-02, 1.65008138e-02]]) # splitting into train and test from sklearn.model_selection import train_test_split xtrain,xtest,ytrain,ytest = train_test_split(X_train_scaled,y_train,test_size=0.3,random_state=42) Repeat all actions on Test Data test df = pd.read csv('mercedes test.csv') test df.head() ID X0 X1 X2 X3 X4 X5 X6 X8 X10 ... X375 X376 X377 X378 X379 X380 X382 X383 X384 X385 0 0 0 0 0 b ai 0 0 0 0 0 0 as c d y i m 0 s as 5 rows × 377 columns X test = test df.drop(['ID'],axis=1) X test.head() X0 X1 X2 X3 X4 X5 X6 X8 X10 X11 ... X375 X376 X377 X378 X379 X380 X382 X383 X384 X385 0 ... 0 0 0 0 0 0 b ai d g d 0 0 ... 0 0 0 1 0 0 0 as 0 0 0 ... 1 0 0 0 0 0 0 0 0 m as 5 rows × 376 columns X test.isna().sum().sum() X test.nunique(axis=0) 49 Χ0 27 X2 4.5 ХЗ 7 X4 X380 X382 X383 X384 2 X385 2 Length: 376, dtype: int64 In [34]: # dropping constant variance columns #select only the numerical variables X test num = X test.iloc[:,8:] X test num.shape Out[34]: (4209, 368) # remove low variance columns from test data X test num = X test num.drop(const var,axis=1) X test num.head() X10 X12 X13 X14 X15 X16 X17 X18 X19 X20 ... X375 X376 X377 X378 X379 X380 X382 X383 X384 X385 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 3 0 0 0 0 0 0 0 ... 0 0 0 0 0 0 0 5 rows × 356 columns In [34]: # selecting categorical variables and convert them to numeric by using Labelencoder X test cat = X test.iloc[:,:8] X test cat.head() X0 X1 X2 X3 X4 X5 X6 X8 W f as d d as for i in X test cat.columns: X_test_cat[i] = le.fit_transform(X_test_cat[i]) X_test_cat.head() X0 X1 X2 X3 X4 X5 X6 X8 **0** 21 23 34 3 26 0 22 8 0 6 24 **2** 21 23 17 3 0 9 9 **3** 21 13 34 3 31 11 13 **4** 45 20 17 2 3 30 8 12 # create the clean table by concating the numerical and the converted categorical tables X_test_clean = pd.concat([X_test_cat, X_test_num], axis=1) X_test_clean.head() X0 X1 X2 X3 X4 X5 X6 X8 X10 X12 ... X375 X376 X377 X378 X379 X380 X382 X383 X384 X385 **0** 21 23 34 0 22 0 ... 0 0 0 0 0 0 0 0 5 3 26 1 0 0 ... 0 **1** 42 6 24 0 ... 9 0 0 0 1 0 0 0 **2** 21 23 17 0 ... 13 34 11 13 0 0 ... 0 0 0 0 0 0 **4** 45 20 17 3 30 8 12 1 0 5 rows × 364 columns # scaling test data X_test_scaled = scaler.transform(X_test_clean) In [40]: # Transform the test Data to Principal Components with 95% variance explained X test pca = pca.transform(X_test_scaled) X_test_pca.shape Out[40]: (4209, 148) Task-5> Predict your test_df values using XGBoost. In [41]: # import xgboost import xgboost print(xgboost.__version__) 1.5.0 In [42]: # apply XGBRegressor model from xgboost import XGBRegressor xg_reg = XGBRegressor(objective ='reg:linear', n_estimators = 100) xg_reg.fit(X_train_clean,y_train) y_pred = xg_reg.predict(xtest) y_pred [11:44:43] WARNING: d:\bld\xgboost-split 1637426510059\work\src\objective\regression obj.cu:188: reg:linear is now deprecated in favor of reg:squarederror. Out[42]: array([102.79558 , 109.84057 , 134.5284 , ..., 106.056435, 99.99118 , 116.85561], dtype=float32) In [43]: # check for the RMSE from sklearn.metrics import mean squared error print('The RMSE is :',np.sqrt(mean_squared_error(ytest,y_pred))) The RMSE is : 21.772053184535892 In [44]: # predicting the test data final_pred = xg_reg.predict(X_test_clean) final pred Out[44]: array([99.44908 , 114.26388 , 100.843636, ..., 91.6147 , 108.98758 , 94.22524], dtype=float32)