In [2]: import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns # algorithms from sklearn.model selection import train test split from sklearn.preprocessing import StandardScaler from sklearn import preprocessing from sklearn.linear model import LinearRegression from sklearn import metrics from scipy import stats import math **Primary Exploratory Data Analysis** In []: co = pd.read csv(r'F:\concretedata.csv') In [3]: co.head() Cement Blast_Furnace_Slag Fly_Ash Water Superplasticizer Coarse_Aggregate Fine_Aggregate Age_in_day 540.0 0.0 0.0 162.0 2.5 1040.0 676.0 1 540.0 0.0 0.0 162.0 1055.0 676.0 2.5 2 332.5 142.5 0.0 228.0 0.0 932.0 594.0 27 3 332.5 228.0 0.0 142.5 0.0 932.0 594.0 36 4 198.6 132.4 192.0 0.0 978.4 825.5 36 0.0 Changing the name of columns co=co.rename(columns = {"Cement":"cement", "Blast Furnace Slag":"bf slag", "Fly Ash":"f In [4]: co.shape Out[6]: (1030, 9) Binder = cement + bf_slag + f_ash co['binder'] = co['cement'] + co['bf slag'] + co['f ash'] In [9]: co['water binder ratio'] = co['water'] / co['binder'] co['coarse fine ratio'] = co['coarse agg'] / co['fine agg'] co = co.drop(['cement','bf_slag','f_ash','water','coarse_agg','fine_agg','binder'], az converting age_days into log format co['log_age'] = np.log(co['age_days']) superplaticizer is not a binding material but act as a water ratio balancer sns.lineplot(x= 'water_binder_ratio', y = 'cc_strength', data =co) Out[17]: <AxesSubplot:xlabel='water_binder_ratio', ylabel='cc_strength'> 80 70 60 strength 50 40 30 20 10 0.3 water_binder_ratio co = co.drop(['age_days'], axis = 1) In [19]: binning the data into different groups co.describe() s_plasticizer cc_strength water_binder_ratio coarse_fine_ratio log_age 1030.000000 1030.000000 1030.000000 1030.000000 1030.000000 count 6.204660 35.817961 0.469233 1.273752 3.165326 mean std 5.973841 16.705742 0.127123 0.185670 1.191450 min 0.000000 2.330000 0.235073 0.858453 0.000000 25% 0.000000 23.710000 0.383916 1.121488 1.945910 50% 6.400000 34.445000 0.471974 1.266055 3.332205 **75**% 10.200000 46.135000 0.561081 1.358146 4.025352 32.200000 82.600000 0.900000 1.874876 5.899897 max co2 = co.copy()sns.displot(x= 's plasticizer', data = co2, kind = 'hist') <seaborn.axisgrid.FacetGrid at 0x287c0d136a0> 400 350 300 250 200 150 100 50 10 15 25 30 s_plasticizer statistic,bin_edges,binnumber = stats.binned_statistic(co2['s_plasticizer'], co2['s_p. co2['sp_bin'] = pd.DataFrame(zip(bin_edges), columns=['sp_bin']) statistic, bin edges, binnumber = stats.binned statistic(co2['cc strength'], co2['cc str co2['cc bin'] = pd.DataFrame(zip(bin edges), columns=['cc bin']) statistic, bin edges, binnumber = stats.binned statistic(co2['water binder ratio'], co2 co2['wb_bin'] = pd.DataFrame(zip(bin_edges), columns=['wb_bin']) statistic, bin edges, binnumber = stats.binned statistic(co2['coarse fine ratio'], co2[co2['cf bin'] = pd.DataFrame(zip(bin edges), columns=['cf bin']) statistic, bin_edges, binnumber = stats.binned_statistic(co2['log_age'], co2['log_age'], co2['lg bin'] = pd.DataFrame(zip(bin edges), columns=['lg bin']) co2 = co2.drop(['s plasticizer','cc strength','water binder ratio','coarse fine ratio MODEL NO 1 train , test = train_test_split(co2, test_size = 0.3 , random_state = 78 , shuffle = ! x train = train.drop(['cc bin'], axis = 1) In [42]: y train = train['cc bin'] x test = test.drop(['cc bin'], axis = 1) y test = test['cc bin'] num col = [x for x in x train.columns] In [45]: for i in num_col: scale = StandardScaler().fit(x_train[[i]]) x train[i] = scale.transform(x train[[i]]) x_test[i] = scale.transform(x_test[[i]]) LR = LinearRegression() model 12 = LR.fit(x train, y train) PERFORMANCE of model no 1 # train y_pred_train = model_12.predict(x_train) print("MSE" , metrics.mean_squared_error(y_train,y_pred_train)) print("R squared", metrics.r2_score(y_train, y_pred_train)) MSE 4.4487214454247777e-29 R squared 1.0 In [52]: # test y_pred_test = model_12.predict(x test) print("MSE" , metrics.mean_squared_error(y_test,y_pred_test)) print("R squared", metrics.r2_score(y_test, y_pred_test)) MSE 4.891831144139941e-29 R squared 1.0 def mean_absolute_percentage_error(y_true, y_pred): y_true, y_pred = np.array(y_true), np.array(y_pred) return np.mean(np.abs((y_true - y_pred) / y_true)) * 100 print('MAPE:', mean_absolute_percentage_error(y_test, y_pred_test)) MAPE: 3.064408755900244e-14 Data Preprocessing for Model NO 2 co3 = co.copy()binning the data and taking mean as statistics x = co3['s_plasticizer'] x = np.sort(co3['s_plasticizer']) In [87]: print(x) [0. 0. 0. ... 32.2 32.2 32.2] converting into bins bin1 = np.zeros((10,103))In [89]: print(bin1) [[0. 0. 0. ... 0. 0. 0.][0. 0. 0. ... 0. 0. 0.] [0. 0. 0. ... 0. 0. 0.] [0. 0. 0. ... 0. 0. 0.] [0. 0. 0. ... 0. 0. 0.][0. 0. 0. ... 0. 0. 0.]] print(type(bin1)) <class 'numpy.ndarray'> print(type(x)) <class 'numpy.ndarray'> print(x) 0. ... 32.2 32.2 32.2] [0. 0. x numpy array In [94]: print(x) [0. 0. 0. ... 32.2 32.2 32.2] convert into 20 bins bin1 = np.zeros((206,5))print(bin1) In [114... [[0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.][0. 0. 0. 0. 0.]] for i in range (0, 1030, 5): k = int(i/5)mean = (x[i] + x[i+1] + x[i+2] + x[i+3] + x[i+4])/5for j in range(5): bin1[k,j] = meanprint("Bin Mean: \n",bin1) Bin Mean: [[0. 0. 0. 0. 0.1 [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [23.4 23.4 23.4 23.4 23.4] [28.2 28.2 28.2 28.2 28.2] [32.2 32.2 32.2 32.2 32.2]] print(type(bin1)) <class 'numpy.ndarray'> convert numpy array into pandas series one_array = bin1.flatten() In [119... print(one_array) [0. 0. 0. ... 32.2 32.2 32.2] s plasticizer = pd.Series(one array) print(s_plasticizer) 0 0.0 0.0 1 0.0 3 0.0 4 0.0 . . . 1025 32.2 1026 32.2 32.2 1027 1028 32.2 1029 32.2 Length: 1030, dtype: float64 In [131...] $x1 = co2['cc_strength']$ print(type(x1)) <class 'pandas.core.series.Series'> x1 = x1.to numpy()In [134... print(type(x1)) <class 'numpy.ndarray'> x1 = np.sort(x1)print(x1) [2.33 3.32 4.57 ... 80.2 81.75 82.6] bin2 = np.zeros((206,5))print(bin2) In [139... [[0. 0. 0. 0. 0.][0. 0. 0. 0. 0.][0. 0. 0. 0. 0.] . . . [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.]] In [141... for i in range(0, 1030 , 5): k = int(i/5)mean = (x1[i] + x1[i+1] + x1[i+2] + x1[i+3] + x1[i+4])/5for j in range(5): bin2[k,j] = meanprint("bin mean \n", bin2) bin mean [[3.966 3.966 3.966 3.966] [6.146 6.146 6.146 6.146] [7.088 7.088 7.088 7.088] [77.2 77.2] 77.2 77.2 77.2 79.2 79.2 79.2 79.2] [79.2 [80.788 80.788 80.788 80.788 80.788]] bin2 = bin2.flatten() In [142... print(bin2) In [143... [3.966 3.966 3.966 ... 80.788 80.788 80.788] cc strength = pd.Series(bin2) In [200... print(cc strength) 0 3.966 1 3.966 2 3.966 3 3.966 4 3.966 . . . 1025 80.788 1026 80.788 1027 80.788 1028 80.788 1029 80.788 Length: 1030, dtype: float64 water binder ratio x2 = co2['water_binder_ratio'] In [146... $x2 = x2.to_numpy()$ In [147... x2 = np.sort(x2)print(type(x2)) <class 'numpy.ndarray'> bin3 = np.zeros((206,5))print(bin3) [[0. 0. 0. 0. 0.][0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.]In [156... | **for** i **in** range(0,1030,5): k = int(i/5)mean = (x2[i] + x2[i+1] + x2[i+2] + x2[i+3] + x2[i+4])/5for j in range(5): bin3[k,j] = meanprint("bin3 mean \n" , bin3) bin3 mean $[[0.23507335 \ 0.23507335 \ 0.23507335 \ 0.23507335]$ $[0.25202971 \ 0.25202971 \ 0.25202971 \ 0.25202971]$ $[0.26669754 \ 0.26669754 \ 0.26669754 \ 0.26669754 \ 0.26669754]$ [0.78546369 0.78546369 0.78546369 0.78546369] [0.81597363 0.81597363 0.81597363 0.81597363] [0.85322034 0.85322034 0.85322034 0.85322034 0.85322034]] bin3 = bin3.flatten() In [158... print(bin3) $[0.23507335 \ 0.23507335 \ 0.23507335 \ \dots \ 0.85322034 \ 0.85322034 \ 0.85322034]$ water_binder_ratio = pd.Series(bin3) In [204... | print(water_binder_ratio) 0 0.235073 1 0.235073 0.235073 0.235073 3 0.235073 1025 0.853220 1026 0.853220 1027 0.853220 1028 0.853220 1029 0.853220 Length: 1030, dtype: float64 coarse fine ratio: In [164... x4 = co2['coarse_fine_ratio'] In [165...] x4 = x4.to_numpy() In [167... print(type(x4)) <class 'numpy.ndarray'> x4 = np.sort(x4)bin4 = np.zeros((206,5))print(bin4) [[0. 0. 0. 0. 0.][0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.] [0. 0. 0. 0. 0.]**for** i **in** range(0,1030,5): k = int(i/5)mean = (x4[i]+x4[i+1]+x4[i+2]+x4[i+3]+x4[i+4])/5for j in range(5): bin4[k,j] = meanprint(bin4) [[0.85845255 0.85845255 0.85845255 0.85845255 0.85845255] [0.9204926 0.9204926 0.9204926 0.9204926 0.9204926] [0.9407463 0.9407463 0.9407463 0.9407463] [1.83523654 1.83523654 1.83523654 1.83523654] [1.83523654 1.83523654 1.83523654 1.83523654 1.83523654] [1.87487603 1.87487603 1.87487603 1.87487603 1.87487603]] bin4 = bin4.flatten() coarse fine ratio = pd.Series(bin4) In [206... print(coarse_fine_ratio) 0 0.858453 1 0.858453 2 0.858453 3 0.858453 0.858453 1025 1.874876 1026 1.874876 1027 1.874876 1028 1.874876 1029 1.874876 Length: 1030, dtype: float64 log_age x5 = co2['log age']In [176... $x5 = x5.to_numpy()$ In [178... print(type(x5)) <class 'numpy.ndarray'> In [179... x5 = np.sort(x5)In [184... print(x5) [0. 0. 1.09861229 ... 5.89989735 5.89989735 5.89989735] bin5 = np.zeros((206,5))for i in range(0,1030,5): k = int(i/5)mean = (x5[i] + x5[i+1] + x5[i+2] + x5[i+3] + x5[i+4])/5for j in range(5): bin5[k,j] = meanIn [183... print(bin5) [[0.65916737 0.65916737 0.65916737 0.65916737 0.65916737] [1.09861229 1.09861229 1.09861229 1.09861229 1.09861229] [1.09861229 1.09861229 1.09861229 1.09861229 1.09861229] [5.89713869 5.89713869 5.89713869 5.89713869] [5.89989735 5.89989735 5.89989735 5.89989735] [5.89989735 5.89989735 5.89989735 5.89989735 5.89989735]] bin5 = bin5.flatten() print(bin5) $[0.65916737 \ 0.65916737 \ 0.65916737 \ \dots \ 5.89989735 \ 5.89989735 \ 5.89989735]$ log_age = pd.Series(bin5) print(log age) 0 0.659167 1 0.659167 0.659167 2 0.659167 3 0.659167 4 . . . 5.899897 1025 1026 5.899897 1027 5.899897 1028 5.899897 1029 5.899897 Length: 1030, dtype: float64 combinning all these dataframe into pandasdataframe In [189... co3 = pd.DataFrame(series_1) print(co3) 0 0 0.0 0.0 2 0.0 3 0.0 4 0.0 1025 32.2 1026 32.2 1027 32.2 1028 32.2 1029 32.2 [1030 rows x 1 columns] data = {"s plasticizer": s plasticizer, "cc strength": cc strength, "water binder ratio co4 = pd.concat(data, axis =1) co4.shape Out[213... (1030, 5) In [214... co4.isnull().sum() Out[214... s_plasticizer 0 cc_strength 0 water_binder_ratio 0 coarse fine_ratio 0 log age 0 dtype: int64 In [219... sns.displot(x= 'log_age', data = co4 , kind = 'hist') <seaborn.axisgrid.FacetGrid at 0x287c11e6910> 400 350 300 250 200 150 100 50 log_age MODEL NO 2 train_1 , test_1 = train_test_split(co4, test_size = 0.3 , random_state = 68 , shuffle x_train = train_1.drop(['cc_strength'], axis = 1) y_train = train_1['cc_strength'] x_test = test_1.drop(['cc_strength'], axis = 1) y_test = test_1['cc_strength'] num_col = [x for x in x_train.columns] for i in num_col: scale = StandardScaler().fit(x_train[[i]]) x_train[i] = scale.transform(x_train[[i]]) x_test[i] = scale.transform(x_test[[i]]) LR = LinearRegression() model_34 = LR.fit(x_train,y_train) Performance of model no 2 # train y_pred_train = model_34.predict(x_train) print("MSE" , metrics.mean_squared_error(y_train,y_pred_train)) print("R squared", metrics.r2_score(y_train, y_pred_train)) MSE 2.3366510619689884 R squared 0.9916149168413269 # test y_pred_test = model_34.predict(x_test) print("MSE" , metrics.mean_squared_error(y_test,y_pred_test)) print("R squared", metrics.r2_score(y_test, y_pred_test)) MSE 2.3923292904040885 R squared 0.991418680756592 def mean_absolute_percentage_error(y_true, y_pred): y_true, y_pred = np.array(y_true), np.array(y_pred) return np.mean(np.abs((y_true - y_pred) / y_true)) * 100 print('MAPE:', mean_absolute_percentage_error(y_test, y_pred_test)) MAPE: 4.166966965309672 NOTE: Both model got more than 99% accuracy in test data