Pandas used to manipulate the dataset import pandas as pd # numpy is used to manipulate the Numerical and dimentional manipulation of import numpy as np import math import seaborn as sns # for visualization purpose using this seaborn # for visualization purpose using this Matplotlib import matplotlib.pyplot as plt %matplotlib inline # SCIKIT Learn is used to perform machine learning (Here from sklearn.model selection import train test split from sklearn.preprocessing import LabelEncoder print("Libraries Loaded...") Libraries Loaded... train = pd.read csv("train.csv") # reading CSV file using pandas train.head() SQBhogar_total SQBedjefe SC v14a refrig v18q v18q1 r4h1 ... SQBescolari SQBage v2a1 hacdor rooms hacapo **0** ID_279628684 190000.0 0 3 0 1849 100 0 NaN 0 ... 100 1 ID_f29eb3ddd 135000.0 0 4 0 1.0 0 4489 144 144 8 0 0 ID_68de51c94 0 0 0 NaN NaN 121 8464 1 ID_d671db89c 180000.0 5 0 1.0 0 ... 289 0 81 16 121 5 0 ... ID_d56d6f5f5 180000.0 0 0 1369 16 121 1.0 121 5 rows × 143 columns test = pd.read csv("test.csv") # reading CSV file using pandas test.head() refrig v18q v18q1 r4h1 ... v2a1 hacdor v14a SQBescolari SQBage SQBhogar_total rooms hacapo age ID_2f6873615 0 5 0 0 9 NaN 0 NaN 4 16 ID_1c78846d2 0 5 0 41 NaN 256 1681 NaN 5 0 9 ID_e5442cf6a NaN 0 0 NaN 1 41 289 1681 ID_a8db26a79 25 0 14 0 1.0 59 256 3481 NaN ID_a62966799 175000.0 0 4 0 121 324 1 1.0 0 18 1 5 rows × 142 columns In [4]: train df=pd.DataFrame(train) train df.head() Out[4]: v14a v18q v18q1 r4h1 ... SQBescolari SQBage SQBhogar_total SQBedjefe v2a1 hacdor refrig rooms hacapo **0** ID_279628684 190000.0 0 3 0 0 100 1849 100 0 NaN 1 0 1 ID_f29eb3ddd 135000.0 0 1.0 0 144 4489 144 8 0 2 ID_68de51c94 0 0 0 0 ... 1 NaN 121 8464 NaN ID_d671db89c 180000.0 0 5 0 1.0 0 81 289 16 121 5 0 ID_d56d6f5f5 180000.0 0 1.0 0 121 1369 16 121 5 rows × 143 columns Relevent Features used to predict the data features=['rooms','v14a','refrig','abastaguano','sanitario1','lugar1','lugar2','lugar3','lugar4','lugar5','lugar6','area1','area2'] test_df=pd.DataFrame(test) test_df.head() v2a1 hacdor rooms hacapo v14a refrig v18q v18q1 r4h1 ... age SQBescolari SQBage SQBhogar_total SQBedjef **0** ID_2f6873615 0 5 0 0 1 ... 0 9 NaN NaN 16 1 ID_1c78846d2 NaN 0 5 0 NaN 41 256 1681 5 ID_e5442cf6a 0 0 NaN 41 289 1681 NaN 3 ID_a8db26a79 NaN 0 1.0 59 256 3481 **4** ID_a62966799 175000.0 0 1.0 18 121 324 5 rows × 142 columns Understand the type of data train df.shape (9557, 143)test df.shape (23856, 142) #train df.info(verbose = True) Q1. Identify the output variable. "Target Column in the csv file Is the Output variable" Q7. Count how many null values are existing in columns for i in train df.columns: if train df[i].isna().sum() >1: print('Total null values in\t',i,'= ',train df[i].isna().sum()) # Summing all the null values Total null values in v2a1 = 6860v18q1 = 7342Total null values in Total null values in $rez_esc = 7928$ Total null values in meaneduc = 5Total null values in SQBmeaned = 5Total number of missing values in the dataset print("Total null values in the dataset", ' = ', train df.isna().sum().sum()) # Summing all the null values in wh Total null values in the dataset = 22140 Missing Data Handlining v2a1 mean=int(train df['v2a1'].mean()) print("Mean value for v2a1 column = ", v2a1 mean) print("\n\t\t\tReplacing the missing values with mean values of the column") train df['v2a1'].fillna(value=v2a1 mean,inplace=True) # Replacing the missing values with mean values of the train df.head() Mean value for v2a1 column = 165231Replacing the missing values with mean values of the column v2a1 hacdor rooms hacapo v14a refrig v18q v18q1 r4h1 ... SQBescolari SQBage SQBhogar_total SQBedjefe SC **0** ID_279628684 190000.0 3 0 NaN 0 ... 1849 1 100 1 ID_f29eb3ddd 135000.0 144 4489 144 **2** ID_68de51c94 165231.0 8 0 0 8464 0 NaN 121 1 **3** ID_d671db89c 180000.0 0 1.0 289 16 121 ID_d56d6f5f5 180000.0 5 0 16 121 1.0 121 1369 5 rows × 143 columns v18q1 mean=int(train df['v18q1'].mean()) print("Mean value for v18q1 column = ", v18q1 mean) print("\n\t\t\tReplacing the missing values with mean values of the column") train df['v18q1'].fillna(value=v18q1 mean,inplace=True) # Replacing the missing values with mean values of train df.head() Mean value for v18q1 column = 1Replacing the missing values with mean values of the column refrig v18q v18q1 r4h1 ... SQBescolari SQBage SQBhogar_total SQBedjefe SC v14a **0** ID_279628684 190000.0 0 3 0 0 1.0 0 ... 100 1849 1 100 1 ID_f29eb3ddd 135000.0 0 4489 1.0 0 144 144 2 ID 68de51c94 165231.0 0 8 0 0 1.0 0 ... 1 0 121 8464 **3** ID_d671db89c 0 1.0 289 121 180000.0 81 16 ID d56d6f5f5 180000.0 0 ... 0 5 0 121 1369 16 121 1.0 5 rows × 143 columns rez_esc_mean=int(train_df['rez_esc'].mean()) print("Mean value for rez_esc column = ", rez_esc_mean) $print("\n\t\t\t)$ tReplacing the missing values with mean values of the column") train_df['rez_esc'].fillna(value=rez_esc_mean,inplace=True) # Replacing the missing values with mean values train_df.head() Mean value for rez esc column = 0Replacing the missing values with mean values of the column refrig v18q v18q1 r4h1 ... SQBescolari SQBage SQBhogar_total SQBedjefe SC v2a1 hacdor rooms hacapo v14a **0** ID_279628684 190000.0 0 3 0 0 1.0 0 ... 100 1849 1 100 **1** ID_f29eb3ddd 135000.0 0 1.0 4489 144 144 0 8 0 0 0 ... 0 **2** ID_68de51c94 165231.0 1.0 8464 1 121 0 **3** ID_d671db89c 180000.0 1.0 81 289 16 121 5 1369 0 0 0 ... ID_d56d6f5f5 180000.0 1.0 121 16 121 5 rows × 143 columns In [14]: meaneduc_mean=int(train_df['meaneduc'].mean()) print("Mean value for meaneduc column = ", meaneduc_mean) print("\n\t\t\tReplacing the missing values with mean values of the column") train_df['meaneduc'].fillna(value=meaneduc_mean,inplace=True) # Replacing the missing values with mean value train df.head() Mean value for meaneduc column = 9 Replacing the missing values with mean values of the column Out[14]: refrig v18q v18q1 r4h1 ... v2a1 hacdor rooms hacapo v14a SQBescolari SQBage SQBhogar_total SQBedjefe SC **0** ID_279628684 190000.0 0 3 0 1.0 0 100 1849 100 1 ID_f29eb3ddd 135000.0 0 4 0 1.0 0 144 4489 144 0 8 0 0 0 ... **2** ID_68de51c94 165231.0 1.0 121 8464 1 0 5 0 **3** ID_d671db89c 180000.0 0 1.0 0 81 289 16 121 5 0 ID_d56d6f5f5 180000.0 0 1.0 0 121 1369 16 121 5 rows × 143 columns SQBmeaned mean=int(train df['SQBmeaned'].mean()) print("Mean value for SQBmeaned column = ", SQBmeaned mean) print("\n\t\t\tReplacing the missing values with mean values of the column") train df['SQBmeaned'].fillna(value=SQBmeaned mean,inplace=True) # Replacing the missing values with mean val train df.head() Mean value for SQBmeaned column = 102 Replacing the missing values with mean values of the column v2a1 hacdor refrig v18q v18q1 r4h1 ... SQBescolari SQBage SQBhogar_total hacapo **0** ID_279628684 190000.0 3 0 100 1849 100 4 0 144 **1** ID_f29eb3ddd 135000.0 1.0 144 4489 2 ID_68de51c94 165231.0 0 8 0 1.0 8464 0 121 ID_d671db89c 5 0 81 289 16 121 180000.0 1.0 0 ID_d56d6f5f5 180000.0 5 0 121 1369 16 121 5 rows × 143 columns Check For Missing Data That Was filled or not for i in train df.columns: if train df[i].isna().sum() >1: $print('Total\ null\ values\ in\t',i,'=',train_df[i].isna().sum())\ \ \textit{\# Summing all\ the\ null\ values}$ else: print("There is no null value (or) Values are replaced with the mean values of tye column") There is no null value (or) Values are replaced with the mean values of tye column Q4. Check whether all members of the house have the same poverty level. print("House Id with different Poverty Level : ", train df.groupby('idhogar')['Target'].unique().head()) House Id with different Poverty Level: idhogar 001ff74ca [4] 003123ec2 [2] 004616164 [2] 004983866 [3] 005905417 [2] Name: Target, dtype: object Enable Encoding method for the catagorical data for i in train_df.columns: if train_df[i].dtype == 'object': print(i,'= ',train_df[i].dtype) Id = object idhogar = object dependency = object edjefe = object edjefa = object idhog =list(train df["idhogar"].unique()) for temp in idhog : (temp, '=', list(train df["idhogar"]).count(temp)) dep_=list(train_df["dependency"].unique()) for temp in dep_: print(temp, '=', list(train_df["dependency"]).count(temp)) no = 17478 = 378yes = 21923 = 236.5 = 1497.25 = 2602 = 730.66666669 = 487.33333334 = 5981.5 = 713.40000001 = 84.75 = 981.25 = 18.2 = 902.5 = 771.2 = 114 = 1001.3333334 = 842.25 = 13.22222222 = 115 = 24.83333331 = 11.80000001 = 186 = 73.5 = 181.66666666 = 8.2857143 = 91.75 = 11.71428573 = 12.16666667 = 7.60000002 = 8ED1=train_df["edjefe"].unique() for temp in ED1: print(temp, '=', list(train df["edjefe"]).count(temp)) 10 = 11112 = 113no = 376211 = 7519 = 48615 = 2854 = 1376 = 18458 = 25717 = 2027 = 23416 = 13414 = 2085 = 22221 = 432 = 19419 = 14yes = 1233 = 30718 = 1913 = 10320 = 7ED2=train df["edjefa"].unique() print(temp, '=', list(train_df["edjefa"]).count(temp)) no = 623011 = 3994 = 13610 = 969 = 23715 = 1887 = 17914 = 12013 = 528 = 21717 = 766 = 9475 = 1763 = 15216 = 11319 = 4yes = 6921 = 512 = 722 = 84 20 = 218 = 3le Id = LabelEncoder() #le idhogar = LabelEncoder() le dependency = LabelEncoder() le edjefe = LabelEncoder() le edjefa = LabelEncoder() #train_df['ID']=le_Id.fit_transform(train_df['ID']) train_df['dependency']=le_dependency.fit_transform(train_df['dependency']) train_df['edjefe']=le_edjefe.fit_transform(train_df['edjefe']) train_df['edjefa']=le_edjefa.fit_transform(train_df['edjefa']) dummy=pd.get_dummies(train_df['v2a1']) Q3. Check if there are any biases in your dataset. In [24]: val= list(train_df[['Target']] .values) uniq=list(np.unique(train_df['Target'])) # This is unique values in the target column $\#print("Unique Values in TARGET colimn are : ", uniq,'\n')$ valu=[] for j in uniq: valu.append(val.count(j)) print('Number of elements in ',j,' = ',val.count(j)) plt.bar(j,valu) Number of elements in 1 = 755Number of elements in 2 = 1597Number of elements in 3 = 1209Number of elements in 4 = 59966000 5000 4000 3000 2000 1000 0 -2.0 2.5 3.0 3.5 4.0 '''correlation= train df.corr() print(correlation['v2a1'].values>-0.5) correlation''' "correlation= train df.corr()\nprint(correlation['v2a1'].values>-0.5)\ncorrelation" sns.heatmap(train_df[['rooms','tipovivi3','coopele','v14a','refrig','abastaguano','sanitario1','lugar1','lugar2 Out[26]: <AxesSubplot:> 1.00 rooms tipovivi3 0.75 coopele vl4a 0.50 refria abastaguano 0.25 sanitario1 0.00 lugar1 lugar2 -0.25lugar3 lugar4 -0.50lugar5 lugar6 -0.75area1 area2 -1.00v14a lugarl coopele sanitario1 In [40]: #features=['rooms','v14a','refrig','abastaguano','sanitario1','lugar1','lugar2','lugar3','lugar4','lugar5','lug x train=train df.drop(['Target','Id','idhogar'],axis=1) # Assigning feartures x train.head() Out[40]: age SQBescolari SQBage SQBhogar_total SQBedjefe v2a1 hacdor rooms hacapo v14a refrig v18q v18q1 r4h1 r4h2 **0** 190000.0 0 3 0 1.0 43 100 1849 100 **1** 135000.0 0 1.0 67 144 4489 144 **2** 165231.0 0 8 0 0 1.0 0 92 121 8464 1 0 **3** 180000.0 1.0 17 81 289 16 121 **4** 180000.0 0 1.0 0 37 121 1369 16 121 5 rows × 140 columns In [41]: y train=train df[['Target']] # Assigning target variable y_train.head() Out[41]: **Target** 0 4 2 4 3 4 In [42]: x train, x test, y train, y test=train test split(x train, y train, test size=0.3, random state=123) print("Train & test data splitted successfully") Train & test data splitted successfully In [43]: train df.head() Out[43]: v2a1 hacdor rooms hacapo v14a refrig v18q v18q1 r4h1 ... SQBescolari SQBage SQBhogar_total SQBedjefe SC **0** ID 279628684 190000.0 0 3 0 1849 100 **1** ID_f29eb3ddd 135000.0 1.0 144 4489 144 0 **2** ID_68de51c94 165231.0 0 8 0 1.0 0 ... 121 8464 0 **3** ID_d671db89c 180000.0 0 1.0 289 121 ID_d56d6f5f5 180000.0 0 5 0 121 1369 16 121 5 rows × 143 columns In [44]: print(x_train.shape) print(x test.shape) print(y train.shape) print(y_test.shape) (6689, 140)(2868, 140)(6689, 1)(2868, 1)In [45]: from sklearn.linear model import LinearRegression from sklearn.metrics import mean squared error lr=LinearRegression(normalize=True) lr.fit(x_train,y_train) print(math.sqrt(mean_squared_error (y_train , lr.predict (x_train)))) print(math.sqrt(mean squared error (y test , lr.predict (x test)))) print('R2 Value/Coefficient of Determination: {}'.format(lr.score (x_test , y_test))) predictions=lr.predict(x_train) print('Predicted_data = ',predictions[:5]) 0.8106221160734605 0.8082601001146387 R2 Value/Coefficient of Determination: 0.3523738002461617 Predicted_data = [[3.2109375] [3.5625 [2.6953125] [3.703125] [3.46875]] Q9. Predict the accuracy using random forest classifier. **Prediction and Acuracy valuvation** In [46]: from sklearn.ensemble import RandomForestClassifier classifier=RandomForestClassifier(random state=123) classifier.fit(x_train,y_train.values.ravel()) Out[46]: RandomForestClassifier(random_state=123) In [47]: y_pred_test = classifier.predict(x_test) y pred test Out[47]: array([1, 4, 4, ..., 4, 4, 2], dtype=int64) In [48]: classifier.score(x test,y test) Out[48]: 0.9079497907949791 Q10. Check the accuracy using random forest with cross validation. In [54]: # Accuracy rate from sklearn.metrics import accuracy score, confusion matrix, classification report print("Accuracy Score : ",accuracy score(y test, y pred test)) Accuracy Score: 0.9079497907949791 print(classification report(y test, y pred test)) precision recall f1-score support 0.93 0.70 0.80 223 0.91 0.86 0.82 468 0.95 0.90 1.00 3 0.81 359 0.95 1818 0.91 2868 0.92 0.81 0.85 2868 macro avg weighted avg 0.91 0.91 0.90 2868 # View confusion matrix for test data and predictions cm=confusion matrix(y test, y pred test) Out[53]: array([[156, 21, 1, 45]**,** [9, 385, 10, 64], 1, 92], 14, 252, 3, 1811]], dtype=int64) print(predictions.var()) 0.3451027003042477 print(y_pred_test.var()) 0.8904561029315937 Q5 Check if there is a house without a family head # parentescol, =1 if household head #tamhog,estadocivil3,parentesco1 yes=(list(train['parentescol'].values)).count(0) no=(list(train['parentesco1'].values)).count(1) print('Total house without family head : ',yes) Total house without family head: 6584