import re import os import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt import datetime as dt accept file path = "Data/accepted 2007 to 2018Q4.csv" #reading accept csv file accept df = pd.read csv(accept file path) C:\Users\glori\anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3146: Dtyp eWarning: Columns (0,19,49,59,118,129,130,131,134,135,136,139,145,146,147) have mixed types. Specify dtype option on import or set low memory=False. has_raised = await self.run_ast_nodes(code_ast.body, cell name, In [4]: reject file path = "Data/rejected 2007 to 2018Q4.csv" #reading reject csv file reject df = pd.read csv(reject file path) list(accept df) Out[6]: ['id', 'member id', 'loan_amnt', 'funded amnt', 'funded amnt inv', 'term', 'int rate', 'installment', 'grade', 'sub grade', 'emp_title' 'emp length', 'home ownership', 'annual inc', 'verification_status', 'issue d', 'loan status', 'pymnt plan', 'url', 'desc', 'purpose', 'title', 'zip code', 'addr_state', 'dti', 'deling 2yrs', 'earliest cr line', 'fico range low', 'fico range high', 'inq last 6mths', 'mths_since_last_delinq', 'mths_since_last_record', 'open_acc', 'pub rec', 'revol_bal' 'revol util', 'total acc', 'initial list status', 'out_prncp', 'out_prncp_inv', 'total_pymnt', 'total_pymnt_inv',
'total_rec_prncp', 'total rec int', 'total rec late fee', 'recoveries', 'collection recovery fee', 'last_pymnt_d', 'last_pymnt_amnt', 'next_pymnt_d', 'last credit_pull_d', 'last fico range high', 'last fico range low', 'collections 12 mths ex med', 'mths since last major derog', 'policy_code', 'application type', 'annual_inc_joint', 'dti_joint', 'verification status joint', 'acc now deling', 'tot coll amt', 'tot_cur_bal', 'open_acc_6m', 'open_act_il', 'open_il_12m',
'open_il_24m', 'mths since rcnt il', 'total bal il', 'il util', 'open_rv_12m', 'open_rv_24m', 'max_bal_bc', 'all util', 'total_rev_hi_lim', 'inq_fi', 'total cu tl', 'inq last 12m', 'acc_open_past_24mths', 'avg_cur_bal', 'bc_open_to_buy', 'bc util', 'chargeoff within 12 mths', 'delinq_amnt', 'mo sin old il acct', 'mo sin old rev tl op', 'mo_sin_rcnt_rev_tl_op', 'mo_sin_rcnt_tl', 'mort_acc', 'mths since recent bc', 'mths since_recent_bc_dlq', 'mths_since_recent_inq', 'mths since recent revol deling', 'num accts ever 120 pd', 'num actv bc tl', 'num_actv_rev_tl', 'num_bc_sats', 'num_bc_tl', 'num il tl', 'num_op_rev_tl', 'num rev accts', 'num rev_tl_bal_gt_0', 'num_sats', 'num_tl_120dpd_2m', 'num_tl_30dpd', 'num_tl_90g_dpd_24m', 'num_tl_op_past_12m', 'pct_tl_nvr_dlq', 'percent bc qt 75', 'pub rec_bankruptcies', 'tax_liens', 'tot_hi_cred_lim', 'total_bal_ex_mort', 'total_bc_limit',
'total_il_high_credit_limit', 'revol_bal_joint', 'sec_app_fico_range_low', 'sec app fico range high', 'sec_app_earliest_cr_line', 'sec_app_inq_last_6mths', 'sec_app_mort_acc' 'sec app open acc', 'sec_app_revol_util',
'sec_app_open_act_il' 'sec app num rev accts', 'sec app chargeoff within 12 mths', 'sec_app_collections_12_mths_ex_med', 'sec_app_mths_since_last_major_derog', 'hardship_flag', 'hardship type', 'hardship_reason' 'hardship status', 'deferral term', 'hardship amount', 'hardship_start_date', 'hardship_end_date', 'payment_plan_start_date', 'hardship_length', 'hardship_dpd', 'hardship loan status', 'orig projected additional accrued interest', 'hardship payoff balance amount', 'hardship_last_payment_amount', 'disbursement_method', 'debt_settlement_flag', 'debt settlement flag date', 'settlement_status', 'settlement date', 'settlement amount', 'settlement percentage', 'settlement term'] list(reject df) Out[7]: ['Amount Requested', 'Application Date', 'Loan Title', 'Risk_Score', 'Debt-To-Income Ratio', 'Zip Code', 'State', 'Employment Length', 'Policy Code'] print(reject_df['Loan Title'].unique()) print('----') print(accept_df['purpose'].unique()) print('----') print(accept_df['title'].unique()) ['Wedding Covered but No Honeymoon' 'Consolidating Debt' 'Want to consolidate my debt' ... 'dougie03' 'freeup' 'Business Advertising Loan'] -----['debt consolidation' 'small business' 'home improvement' 'major purchase' 'credit card' 'other' 'house' 'vacation' 'car' 'medical' 'moving' 'renewable_energy' 'wedding' 'educational' nan] ['Debt consolidation' 'Business' nan ... 'takeitaway' 'Creditt Card Loan' 'debt reduction/hone updates'] #needed columns from the accepted data clean_accept_df = accept_df[['loan_amnt', 'int_rate', 'issue d', 'fico_range_high', 'dti', 'zip_code', 'addr_state', 'emp_length', 'purpose']] #columns from the rejected data clean reject df = reject df[['Amount Requested', 'Application Date', Risk Score', 'Debt-To-Income Ratio', 'Zip Code', 'State', 'Employment Length', 'Loan Title']] clean_accept_df.shape (2260701, 9) clean_reject_df.shape (27648741, 8)#added target value to the accepted data clean accept df['Loan Status'] = "Accepted" clean_accept_df.head() <ipython-input-13-841988178c2f>:2: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value instead See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/use r_guide/indexing.html#returning-a-view-versus-a-copy clean accept df['Loan Status'] = "Accepted" loan_amnt int_rate issue_d fico_range_high dti zip_code addr_state emp_length purpose Dec-0 3600.0 13.99 679.0 190xx 5.91 PΑ 10+ years debt_consolidation 2015 Dec-1 24700.0 11.99 719.0 16.06 577xx SD small_business 10+ years 2015 Dec-2 20000.0 699.0 home_improvement 10.78 10.78 605xx IL 10+ years 2015 Dec-35000.0 14.85 3 789.0 17.06 076xx NJ 10+ years debt_consolidation 2015 Dec-4 10400.0 22.45 699.0 25.37 174xx PΑ 3 years major_purchase 2015 In [14]: # adding a target value to the rejected data clean_reject_df['Loan Status'] = "Rejected" clean reject df.head() Out[14]: Debt-To-**Employment Application** Zip Amount Loan Risk_Score Income State **Loan Title** Requested Date Code Length **Status** Ratio Wedding 2007-05-0 1000.0 693.0 481xx 10% NM 4 years Covered but No Rejected 26 Honeymoon 2007-05-Consolidating 1000.0 < 1 year 1 703.0 10% 010xx MA Rejected 26 Debt Want to 2007-05-2 11000.0 715.0 consolidate my 10% 212xx MD 1 year Rejected 27 debt 2007-05-3 6000.0 017xx Rejected 698.0 38.64% waksman MA < 1 year 27 2007-05-9.43% 209xx 4 1500.0 509.0 mdrigo Rejected < 1 year 27 clean accept df['int rate'].describe() #interest rates in accept data (min 5.321%, man Out[15]: count 2.260668e+06 1.309283e+01 4.832138e+00 std 5.310000e+00 min 25% 9.490000e+00 50% 1.262000e+01 1.599000e+01 75% 3.099000e+01 max Name: int rate, dtype: float64 #histogram correlation b/t fico score and Risk Score plt.figure(figsize=(10,6)) clean accept df["fico range high"].hist(alpha=.5, color='green', bins=25), clean rejection plt.show() 2.5 2.0 1.5 1.0 0.5 0.0 400 600 800 1000 200 #correlation b/t Risk Score and fico plt.figure(figsize=(10,6)) clean reject df["Risk Score"].hist(alpha=.5, color='red', bins=25), clean accept df[": plt.show() le6 2.5 2.0 1.5 1.0 0.5 0.0 200 400 600 800 1000 sns.jointplot(x="fico_range_high", y="loan_amnt", data=clean_accept_df, color='purple Out[18]: <seaborn.axisgrid.JointGrid at 0x291dc103e50> 40000 35000 30000 25000 loan amnt 20000 15000 10000 5000 700 800 650 750 850 fico_range_high sns.jointplot(x="fico_range_high", y="int_rate", data=clean_accept_df, color='purple' Out[19]: <seaborn.axisgrid.JointGrid at 0x29197979d60> 30 25 20 int_rate 15 10 650 800 700 750 850 fico_range_high print(clean_accept_df.info()) <class 'pandas.core.frame.DataFrame'> RangeIndex: 2260701 entries, 0 to 2260700 Data columns (total 10 columns): Column Dtype 0 loan amnt float64 int rate float64 1 2 issue d object 3 fico range high float64 float64 zip code object addr state object emp length object purpose object 9 Loan Status object dtypes: float64(4), object(6) memory usage: 172.5+ MB print(clean reject df.info()) <class 'pandas.core.frame.DataFrame'> RangeIndex: 27648741 entries, 0 to 27648740 Data columns (total 9 columns): Column Dtype 0 Amount Requested float64 1 Application Date object 2 Risk Score float64 3 Debt-To-Income Ratio object Zip Code object 5 object Employment Length 6 object 7 Loan Title object Loan Status object dtypes: float64(2), object(7) memory usage: 1.9+ GB None #accepted issue date as date clean_accept_df['issue_d'] = pd.to_datetime(clean_accept_df['issue_d']) clean accept df.head() <ipython-input-22-51609ae8d551>:2: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row_indexer,col_indexer] = value instead See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/use r guide/indexing.html#returning-a-view-versus-a-copy clean_accept_df['issue_d'] = pd.to_datetime(clean_accept_df['issue_d']) loan_amnt int_rate issue_d fico_range_high dti zip_code addr_state emp_length purpose 2015-190xx 0 3600.0 13.99 679.0 5.91 PA 10+ years debt_consolidation 12-01 2015-24700.0 11.99 719.0 16.06 577xx SD 10+ years small_business 12-01 2015-2 20000.0 10.78 10.78 699.0 605xx IL 10+ years home_improvement 12-01 2015-3 35000.0 14.85 789.0 17.06 10+ years debt_consolidation 076xx 12-01 2015-4 10400.0 22.45 699.0 25.37 PΑ 174xx 3 years major_purchase 12-01 #2018 accept data accept filtered df = clean accept df[clean accept df['issue d'].dt.year==2018] accept filtered df.head() loan_amnt int_rate issue_d fico_range_high dti zip_code addr_state emp_length pur 2018-421097 5000.0 20.39 669.0 21.80 740xx OK 8 years 03-01 2018-421098 15000.0 9.92 704.0 18.29 337xx FL 2 years debt_consolic 03-01 2018-30.79 669.0 43.97 421099 11200.0 030xx NH < 1 year me 03-01 2018-12.89 421100 25000.0 21.85 669.0 361xx ΑL 10+ years debt_consolic 03-01 2018-421101 3000.0 7.34 764.0 0.58 988xx WA 9 years major_purc 03-01 #see the array accept filtered df ['issue d'].dt.year.unique() Out[51]: array([2018], dtype=int64) #rejected application date as date clean_reject_df['Application Date'] = pd.to_datetime(clean_reject_df['Application Date') clean_reject_df.head() Debt-To-**Employment** Amount Application Zip Loan State Risk_Score Income **Loan Title** Requested Date Code Length **Status** Ratio Wedding 2007-05-1000.0 0 693.0 10% 481xx NM Covered but No Rejected 4 years 26 Honeymoon Consolidating 2007-05-703.0 1 1000.0 10% 010xx Rejected MA < 1 year 26 Debt Want to 2007-05-212xx 2 11000.0 715.0 10% consolidate my MD 1 year Rejected 27 debt 2007-05-38.64% 3 6000.0 698.0 017xx MA < 1 year waksman Rejected 27 2007-05-4 1500.0 < 1 year 509.0 9.43% 209xx MD mdrigo Rejected 27 #2018 reject data reject filtered df = clean reject df[clean reject df['Application Date'].dt.year==2018 reject filtered df.head() Debt-Amount Application To-Zip **Employment** Loan State **Loan Title** Risk_Score Requested Date Income Code Length Status Ratio Debt 4404427 3000.0 2018-07-01 NaN 100% 925xx CA Rejected < 1 year consolidation Major 4404428 40000.0 2018-07-01 NaN 7.45% 335xx FL < 1 year Rejected purchase Debt 4404429 16000.0 2018-07-01 34.93% PΑ NaN 156xx < 1 year Rejected consolidation Debt 40000.0 2018-07-01 957xx 4404430 NaN 27.87% CA < 1 year Rejected consolidation **Business** Rejected 4404431 300000.0 2018-07-01 NaN -1% 258xx TN < 1 vear Loan In [54]: #see the array reject filtered df ['Application Date'].dt.year.unique() array([2018], dtype=int64) Out[54]: reject filtered df.rename(columns=({"Amount Requested": "loan amt", "Application Date": "issue d", "Risk Score" reject filtered df.head() C:\Users\glori\anaconda3\lib\site-packages\pandas\core\frame.py:4296: SettingWithCopyW arning: A value is trying to be set on a copy of a slice from a DataFrame See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/use r guide/indexing.html#returning-a-view-versus-a-copy return super().rename(L dti zip_code addr_state emp_length loan_amt issue_d fico_range_high purpose Sta 2018-Debt 4404427 3000.0 NaN 100% 925xx CA < 1 year Rejec 07-01 consolidation 2018-Major 4404428 40000.0 7.45% 335xx FL NaN Rejec < 1 year 07-01 purchase 2018-Debt 4404429 16000.0 NaN 34.93% 156xx PΑ < 1 year Rejec consolidation 07-01 2018-Debt 4404430 40000.0 27.87% CA NaN 957xx Rejec < 1 year consolidation 07-01 2018-**Business** 4404431 300000.0 Rejec NaN -1% 258xx ΤN < 1 year 07-01 accept filtered df.head() loan_amnt int_rate issue_d fico_range_high dti zip_code addr_state emp_length pur 2018-421097 5000.0 20.39 669.0 21.80 740xx OK 8 years 03-01 2018-704.0 18.29 421098 15000.0 9.92 337xx FL 2 years debt_consolic 03-01 2018-421099 11200.0 30.79 669.0 43.97 030xx NΗ < 1 year mє 03-01 2018-421100 25000.0 21.85 12.89 361xx 669.0 AL 10+ years debt_consolic 03-01 2018-421101 3000.0 7.34 764.0 0.58 988xx WA 9 years major_purc 03-01 #categorical data to numerical data (accept) accept filtered df.loc[(accept filtered df.emp length == '< 1 year'), "emp length"] = accept_filtered_df.loc[(accept_filtered_df.emp_length == '1 year'), "emp_length"] = 1 accept_filtered_df.loc[(accept_filtered_df.emp_length == '2 years'), "emp length"] = accept_filtered_df.loc[(accept_filtered_df.emp_length == '3 years'), "emp length"] = accept filtered df.loc[(accept filtered df.emp length == '4 years'), "emp length"] = accept_filtered_df.loc[(accept_filtered_df.emp_length == '5 years'), "emp length"] = accept_filtered_df.loc[(accept_filtered_df.emp_length == '6 years'), "emp length"] = accept filtered df.loc[(accept filtered df.emp length == '7 years'), "emp length"] accept_filtered_df.loc[(accept_filtered_df.emp_length == '8 years'), "emp length"] = accept filtered df.loc[(accept filtered df.emp length == '9 years'), "emp length"] = accept filtered df.loc[(accept filtered df.emp length == '10+ years'), "emp length"] C:\Users\glori\anaconda3\lib\site-packages\pandas\core\indexing.py:1765: SettingWithCo pyWarning: A value is trying to be set on a copy of a slice from a DataFrame. Try using .loc[row indexer,col indexer] = value instead See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/use r_guide/indexing.html#returning-a-view-versus-a-copy isetter(loc, value) #categorical data to numerical data (reject) reject_filtered_df.loc[(reject_filtered_df.emp_length == '< 1 year'), "emp_length"] = reject_filtered_df.loc[(reject_filtered_df.emp_length == '1 year'), "emp_length"] = 1 reject_filtered_df.loc[(reject_filtered_df.emp_length == '2 years'), "emp_length"] = reject_filtered_df.loc[(reject_filtered_df.emp_length == '3 years'), "emp length"] = "emp length"] = reject_filtered_df.loc[(reject_filtered_df.emp_length == '4 years'), reject_filtered_df.loc[(reject_filtered_df.emp_length == '5 years'), "emp length"] = "emp length"] = reject_filtered_df.loc[(reject_filtered_df.emp_length == '6 years'), reject_filtered_df.loc[(reject_filtered_df.emp_length == '7 years'), "emp length"] reject_filtered_df.loc[(reject_filtered_df.emp_length == '8 years'), reject_filtered_df.loc[(reject_filtered_df.emp_length == '9 years'), "emp_length"] reject filtered df.loc[(reject filtered df.emp length == '10+ years'), "emp length"] final_accept_df = accept_filtered_df.drop(['zip_code', 'issue_d'], axis = 1) final_reject_df = reject_filtered_df.drop(['zip_code', 'issue_d'], axis = 1) print(final accept df.info()) <class 'pandas.core.frame.DataFrame'> Int64Index: 495242 entries, 421097 to 1611876 Data columns (total 8 columns): # Column Non-Null Count Dtype --------loan_amnt 495242 non-null float64 int_rate 495242 non-null float64 fico_range_high 495242 non-null float64 dti 494110 non-null float64 0 dti 494110 non-null float64 addr_state 495242 non-null object emp_length 453255 non-null object purpose 495242 non-null object Loan Status 495242 non-null object dtypes: float64(4), object(4) memory usage: 34.0+ MB None print(final reject df.info()) <class 'pandas.core.frame.DataFrame'> Int64Index: 9496782 entries, 4404427 to 19699075 Data columns (total 7 columns): # Column Dtype 0 loan amt float64 1 fico range high float64 object object dti 3 addr_state emp_length object purpose object Loan Status object 5 6 dtypes: float64(2), object(5) memory usage: 579.6+ MB None In [64]: #drop na values in accept df final_accept_df = final_accept_df.dropna() final_accept_df.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 453176 entries, 421097 to 1611876 Data columns (total 8 columns): # Column Non-Null Count Dtype -----fico_range_high 453176 non-null float64 dti 453176 non-null float64
addr_state 453176 non-null object
emp_length 453176 non-null object 5 purpose 453176 non-null object Loan Status 453176 non-null object purpose 6 dtypes: float64(4), object(4) memory usage: 31.1+ MB #drop na values in reject df final reject df = final_reject_df.dropna() final reject df.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 628344 entries, 4404458 to 19699073 Data columns (total 7 columns): Non-Null Count Dtype # Column ----loan amt 628344 non-null float64 0 1 fico_range_high 628344 non-null float64 628344 non-null object 628344 non-null object dti 2 addr state 3 object emp_length 628344 non-null 5 628344 non-null purpose object Loan Status object 628344 non-null dtypes: float64(2), object(5) memory usage: 38.4+ MB filenames = [final accept df, final reject df] combined data = pd.concat(filenames) combined_data.head() Loan lo loan_amnt int_rate fico_range_high dti addr_state emp_length purpose **Status** 421097 5000.0 20.39 669.0 21.8 OK other Accepted 704.0 421098 15000.0 9.92 18.29 FL debt_consolidation Accepted 421099 11200.0 30.79 669.0 43.97 NH medical Accepted 421100 25000.0 21.85 669.0 12.89 ΑL debt_consolidation Accepted 421101 3000.0 7.34 764.0 0.58 WA 9 major_purchase Accepted combined data.to csv('Data/model 1 combine data.csv') In [79]: combined data.shape (1081520, 9)accept filtered df.to csv('Data/model 1 accept data.csv') In [74]: accept filtered df.head(5) Out[74]: loan_amnt int_rate issue_d fico_range_high dti zip_code addr_state emp_length pur 2018-669.0 21.80 421097 5000.0 20.39 740xx OK 8 03-01 2018-421098 15000.0 9.92 704.0 18.29 debt_consolic 337xx FL 03-01 2018-421099 11200.0 30.79 669.0 43.97 030xx NH 0 me 03-01 2018-421100 25000.0 21.85 669.0 12.89 361xx ΑL debt_consolic 03-01 2018-421101 3000.0 7.34 0.58 988xx 764.0 WA 9 major_purc 03-01 accept filtered df.shape (495242, 10)reject_filtered_df.to_csv('Data/model_1_reject_data.csv') reject filtered df.head(5) L loan_amt issue_d fico_range_high dti zip_code addr_state emp_length purpose Sta 2018-Debt 4404427 3000.0 NaN 100% 925xx CA Rejec 07-01 consolidation Major 7.45% 4404428 40000.0 NaN 335xx Rejec 07-01 purchase 2018-Debt 4404429 16000.0 NaN 34.93% 156xx PΑ Rejec consolidation 07-01 2018-Debt 4404430 40000.0 27.87% Rejec NaN 957xx CA consolidation 07-01 2018-**Business 4404431** 300000.0 NaN -1% 258xx ΤN Rejec 07-01 Loan In [78]: reject filtered df.shape (9496782, 9) Out[78]: !pip install sklearn --upgrade