Loading Data In [1]: import numpy as np import pandas as pd import matplotlib.pyplot as plt %matplotlib inline import seaborn as sns import warnings warnings.filterwarnings('ignore') In [2]: df = pd.read_csv("C:/Users/User/Desktop/application record.csv") df.head() ID CODE_GENDER FLAG_OWN_CAR FLAG_OWN_REALTY CNT_CHILDREN AMT_INCOME_TOTAL NAME_INCOME_TYPE NAME_EDUCATION Out[2]: 0 5008804 M 427500.0 Working Higher e **1** 5008805 427500.0 Working Higher e Secondary / s **2** 5008806 112500.0 Working Secondary / s **3** 5008808 270000.0 Commercial associate Secondary / s **4** 5008809 0 Ν 270000.0 Commercial associate df.shape In [3]: (438557, 18)Out[3]: df.info() In [4]: <class 'pandas.core.frame.DataFrame'> RangeIndex: 438557 entries, 0 to 438556 Data columns (total 18 columns): # Column Non-Null Count Dtype -----0 ID 438557 non-null int64
1 CODE_GENDER 438557 non-null object
2 FLAG_OWN_CAR 438557 non-null object
3 FLAG_OWN_REALTY 438557 non-null object
4 CNT_CHILDREN 438557 non-null int64 5 AMT_INCOME_TOTAL 438557 non-null float64 6 NAME_INCOME_TYPE 438557 non-null object 7 NAME EDUCATION_TYPE 438557 non-null object 8 NAME FAMILY_STATUS 438557 non-null object 9 NAME_HOUSING_TYPE 438557 non-null object

 10
 DAYS_BIRTH
 438557 non-null int64

 11
 DAYS_EMPLOYED
 438557 non-null int64

 12
 FLAG_MOBIL
 438557 non-null int64

 12 FLAG_MOBIL 438557 non-null int64 13 FLAG_WORK_PHONE 438557 non-null int64 14 FLAG_PHONE 438557 non-null int64 15 FLAG EMAIL 438557 non-null int64 16 OCCUPATION_TYPE 304354 non-null object 17 CNT_FAM_MEMBERS 438557 non-null float64 dtypes: float64(2), int64(8), object(8) memory usage: 60.2+ MB In [5]: credit_df = pd.read_csv("C:/Users/User/Desktop/credit_record.csv") credit df.head() Out[5]: ID MONTHS_BALANCE STATUS **0** 5001711 0 **1** 5001711 0 2 5001711 -2 0 **3** 5001711 **4** 5001712 credit df.shape (1048575, 3)Out[6]: In [7]: credit_df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 1048575 entries, 0 to 1048574 Data columns (total 3 columns): # Column Non-Null Count 0 ID 1048575 non-null int64 1 MONTHS_BALANCE 1048575 non-null int64 2 STATUS 1048575 non-null object dtypes: int64(2), object(1) memory usage: 24.0+ MB In [8]: # dropping occupation type which has many null values df.drop('OCCUPATION TYPE', axis=1, inplace=True) In [9]: # Checking duplicates in 'ID' column len(df['ID']) - len(df['ID'].unique()) Out[9]: # Dropping duplicate entries from ID column In [10]: df = df.drop_duplicates('ID', keep='last') In [11]: # Checking Non-Numerical Columns cat_columns = df.columns[(df.dtypes =='object').values].tolist() cat columns Out[11]: ['CODE_GENDER', 'FLAG OWN CAR', 'FLAG OWN REALTY', 'NAME INCOME TYPE', 'NAME EDUCATION TYPE', 'NAME FAMILY STATUS', 'NAME HOUSING TYPE'] In [12]: # Checking Numerical Columns df.columns[(df.dtypes !='object').values].tolist() ['ID', Out[12]: 'CNT CHILDREN', 'AMT INCOME TOTAL', 'DAYS BIRTH', 'DAYS EMPLOYED', 'FLAG MOBIL', 'FLAG WORK PHONE', 'FLAG PHONE', 'FLAG EMAIL', 'CNT FAM MEMBERS'] In [13]: # Checking unique values from Categorical Columns for i in df.columns[(df.dtypes =='object').values].tolist(): print(df[i].value counts()) CODE GENDER 294412 M 144098 Name: CODE GENDER, dtype: int64 FLAG OWN CAR 275428 Y 163082 Name: FLAG OWN CAR, dtype: int64 FLAG OWN REALTY Y 304043 N 134467 Name: FLAG OWN REALTY, dtype: int64 NAME INCOME TYPE Working Commercial associate 100739 Pensioner State servant 75483 36184 State servant Student Name: NAME INCOME TYPE, dtype: int64 NAME EDUCATION TYPE Secondary / secondary special Higher education 117509 Incomplete higher 14849 Lower secondary 4051 312 Academic degree Name: NAME EDUCATION TYPE, dtype: int64 NAME FAMILY STATUS Married 55268 36524 Single / not married Civil marriage Separated 27249 Widow Name: NAME FAMILY STATUS, dtype: int64 NAME HOUSING TYPE House / apartment 393788 With parents 19074 19074 14213 Municipal apartment Rented apartment 5974
Office apartment 3922 Co-op apartment Name: NAME HOUSING TYPE, dtype: int64 In [14]: # Checking unique values from Numerical Columns df['CNT CHILDREN'].value_counts() 0 304038 Out[14]: 1 88518 2 39879 3 5430 4 486 133 5 7 9 12 6 14 19 1 Name: CNT CHILDREN, dtype: int64 In [15]: # Checking Min , Max values from 'DAYS BIRTH' column print('Min DAYS BIRTH :', df['DAYS BIRTH'].min(),'\nMax DAYS BIRTH :', df['DAYS BIRTH'].max()) Min DAYS BIRTH : -25201 Max DAYS_BIRTH : -7489 In [16]: # Converting 'DAYS_BIRTH' values from Day to Years df['DAYS_BIRTH'] = round(df['DAYS BIRTH']/-365,0) df.rename(columns={'DAYS_BIRTH':'AGE_YEARS'}, inplace=True) # Checking unique values greater than 0 df[df['DAYS_EMPLOYED']>0]['DAYS_EMPLOYED'].unique() Out[16]: array([365243], dtype=int64) # As mentioned in document, if 'DAYS_EMPLOYED' is positive no, it means person currently unemployed, hence repl In [17]: df['DAYS_EMPLOYED'].replace(365243, 0, inplace=True) # Converting 'DAYS EMPLOYED' values from Day to Years In [18]: df['DAYS EMPLOYED'] = abs(round(df['DAYS EMPLOYED']/-365,0)) df.rename(columns={'DAYS EMPLOYED':'YEARS EMPLOYED'}, inplace=True) In [19]: df['FLAG MOBIL'].value counts() 1 438510 Out[19]: Name: FLAG MOBIL, dtype: int64 In [20]: # As all the values in column are 1, hence dropping column df.drop('FLAG MOBIL', axis=1, inplace=True) In [21]: df['FLAG_WORK_PHONE'].value_counts() 348118 Out[21]: 90392 Name: FLAG WORK PHONE, dtype: int64 In [22]: # This column only contains 0 & 1 values for Mobile no submitted, hence dropping column df.drop('FLAG WORK PHONE', axis=1, inplace=True) In [23]: df['FLAG_PHONE'].value counts() 312323 Out[23]: 126187 Name: FLAG PHONE, dtype: int64 In [24]: # This column only contains 0 & 1 values for Phone no submitted, hence dropping column df.drop('FLAG PHONE', axis=1, inplace=True) In [25]: df['FLAG_EMAIL'].value_counts() 391062 Out[25]: 1 47448 Name: FLAG EMAIL, dtype: int64 In [26]: # This column only contains 0 & 1 values for Email submitted, hence dropping column df.drop('FLAG EMAIL', axis=1, inplace=True) In [27]: df['CNT_FAM_MEMBERS'].value counts() 233867 2.0 Out[27]: 1.0 84483 3.0 77119 4.0 37351 5.0 5081 6.0 459 124 7.0 9.0 11.0 14.0 8.0 15.0 20.0 Name: CNT FAM MEMBERS, dtype: int64 In [28]: df.head() ID CODE GENDER FLAG OWN CAR FLAG OWN REALTY CNT CHILDREN AMT INCOME TOTAL NAME INCOME TYPE NAME EDUCATION Out[28]: 0 5008804 Υ Υ 0 427500.0 Working Higher e M **1** 5008805 0 427500.0 Working Higher e M Secondary / s 2 5008806 Μ 0 112500.0 Working Secondary / s **3** 5008808 270000.0 Ν Commercial associate Secondary / s 4 5008809 270000.0 Commercial associate Visualization #create plot to detect outliers In [29]: sns.boxplot(df['CNT CHILDREN']) <AxesSubplot:xlabel='CNT CHILDREN'> Out[29]: 10.0 12.5 15.0 17.5 0.0 2.5 5.0 7.5 CNT_CHILDREN sns.boxplot(df['AMT INCOME TOTAL']) In [30]: <AxesSubplot:xlabel='AMT_INCOME_TOTAL'> Out[30]: AMT_INCOME_TOTAL sns.boxplot(df['AGE YEARS']) In [31]: <AxesSubplot:xlabel='AGE YEARS'> Out[31]: 50 AGE_YEARS sns.boxplot(df['YEARS EMPLOYED']) <AxesSubplot:xlabel='YEARS EMPLOYED'> 10 40 20 30 YEARS_EMPLOYED sns.boxplot(df['CNT FAM MEMBERS']) In [33]: <AxesSubplot:xlabel='CNT FAM MEMBERS'> Out[33]: 2.5 5.0 7.5 10.0 12.5 15.0 17.5 20.0 CNT_FAM_MEMBERS **Data cleansing (Removing Outliers)** In [34]: high bound = df['CNT CHILDREN'].quantile(0.999) print('high bound :', high bound) low bound = df['CNT CHILDREN'].quantile(0.001) print('low bound :', low bound) high bound: 4.0 low bound : 0.0 In [35]: df = df[(df['CNT_CHILDREN']>=low_bound) & (df['CNT_CHILDREN']<=high_bound)]</pre> In [36]: high_bound = df['AMT_INCOME_TOTAL'].quantile(0.999) print('high_bound :', high_bound) low_bound = df['AMT_INCOME_TOTAL'].quantile(0.001) print('low_bound :', low_bound) high bound : 990000.0 low bound : 36000.0 In [37]: df = df[(df['AMT_INCOME_TOTAL']>=low_bound) & (df['AMT_INCOME_TOTAL']<=high bound)]</pre> In [38]: high_bound = df['YEARS_EMPLOYED'].quantile(0.999) print('high bound :', high bound) low bound = df['YEARS EMPLOYED'].quantile(0.001) print('low_bound :', low_bound) high bound: 40.0 low bound : 0.0 In [39]: df = df[(df['YEARS_EMPLOYED']>=low_bound) & (df['YEARS EMPLOYED']<=high bound)]</pre> In [40]: high_bound = df['CNT_FAM_MEMBERS'].quantile(0.999) print('high_bound :', high_bound) low bound = df['CNT FAM MEMBERS'].quantile(0.001) print('low_bound :', low_bound) high bound : 6.0 low bound : 1.0 In [41]: df = df[(df['CNT_FAM_MEMBERS']>=low_bound) & (df['CNT_FAM_MEMBERS']<=high bound)]</pre> df.head() In [42]: ID CODE_GENDER FLAG_OWN_CAR FLAG_OWN_REALTY CNT_CHILDREN AMT_INCOME_TOTAL NAME_INCOME_TYPE NAME_EDUCATION Out[42]: **0** 5008804 Υ 0 Μ 427500.0 Working Higher e **1** 5008805 Μ 427500.0 Working Higher e Secondary / s 2 5008806 Υ 0 M 112500.0 Working Secondary / s **3** 5008808 Ν 270000.0 Commercial associate Secondary / s 4 5008809 Commercial associate Ν 0 270000.0 On File - Credit Record.csv credit df.head() In [43]: Out[43]: ID MONTHS_BALANCE STATUS **0** 5001711 0 Χ **1** 5001711 2 5001711 -2 0 **3** 5001711 **4** 5001712 0 C df.isnull().sum() In [44]: 0 ID Out[44]: CODE GENDER 0 FLAG OWN CAR 0 FLAG OWN REALTY 0 CNT CHILDREN 0 AMT INCOME TOTAL NAME INCOME TYPE NAME EDUCATION TYPE NAME FAMILY STATUS NAME HOUSING TYPE AGE YEARS YEARS EMPLOYED 0 CNT FAM MEMBERS 0 dtype: int64 credit_df['STATUS'].value_counts() In [45]: 442031 Out[45]: 0 383120 209230 Χ 1 11090 5 1693 2 868 3 320 4 223 Name: STATUS, dtype: int64 In [46]: # categorizing 'STATUS' column to binary classification 0 : Good Client and 1 : bad client credit_df['STATUS'].replace(['C', 'X'],0, inplace=True) In [47]: credit_df['STATUS'].replace(['2','3','4','5'],1, inplace=True) In [48]: credit_df['STATUS'] = credit_df['STATUS'].astype('int') In [49]: credit_df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 1048575 entries, 0 to 1048574 Data columns (total 3 columns): # Column Non-Null Count Dtype 1048575 non-null int64 0 ID 1 MONTHS_BALANCE 1048575 non-null int64 2 STATUS 1048575 non-null int32 dtypes: int32(1), int64(2) memory usage: 20.0 MB In [50]: credit_df['STATUS'].value_counts(normalize=True) *100 98.646353 Out[50]: 1.353647 Name: STATUS, dtype: float64 In [51]: credit_df_trans = credit_df.groupby('ID').agg(max).reset_index() In [52]: credit_df_trans.drop('MONTHS_BALANCE', axis=1, inplace=True) credit df trans.head() **ID STATUS** Out[52]: **0** 5001711 **1** 5001712 **2** 5001713 **3** 5001714 **4** 5001715 In [53]: credit_df_trans['STATUS'].value_counts(normalize=True)*100 88.365771 Out[53]: 11.634229 Name: STATUS, dtype: float64 **Merging Dataframes** In [92]: # merging the two datasets based on 'ID' data = pd.merge(df, credit df trans, on='ID', how='inner') data.shape (36326, 14)Out[92]: data.shape In [55]: (36326, 14)Out[55]: data.drop('ID', axis=1, inplace=True) In [56]: # checking if there are still duplicate rows in Final Dataframe In [57]: len(data) - len(data.drop duplicates()) 25268 Out[57]: In [58]: # Dropping duplicate records data = data.drop duplicates() data.reset_index(drop=True ,inplace=True) In [59]: data.shape (11058, 13)Out[59]: data.isnull().sum() In [60]: 0 CODE GENDER Out[60]: FLAG OWN CAR 0 FLAG OWN REALTY 0 CNT CHILDREN 0 AMT INCOME TOTAL 0 NAME INCOME TYPE 0 NAME EDUCATION TYPE 0 NAME FAMILY STATUS 0 NAME HOUSING TYPE 0 AGE YEARS 0 YEARS EMPLOYED 0 CNT FAM MEMBERS 0 STATUS 0 dtype: int64 data['STATUS'].value_counts(normalize=True) *100 In [61]: 78.513294 Out[61]: 21.486706 Name: STATUS, dtype: float64 **Feature Engineering** data.head() In [62]: CODE_GENDER FLAG_OWN_CAR FLAG_OWN_REALTY CNT_CHILDREN AMT_INCOME_TOTAL NAME_INCOME_TYPE NAME_EDUCATION_TYPE Out[62]: 0 Υ Υ 0 427500.0 M Working Higher education Secondary / secondary 112500.0 1 M Working special Secondary / secondary 2 F 0 Commercial associate 270000.0 Ν special 3 Ν 283500.0 Pensioner Higher education 4 Υ 0 M 270000.0 Working Higher education In [63]: cat_columns = data.columns[(data.dtypes =='object').values].tolist() cat columns ['CODE GENDER', Out[63]: 'FLAG OWN CAR', 'FLAG OWN REALTY', 'NAME INCOME TYPE', 'NAME EDUCATION TYPE', 'NAME FAMILY STATUS', 'NAME_HOUSING_TYPE'] #Converting all Non-Numerical Columns to Numerical In [64]: from sklearn.preprocessing import LabelEncoder for col in cat_columns: globals()['LE_{{}}'.format(col)] = LabelEncoder() data[col] = globals()['LE {}'.format(col)].fit transform(data[col]) CODE_GENDER FLAG_OWN_CAR FLAG_OWN_REALTY CNT_CHILDREN AMT_INCOME_TOTAL NAME_INCOME_TYPE NAME_EDUCATION_TYPE Out[64]: 0 427500.0 1 1 1 0 4 1 1 0 112500.0 2 0 0 1 0 270000.0 0 4 0 0 283500.0 1 270000.0 for col in cat columns: In [65]: print(col , " : ", globals()['LE_{{}}'.format(col)].classes_) CODE GENDER : ['F' 'M'] FLAG OWN CAR : ['N' 'Y'] : ['N' 'Y'] FLAG OWN REALTY : ['Commercial associate' 'Pensioner' 'State servant' 'Student' 'Working'] NAME INCOME TYPE NAME_EDUCATION_TYPE : ['Academic degree' 'Higher education' 'Incomplete higher' 'Lower secondary' 'Secondary / secondary special'] NAME FAMILY STATUS : ['Civil marriage' 'Married' 'Separated' 'Single / not married' 'Widow'] NAME_HOUSING_TYPE : ['Co-op apartment' 'House / apartment' 'Municipal apartment' 'Office apartment' 'Rented apartment' 'With parents'] data.corr() In [66]: CODE_GENDER FLAG_OWN_CAR FLAG_OWN_REALTY CNT_CHILDREN AMT_INCOME_TOTAL NAME_INCOME_TYPE Out[66]: CODE_GENDER 1.000000 0.348307 -0.052647 0.052351 0.199358 0.072725 -0.000126 FLAG_OWN_CAR 0.348307 1.000000 0.087407 0.218026 0.030008 FLAG_OWN_REALTY -0.052647 -0.000126 1.000000 0.001740 0.036549 -0.034673 CNT_CHILDREN 0.052351 0.087407 0.001740 1.000000 0.027806 0.086772 AMT_INCOME_TOTAL 0.199358 0.218026 0.027806 0.036549 1.000000 -0.083800 NAME_INCOME_TYPE 0.072725 0.030008 1.000000 -0.034673 0.086772 -0.083800 NAME_EDUCATION_TYPE 0.013080 -0.085706 0.003771 -0.025059 -0.231827 0.077636 NAME_FAMILY_STATUS -0.075315 -0.125707 -0.011871 -0.015623 -0.029347 -0.166742NAME HOUSING TYPE 0.036582 0.052185 -0.011203 -0.178070 0.010091 -0.014723 AGE_YEARS -0.157507 -0.106634 0.121351 -0.326642 -0.044799 -0.191870 YEARS_EMPLOYED -0.043657 0.006570 -0.011665 0.033979 0.079548 0.163301 CNT_FAM_MEMBERS 0.078853 0.138978 0.009698 0.025272 0.080072 0.884676 **STATUS** 0.002039 -0.003826 -0.022887 0.004820 0.020289 -0.002579 top corr features = data.corr().index plt.figure(figsize=(20, 20)) sns.heatmap(data[top_corr_features].corr(),annot=True,cmap="RdYlGn") <AxesSubplot:> Out[67]: CODE GENDER 0.35 -0.053 0.052 0.2 0.073 0.013 -0.075 0.052 -0.16 -0.0440.079 0.002 FLAG_OWN_CAR -0.00013 0.03 -0.13 -0.0038 0.8 FLAG_OWN_REALTY -0.0017 0.037 -0.035 -0.012 0.12 -0.053 -0.00013 0.0038 -0.012 0.0097 -0.0230.6 CNT_CHILDREN -0.087 0.0017 0.034 AMT_INCOME_TOTAL -0.037 0.028 -0.084 -0.016 0.22 -0.015-0.0450.08 0.025 0.02 0.4 NAME_INCOME_TYPE --0.035 -0.084 -0.0026 0.013 -0.025 0.078 0.0048 -0.032 0.15 -0.016 NAME_EDUCATION_TYPE -0.0860.0038 -0.021 -0.0064 - 0.2 -0.13 -0.012 -0.016 0.04 0.077 -0.055 0.0035 NAME_FAMILY_STATUS - 0.0 0.052 -0.011 0.01 -0.015 0.037 -0.032 0.04 -0.042 -0.022 0.016 NAME_HOUSING_TYPE AGE_YEARS -0.16 -0.11 0.12 -0.045 0.15 0.017 -0.023 -0.2 -0.044 0.0066 -0.012 0.034 0.08 0.16 -0.016 -0.055 -0.042 0.017 0.055 YEARS_EMPLOYED -8.1e-05 CNT_FAM_MEMBERS 0.14 0.0097 0.08 -0.021 -0.022 0.055 -0.4STATUS 0.002 -0.0038 -0.023 0.0048 -0.0026 -0.0064 CODE_GENDER FLAG_OWN_CAR CNT_CHILDREN NAME_EDUCATION_TYPE YEARS_EMPLOYED CNT_FAM_MEMBERS NAME_HOUSING_TYPE features = data.drop(['STATUS'], axis=1) label = data['STATUS'] features.head() In [69]: CODE_GENDER FLAG_OWN_CAR FLAG_OWN_REALTY CNT_CHILDREN AMT_INCOME_TOTAL NAME_INCOME_TYPE NAME_EDUCATION_TYPE Out[69]: 0 427500.0 1 0 112500.0 1 0 4 2 0 0 1 0 270000.0 0 4 3 0 0 0 283500.0 4 1 1 1 0 270000.0 4 1 label.head() In [70]: 1 Out[70]: 0 2 0 3 0 Name: STATUS, dtype: int32 **Machine Learning Model** Random Forest before Balancing In [71]: from sklearn.model_selection import train test split x_train, x_test, y_train, y_test = train_test_split(features, test size=0.2, random state = 15)In [82]: # Using RandomForest Classifier to build the forecasting model from sklearn.ensemble import RandomForestClassifier rf=RandomForestClassifier(random_state=20, n_estimators=100, criterion='entropy') rf.fit(x_train, y_train) RandomForestClassifier(criterion='entropy', random state=20) Out[82]: In [83]: print('Random Forest Accuracy : ', rf.score(x_test, y test)*100, '%') from sklearn.metrics import confusion matrix prediction = rf.predict(x test) cf matrix=confusion matrix(y test, prediction) print('\nConfusion matrix :') print(confusion matrix(y test, prediction)) from sklearn.metrics import classification report print('\nClassification report:') print(classification_report(y_test, prediction)) Random Forest Accuracy : 64.46654611211574 % Confusion matrix : [[1420 317] [469 6]] Classification report: precision recall f1-score support 0 0.75 0.82 0.78 1737 0.02 0.01 0.02 475 0.64 2212 accuracy 0.39 0.42 0.40 2212 macro avg 2212 weighted avg 0.59 0.64 0.62 In [128... condition = ["TN", "FP", "FN", "TP"] counts = ['{0:0.0f}'.format(value) for value in cf matrix.flatten()] percent = ['{0:.2%}'.format(value) for value in cf matrix.flatten()/np.sum(cf matrix)] labels = $[f'{v1}\n{v2}\n{v3}'$ for v1, v2, v3 in zip(condition,counts,percent)] labels = np.asarray(labels).reshape(2,2) sns.heatmap(cf matrix, annot=labels, fmt='', cmap='Blues') <AxesSubplot:> Out[128]: 1400 1200 FP 317 14.33% 1000 800 600 ΤP 400 469 21.20% 200 In [133... # Getting the auc score from sklearn.metrics import roc_auc_score prob rf=rf.predict proba(x test) roc_auc_score_rf=roc_auc_score(y_test,prob_rf[:,1]) print('The AUC score is: '+ str(roc_auc_score_rf)) The AUC score is: 0.24015028936763327 **Decision Tree before Balancing** In [116... **from** sklearn **import** tree dt =tree.DecisionTreeClassifier(criterion='entropy') dt.fit(x_train, y_train)







