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 gt 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', '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, 8) 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-13-5e810bf9972e>: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" Loar loan\_amnt issue\_d fico\_range\_high dti zip\_code addr\_state emp\_length purpose Status Dec-0 3600.0 679.0 5.91 190xx PΑ 10+ years debt\_consolidation Accepted 2015 Dec-1 24700.0 719.0 16.06 577xx SD 10+ years small\_business Accepted 2015 Dec-2 20000.0 699.0 10.78 605xx IL 10+ years home\_improvement 2015 Dec-3 35000.0 789.0 17.06 076xx debt\_consolidation NJ 10+ years Accepted 2015 Dec-4 699.0 10400.0 25.37 174xx PΑ 3 years major\_purchase Accepted 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-**Application Employment** Amount Zip Loan Risk\_Score Income State **Loan Title** Requested Date Code Length **Status** Ratio Wedding 2007-05-0 1000.0 10% 693.0 481xx NM Covered but No Rejected 4 years 26 Honeymoon 2007-05-Consolidating 1000.0 1 703.0 10% 010xx MA < 1 year Rejected 26 Debt Want to 2007-05-2 11000.0 715.0 10% 212xx MD consolidate my 1 year Rejected 27 debt 2007-05-3 6000.0 698.0 38.64% 017xx Rejected MA < 1 year waksman 27 2007-05-9.43% 209xx 4 1500.0 509.0 MD < 1 year mdrigo Rejected 27 print(clean accept df.info()) <class 'pandas.core.frame.DataFrame'> RangeIndex: 2260701 entries, 0 to 2260700 Data columns (total 9 columns): Column Dtype 0 loan amnt float64 issue d 1 object fico range high float64 float64 3 4 zip\_code object 5 addr\_state object 6 emp\_length object purpose 7 object 8 Loan Status object dtypes: float64(3), object(6) memory usage: 155.2+ MB None 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 Risk Score float64 Debt-To-Income Ratio object Zip Code object 5 object 6 Employment Length object Loan Title 7 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-18-4db1f5c7d393>: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']) Loar loan\_amnt issue\_d fico\_range\_high dti zip\_code addr\_state emp\_length purpose **Status** 2015-3600.0 0 5.91 190xx 679.0 PA 10+ years debt\_consolidation Accepted 12-01 2015-24700.0 719.0 16.06 577xx SD 10+ years small\_business Accepted 12-01 2015-20000.0 699.0 10.78 2 605xx IL 10+ years home\_improvement Accepted 12-01 2015-35000.0 789.0 17.06 076xx 10+ years debt\_consolidation Accepted 12-01 2015-4 10400.0 699.0 25.37 PΑ major\_purchase Accepted 174xx 3 years 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 issue\_d fico\_range\_high dti zip\_code addr\_state emp\_length purpose 2018-5000.0 421097 669.0 21.80 740xx OK other Ac 8 years 03-01 2018-421098 15000.0 704.0 18.29 337xx FL 2 years debt\_consolidation 03-01 2018-421099 11200.0 669.0 43.97 030xx NH medical < 1 year Ac 03-01 2018-25000.0 421100 669.0 12.89 361xx AL 10+ years debt\_consolidation 03-01 2018-421101 3000.0 764.0 0.58 988xx WA 9 years major\_purchase Ac 03-01 #see the array accept filtered df ['issue d'].dt.year.unique() 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-**Amount** Application Zip **Employment** Loan Risk\_Score Income State **Loan Title** Requested Date Code Length **Status** Ratio Wedding 2007-05-0 1000.0 693.0 10% 481xx NM Covered but No Rejected 4 years 26 Honeymoon Consolidating 2007-05-1 1000.0 010xx Rejected 703.0 10% MA < 1 year 26 Debt Want to 2007-05-2 11000.0 715.0 10% 212xx MD 1 year consolidate my Rejected 27 debt 2007-05-3 6000.0 698.0 38.64% 017xx MA Reiected < 1 year waksman 27 2007-05-4 1500.0 509.0 9.43% 209xx MD < 1 year mdrigo Rejected 27 #2018 reject data reject filtered df = clean reject df[clean reject df['Application Date'].dt.year == 201 reject filtered df.head() Debt-**Amount** Application To-Zip **Employment** Loan Risk\_Score State Loan Title 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 Rejected < 1 year purchase Debt 4404429 16000.0 2018-07-01 34.93% 156xx PΑ Rejected NaN < 1 year consolidation Debt 4404430 40000.0 2018-07-01 NaN 27.87% 957xx CA Rejected < 1 year consolidation **Business** 4404431 300000.0 2018-07-01 NaN -1% 258xx TN < 1 year Rejected Loan #see the array reject\_filtered\_df ['Application Date'].dt.year.unique() Out[23]: array([2018], dtype=int64) accept\_filtered\_df = accept\_filtered\_df.rename(columns={"loan\_amnt": "Amount\_Requested") "issue\_d": "Application\_Date", "fico\_range\_high": "Risk\_Score", "dti": "debt to income ratio", "addr state": "State", "emp length": "Employment\_Length", "purpose": "Purpose", "Loan Status": "Loan Status", }) reject filtered df = reject filtered df.rename(columns={"Amount Requested": "Amount Re "Application Date": "Application Date", "Risk Score": "Risk Score", "Debt-To-Income Ratio": "debt to income ratio", "Zip Code": "zip code", "Employment Length": "Employment Length", "Loan Title": "Purpose", "Loan Status": "Loan Status" accept filtered df.head() Amount\_Requested Application\_Date Risk\_Score debt\_to\_income\_ratio zip\_code State Employmer 421097 5000.0 2018-03-01 669.0 21.80 740xx OK 421098 15000.0 FΙ 2018-03-01 704.0 18.29 337xx 421099 11200.0 2018-03-01 669.0 43.97 030xx NH 421100 25000.0 2018-03-01 669.0 12.89 361xx ΑL 421101 3000.0 764.0 WA 2018-03-01 0.58 988xx reject filtered df.head() Amount\_Requested Application\_Date Risk\_Score debt\_to\_income\_ratio zip\_code State Employm 4404427 3000.0 2018-07-01 100% 925xx NaN CA 4404428 40000.0 2018-07-01 NaN FL 7.45% 335xx 4404429 16000.0 2018-07-01 NaN 34.93% 156xx PΑ 4404430 40000.0 2018-07-01 NaN 27.87% 957xx CA 4404431 300000.0 2018-07-01 NaN -1% 258xx TN #categorical data to numerical data (accept) accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '< 1 year'), "Employment\_Length"> '< 1 year')</pre> accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '1 year'), "Employment accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '2 years'), "Employment\_Length == '2 years'), accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '3 years'), "Employment\_Length == '3 years'), accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '4 years'), "Employment\_Length == '4 years'), accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '5 years'), "Employment\_Length == '5 years'), "Employment\_Length == '5 years'), "Employment\_Length == '5 years') accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '6 years'), "Employment\_Length == '6 years'), accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '7 years'), "Employment\_Length == '7 years'), accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '8 years'), "Employment\_Length == '8 years'), accept\_filtered\_df.loc[(accept\_filtered\_df.Employment\_Length == '9 years'), "Employment\_Length == '9 years'), accept filtered df.loc[(accept filtered df.Employment Length == '10+ years'), "Employment Length" #categorical data to numerical data (reject) reject filtered df.loc[(reject filtered df.Employment Length == '< 1 year'), "Employment Length" reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '1 year'), "Employment reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '2 years'), "Employment\_Length" reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '3 years'), "Employment\_ reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '4 years'), "Employment reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '5 years'), "Employment reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '6 years'), "Employment reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '7 years'), "Employment reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '8 years'), "Employment reject\_filtered\_df.loc[(reject\_filtered\_df.Employment\_Length == '9 years'), "Employment reject filtered df.loc[(reject filtered df.Employment Length == '10+ years'), "Employment Length" print(reject\_filtered\_df["Employment\_Length"].unique()) print(accept filtered df["Employment Length"].unique()) [0 nan 2 1 5 9 10 8 3 4 7 6] [8 2 0 10 9 nan 5 4 1 7 3 6] #fix dti in reject filtered df to match dti in accept df reject\_filtered\_df["debt\_to\_income\_ratio"] = pd.to\_numeric(reject\_filtered\_df["debt\_to\_ reject filtered df Amount\_Requested Application\_Date Risk\_Score debt\_to\_income\_ratio zip\_code State Employn 4404427 3000.0 2018-07-01 NaN 100.00 925xx CA 335xx 4404428 40000.0 2018-07-01 NaN 7.45 FL 4404429 16000.0 2018-07-01 NaN 34.93 PΑ 156xx 4404430 40000.0 2018-07-01 NaN 27.87 957xx CA 258xx 4404431 300000.0 2018-07-01 NaN -1.00 TN 19699071 17000.0 2018-06-30 26.60 NaN 301xx GA 19699072 2000.0 2018-06-30 NaN 0.00 117xx NY 19699073 2500.0 2018-06-30 567.0 0.00 366xx 19699074 4500.0 2018-06-30 NaN 4.65 780xx TΧ 19699075 2500.0 2018-06-30 7.33 951xx CA NaN 9496782 rows × 9 columns final\_accept\_df = accept\_filtered\_df.drop(['zip\_code', 'Application\_Date'], axis = 1) final\_reject\_df = reject\_filtered\_df.drop(['zip\_code', 'Application Date'], axis = 1) print(final accept df.info()) <class 'pandas.core.frame.DataFrame'> Int64Index: 495242 entries, 421097 to 1611876 Data columns (total 7 columns): Column Non-Null Count Dtype \_\_\_ Amount\_Requested 495242 non-null float64 0 Risk Score 495242 non-null float64 1 debt\_to\_income\_ratio 494110 non-null float64 3 495242 non-null object 453255 non-null object Employment\_Length 4 Purpose 495242 non-null 5 object Loan Status 495242 non-null object dtypes: float64(3), object(4) memory usage: 30.2+ MB None In [34]: print(final reject df.info()) <class 'pandas.core.frame.DataFrame'> Int64Index: 9496782 entries, 4404427 to 19699075 Data columns (total 7 columns): # Column Amount\_Requested float64 0 float64 Risk Score 1 debt\_to\_income\_ratio float64 3 4 Employment\_Length object Purpose object object 5 6 Loan Status object dtypes: float64(3), object(4) memory usage: 579.6+ MB None #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 7 columns): Non-Null Count Dtype # Column ----------Amount\_Requested 453176 non-null float64 Risk\_Score 453176 non-null float64 0 Risk\_Score 453176 non-null float64 debt\_to\_income\_ratio 453176 non-null float64 State 453176 non-null object Employment\_Length 453176 non-null object 1 453176 non-null object Purpose 4531/6 non-null object Loan Status dtypes: float64(3), object(4) memory usage: 27.7+ 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): # Column Non-Null Count Dtype Amount\_Requested 628344 non-null float64 Risk\_Score 628344 non-null float64 0 debt\_to\_income\_ratio 628344 non-null float64 628344 non-null object Employment\_Length 628344 non-null object Purpose 628344 non-null object 5 628344 non-null object 6 Loan\_Status dtypes: float64(3), object(4) memory usage: 38.4+ MB filenames = [final accept df, final reject df] combined data = pd.concat(filenames) combined data["Employment Length"].astype(str).astype(float) Out[38]: 421097 421098 2.0 421099 0.0 421100 10.0 421101 9.0 19698926 5.0 0.0 19698936 19698996 0.0 19699052 19699073 Name: Employment\_Length, Length: 1081520, dtype: float64 combined data Amount\_Requested Risk\_Score debt\_to\_income\_ratio State Employment\_Length Purpo: 421097 5000.0 669.0 21.80 8 OK oth 15000.0 704.0 421098 18.29 FL 2 debt\_consolidation 421099 11200.0 669.0 43.97 NH medic 12.89 421100 25000.0 669.0 ΑL debt\_consolidation 421101 3000.0 764.0 0.58 WA 9 major\_purcha 535.0 19698926 3000.0 1.54 CT 5 Medical expens 19698936 6000.0 580.0 5.55 TX 0 Car financii 19698996 Debt consolidation 3000.0 561.0 0.00 OR 19699052 2500.0 573.0 8.55 NY Oth 19699073 2500.0 567.0 0.00 AL 2 Car financii 1081520 rows × 7 columns In [40]: combined data.shape (1081520, 7)Out[40]: In [41]: combined data.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 1081520 entries, 421097 to 19699073 Data columns (total 7 columns): Column Non-Null Count Dtype 1081520 non-null float64 Amount Requested 0 1081520 non-null float64 Risk Score 1 debt to income ratio 1081520 non-null float64 3 1081520 non-null object 1081520 non-null object 4 Employment Length 1081520 non-null object 5 Purpose Loan Status 1081520 non-null object dtypes: float64(3), object(4) memory usage: 66.0+ MB In [42]: combined data.to csv('Data/model 1 combine data.csv') In [43]: !pip install sklearn --upgrade Requirement already up-to-date: sklearn in c:\users\glori\anaconda3\lib\site-packages Requirement already satisfied, skipping upgrade: scikit-learn in c:\users\glori\anacon da3\lib\site-packages (from sklearn) (0.23.2) Requirement already satisfied, skipping upgrade: scipy>=0.19.1 in c:\users\glori\anaco  $\verb|nda3|lib|site-packages| (from scikit-learn->sklearn)| (1.5.2)$ Requirement already satisfied, skipping upgrade: joblib>=0.11 in c:\users\glori\anacon da3\lib\site-packages (from scikit-learn->sklearn) (0.17.0) Requirement already satisfied, skipping upgrade: numpy>=1.13.3 in c:\users\glori\anaco nda3\lib\site-packages (from scikit-learn->sklearn) (1.19.2) Requirement already satisfied, skipping upgrade: threadpoolctl>=2.0.0 in c:\users\glor i\anaconda3\lib\site-packages (from scikit-learn->sklearn) (2.1.0)