



Guidelines for PGPDSE FT Capstone Project – Interim Report

Project Group Info:

BATCH DETAILS	DSE online January 2022
TEAM MEMBERS	<ul style="list-style-type: none">• Sree Ganesha Chellappa• Faishal Sharif• Hasitha vaddi• Nikhil Nivalagi• Ashish Shukla.
DOMAIN OF PROJECT	Predictive Analysis in Finance
PROJECT TITLE	Loan Default Prediction
GROUP NUMBER	GROUP – 4
TEAM LEADER	Hasitha vaddi
MENTOR NAME	Mrs Vidhya Kannaiah

Abstract:

The bank Indessa has not been performing well for the last three quarters and would like to improve their bank's performance by reducing their NPA's. Since the data which is collected is very messy, identification of defaulters who are the major cause which increases the NPA's has become difficult. So we are devising a model which can predict whether the person who is a loan applicant or so already has a loan will be a potential defaulter in the future.

Objectives:

Bank Indessa faces challenges in identifying the loan applicants of those who can be potential defaulters. Any discrepancies in data faced by the employees against a set of standards can indicate that the loan applicant can be a potential defaulter and can mark the applicants for any differences against some set standards if they have crossed the set number of markings (parameters which have been found). They can be identified as a potential defaulter and the approval of the loan can be rejected.

The objective of this Capstone Project is to predict the optimal parameters for the identification of potential defaulters of a loan to meet the stockholder confidence in their company and to carry out any future increase of NPA's (non-performing assets) and increase the company's share in the market.

Industry Review:

Commercial lending options provide flexible long-term lending, which is a major commercial lending market driver.

In addition, payment collection collaborations between digital lending organizations and FinTech companies are expected to grow in the market.

The commercial lending market size was valued at USD 8,823.53 Billion in 2020 and is projected to reach USD 29,379.83 Billion by 2030, growing at a CAGR of 13.1% from 2021 to 2030.

Many business owners had to take out commercial loans to keep their businesses afloat as COVID-19 cases continued to rise and more restrictions were

imposed during the pandemic. This resulted in a surge in commercial lending market growth.

Commercial lending offers the lowest interest rates on all loan options, allowing business owners to get needed funds while keeping overhead costs low. Borrowers who choose fixed monthly repayments can use them accurately in their business.

planning and forecasting, allowing them to structure their business finance with a bit more certainty.

Furthermore, commercial lending payment plans are typically for several years, allowing a company to focus on other important business matters such as sales, overhead management, and employee training. As a result, this is a significant driving force in the commercial lending market.

Problem statement

Bank Indessa has not done well in the last 3 quarters. Their NPAs (Non-Performing Assets) has reached an all-time high. It is starting to lose the confidence of its investors. As a result, its stock has fallen by 20% in the previous quarter alone.

After careful analysis, it was found that the majority of NPA was contributed by loan defaulters. With the messy data collected over the years, this bank has decided to use machine learning to figure out a way to find these defaulters and devise a plan to reduce them.

This bank uses a pool of investors to sanction their loans.

We will help this bank by predicting the probability that a member will default.

Project Outcome:

By implementing the resultant models built using the above methods, we can suggest the ideal standards or parameters required on which the loan defaulters can be identified in a short duration to reduce the increase in NPAs and improve efficiency.

Dataset and Domain:

Data Dictionary:

Real-time bank dataset obtained from an organization.

The indessa.xlsx This data set comprises information captured in December 2016.

The dataset has 5,32,428 records and 45 attributes

The attribute/feature/column names are given below:

```
Index(['member_id', 'loan_amnt', 'funded_amnt', 'funded_amnt_inv', 'term',  
      'batch_enrolled', 'int_rate', 'grade', 'sub_grade', 'emp_title', 'emp_length',  
      'home_ownership', 'annual_inc', 'verification_status', 'pymnt_plan', 'desc', 'purpose',  
      'title', 'zip_code', 'addr_state', 'dti', 'delinq_2yrs', 'inq_last_6mths',  
      'mths_since_last_delinq', 'mths_since_last_record', 'open_acc', 'pub_rec',
```

'revol_bal', 'revol_util', 'total_acc', 'initial_list_status',
'total_rec_int', 'total_rec_late_fee', 'recoveries', 'collection_recovery_fee',
'collections_12_mths_ex_med', 'mths_since_last_major_derog',
'application_type', 'verification_status_joint', 'last_week_pay', 'acc_now_delinq',
'tot_coll_amt', 'tot_cur_bal', 'total_rev_hi_lim', 'loan_status'])

There are 27 numerical and 18 object dtypes

dtype - numerical - ['member_id', 'loan_amnt', 'funded_amnt', 'funded_amnt_inv', 'int_rate',
'annual_inc', 'dti', 'delinq_2yrs', 'inq_last_6mths', 'mths_since_last_delinq',
'mths_since_last_record', 'open_acc', 'pub_rec', 'revol_bal', 'revol_util', 'total_acc', 'total_rec_int',
'total_rec_late_fee', 'recoveries', 'collection_recovery_fee', 'collections_12_mths_ex_med',
'mths_since_last_major_derog', 'acc_now_delinq', 'tot_coll_amt', 'tot_cur_bal', 'total_rev_hi_lim',
'loan_status']

dtype - object - ['term', 'batch_enrolled', 'grade', 'sub_grade', 'emp_title', 'emp_length',
'home_ownership', 'verification_status', 'pymnt_plan', 'desc', 'purpose', 'title', 'zip_code',
'addr_state', 'initial_list_status', 'application_type', 'verification_status_joint', 'last_week_pay']

Variable	Description	Dtype
member_id	unique ID assigned to each member	int64
loan_amnt	loan amount (\$) applied by the member	int64
funded_amnt	loan amount (\$) sanctioned by the bank	int64
funded_amnt_inv	loan amount (\$) sanctioned by the investors	float64
term	term of loan (in months)	object
batch_enrolled	batch numbers allotted to members	object

Variable	Description	Dtype
int_rate	interest rate (%) on loan	float64
grade	grade assigned by the bank	object
sub_grade	grade assigned by the bank	object
emp_title	job / Employer title of member	object
emp_length	employment length, where 0 means less than one year and 10 means ten or more years	object
home_ownership	status of home ownership	object
annual_inc	annual income (\$) reported by the member	float64
verification_status	status of income verified by the bank	object
pymnt_plan	indicates if any payment plan has started against loan	object
desc	loan description provided by member	object
purpose	purpose of loan	object
title	loan title provided by member	object

Variable	Description	Dtype
zip_code	first three digits of area zip code of member	object
addr_state	living state of member	object
dti	ratio of member's total monthly debt repayment excluding mortgage divided by self reported monthly income	float64
delinq_2yrs	number of 30+ days delinquency in past 2 years	float64
inq_last_6mths	number of inquiries in last 6 months	float64
mths_since_last_delinq	number of months since last delinq	float64
mths_since_last_record	number of months since last public record	float64
open_acc	number of open credit line in member's credit line	float64
pub_rec	number of derogatory public records	float64
revol_bal	total credit revolving balance	float64
revol_util	amount of credit a member is using relative to revol_bal	float64

Variable	Description	Dtype
total_acc	total number of credit lines available in members credit line	float64
initial_list_status	unique listing status of the loan - W(Waiting), F(Forwarded)	object
total_rec_int	interest received till date	float64
total_rec_late_fee	Late fee received till date	float64
recoveries	post charge off gross recovery	float64
collection_recovery_fee	post charge off collection fee	float64
collections_12_mths_ex_med	number of collections in last 12 months excluding medical collections	float64
mths_since_last_major_derog	months since most recent 90 day or worse rating	float64
application_type	indicates when the member is an individual or joint	object
verification_status_joint	indicates if the joint members income was verified by the bank	object
last_week_pay	indicates how long (in weeks) a member has paid EMI after batch enrolled	object

Variable	Description	Dtype
acc_now_delinq	number of accounts on which the member is delinquent	float64
tot_coll_amt	total collection amount ever owed	float64
tot_cur_bal	total current balance of all accounts	float64
total_rev_hi_lim	total revolving credit limit	float64
loan_status	status of loan amount, 1 = Defaulter, 0 = Non Defaulters	int64

Nulls

BEFORE

member_id	0
loan_amnt	0
funded_amnt	0
funded_amnt_inv	0
term	0
batch_enrolled	85149
int_rate	0
grade	0
sub_grade	0
emp_title	30833

emp_length	26891
home_ownership	0
annual_inc	3
verification_status	0
pymnt_plan	0
desc	456829
purpose	0
title	90
zip_code	0
addr_state	0
dti	0
delinq_2yrs	16
inq_last_6mths	16
mths_since_last_delinq	272554
mths_since_last_record	450305
open_acc	16
pub_rec	16
revol_bal	0
revol_util	287
total_acc	16
initial_list_status	0
total_rec_int	0
total_rec_late_fee	0
recoveries	0
collection_recovery_fee	0
collections_12_mths_ex_med	95
mths_since_last_major_derog	399448
application_type	0
verification_status_joint	532123
last_week_pay	0
acc_now_delinq	16
tot_coll_amt	42004
tot_cur_bal	42004

total_rev_hi_lim	42004
loan_status	0

AFTER

loan_amnt	0
funded_amnt	0
funded_amnt_inv	0
term	0
int_rate	0
emp_length	0
annual_inc	0
verification_status	0
dti	0
delinq_2yrs	0
inq_last_6mths	0
mths_since_last_delinq	0
open_acc	0
pub_rec	0
revol_bal	0
revol_util	0
total_acc	0
initial_list_status	0
total_rec_int	0
total_rec_late_fee	0
recoveries	0
collection_recovery_fee	0
collections_12_mths_ex_med	0
last_week_pay	0
acc_now_delinq	0
tot_coll_amt	0
tot_cur_bal	0
total_rev_hi_lim	0
loan_status	0
grade_num	0
addr_state_NE	0
addr_state_S	0
addr_state_W	0

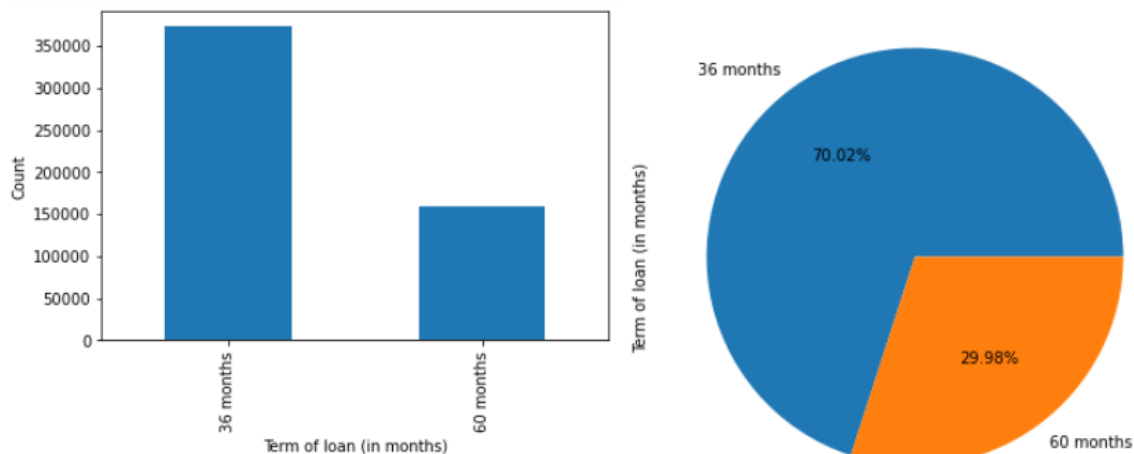
Univariate Analysis:

Analysis of Categorical Variables:

From the below plots, each of the categorical variables is analysed individually.

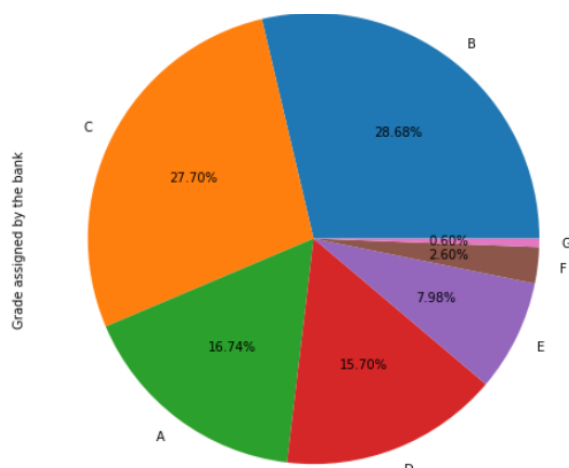
Term of Loan Attribute:

From the above plot, we can observe that 70% of members need the loan for the duration of 36 months and the remaining 30% of members for 60 months.



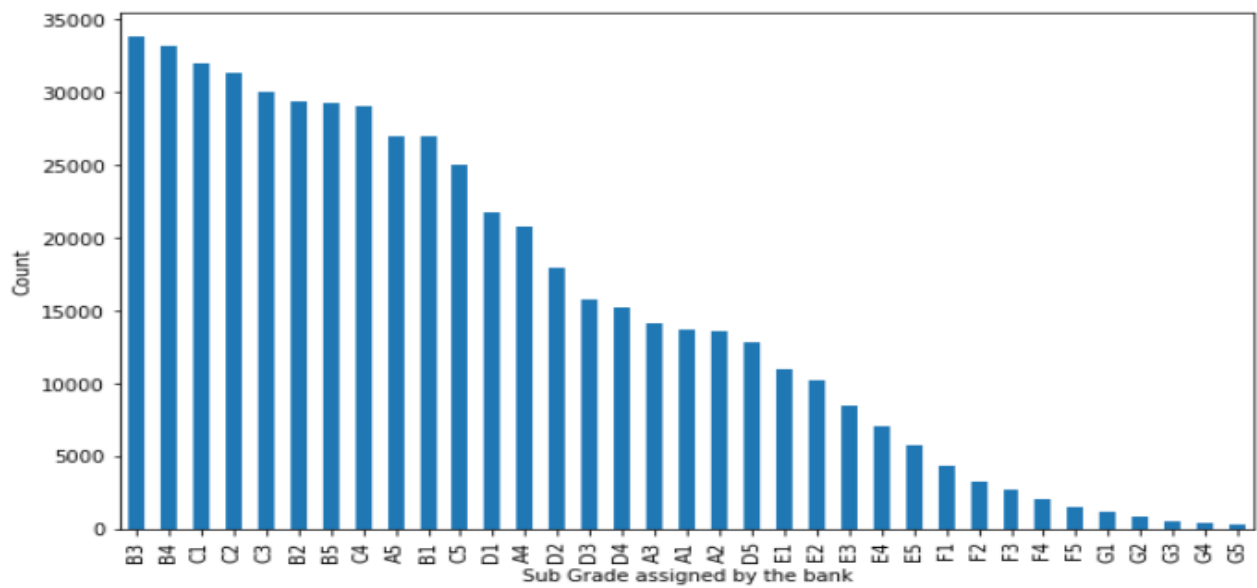
Grade Attribute:

From the above plot, we can observe that the highest no. of members belongs to Grade B and the lowest no. of members from Grade G. Also the no. of members from Grade A and D is approximately the same.

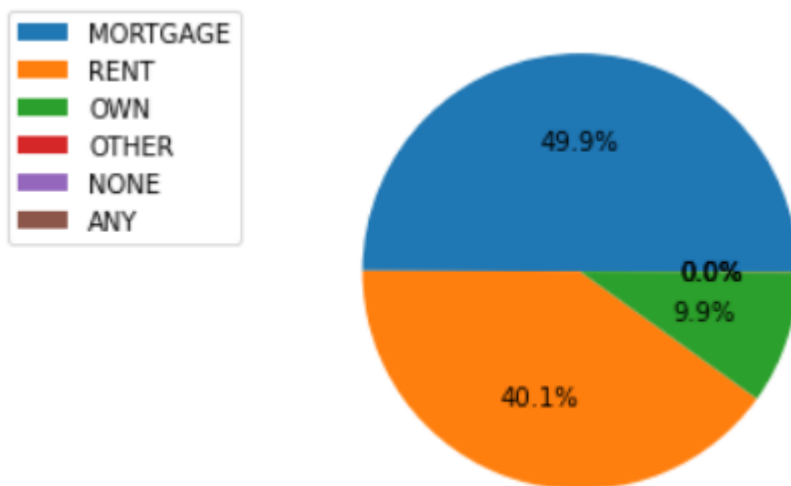


Sub Grade Attribute:

From the above plot, we can observe that the maximum number of members are assigned in Sub Grade B3 and the lowest members are assigned in Sub Grade G5. More than 30,000 members are assigned in these 5 Subgrades B3, B4, C1, C2, and C3.

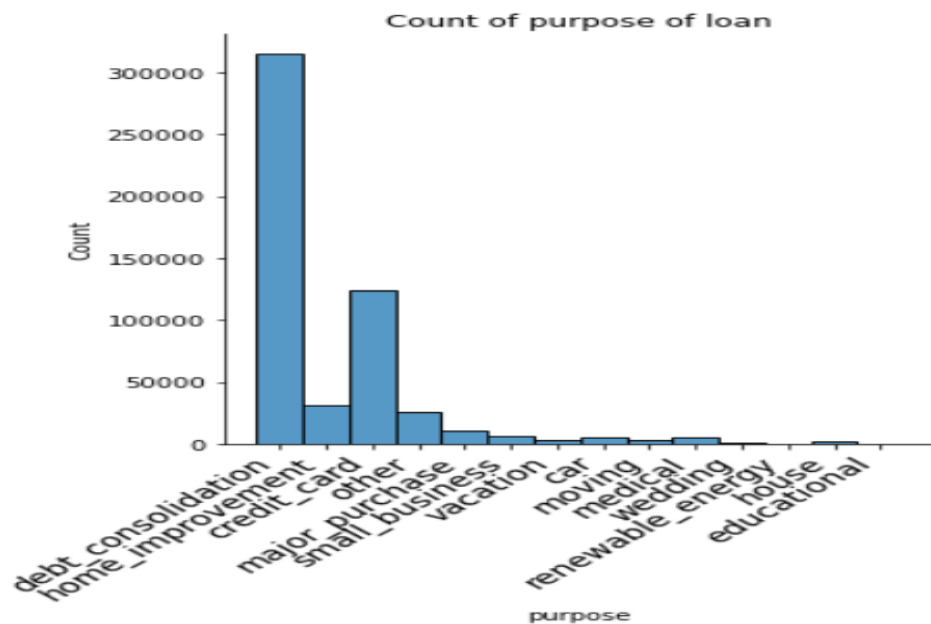


Home Ownership Attribute:



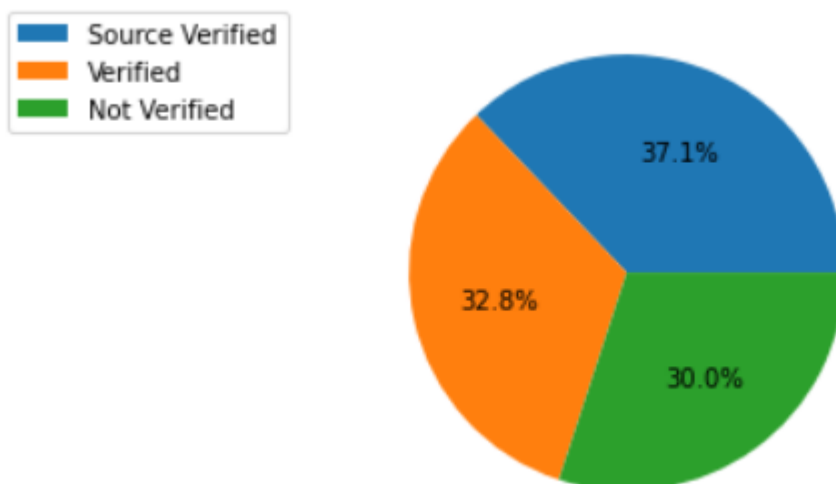
From the above plot, we can observe that almost 50% of home ownership status is of the 'Mortgage' type. However, mortgage, rent, and own totally comprise the 99.99% home ownership status of the total members.

Purpose Attribute:



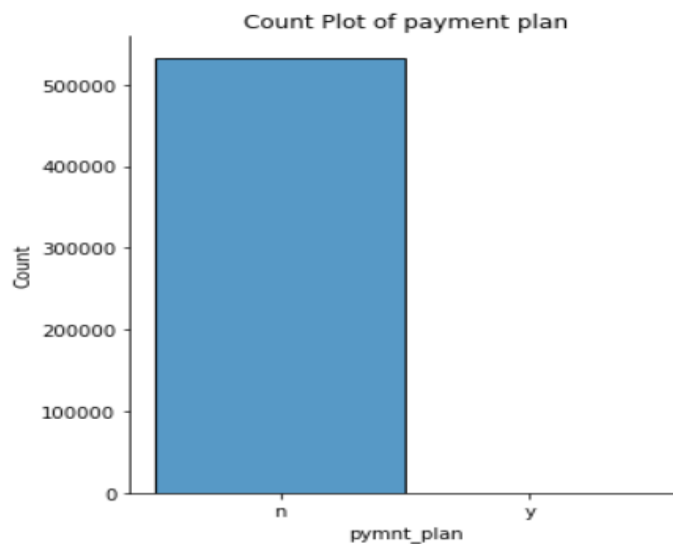
- From the above count plot, it can be observed that the maximum number of members (60%) have the purpose of the loan as 'debt_consolidation'.
- More than 1lakh employees have put their purpose of the loan as 'credit card'.

Verification Status Attribute:



- From the above pie plot, it can be observed that the maximum status of income verified by the bank is 'source verified'. However not verified members are also in significant numbers. All together all the three categories contribute significantly as all three have members present in good numbers.

Payment Plan attribute:



From the above count plot, we can observe that almost all the members don't have any payment plan started against the loan.

last_week_pay: it has too many categories to use

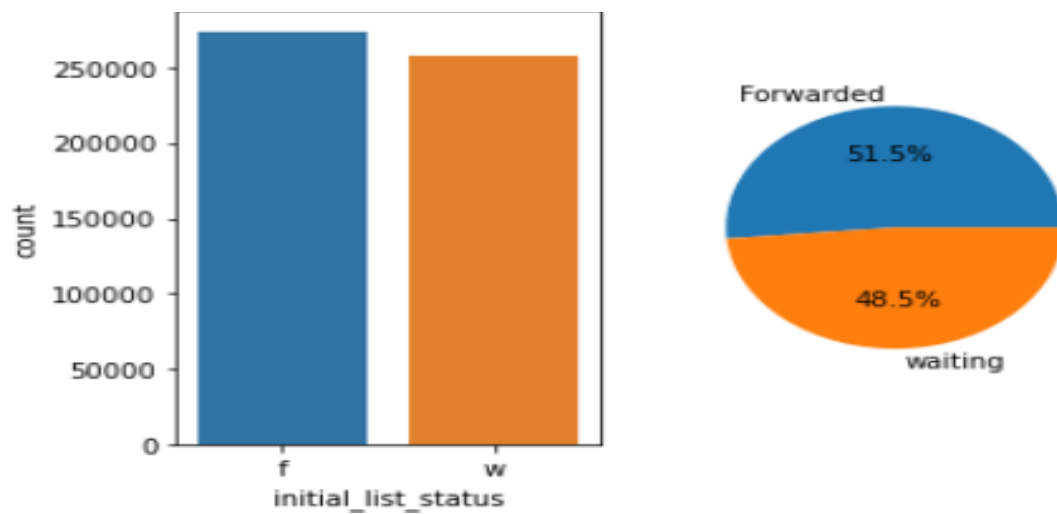
Employee Length Attribute:



- From the above distribution plot we can assume that the distribution is highest at 10 years. That means the max members have an employment length of 10 years.
- Similarly, the above boxplot tells us that the mean employment length of 6 years approximately.
- Majority of the members have employment length in the range of 3 years to 10 years.

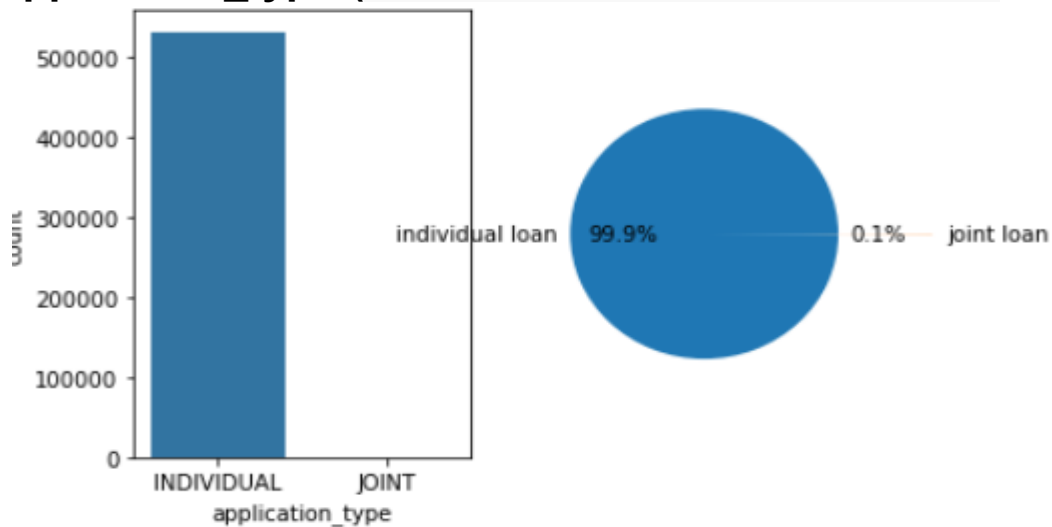
initial_list_status:

unique listing status of the loan - W(Waiting), F(Forwarded)



- forwarded loan status() occupies 51.5 per cent of the total data count
- waiting loan status() occupies 48.5 per cent of the total data count
- unique listing status of the loan - W(Waiting), F(Forwarded)
- forwarded loan status() (274018) is slightly greater than waiting for status ()

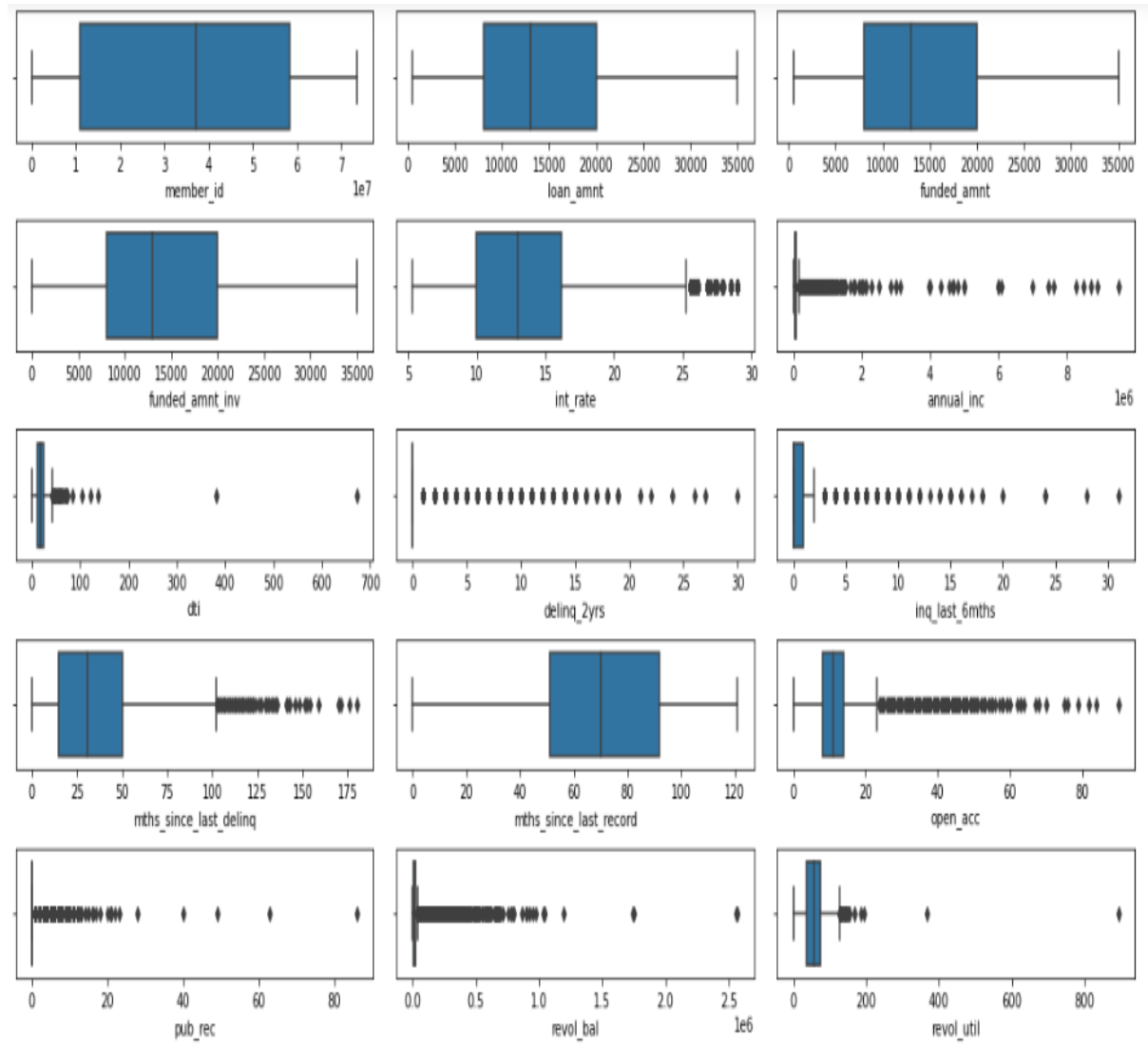
Application_type: (indicates when the member is an individual or joint)



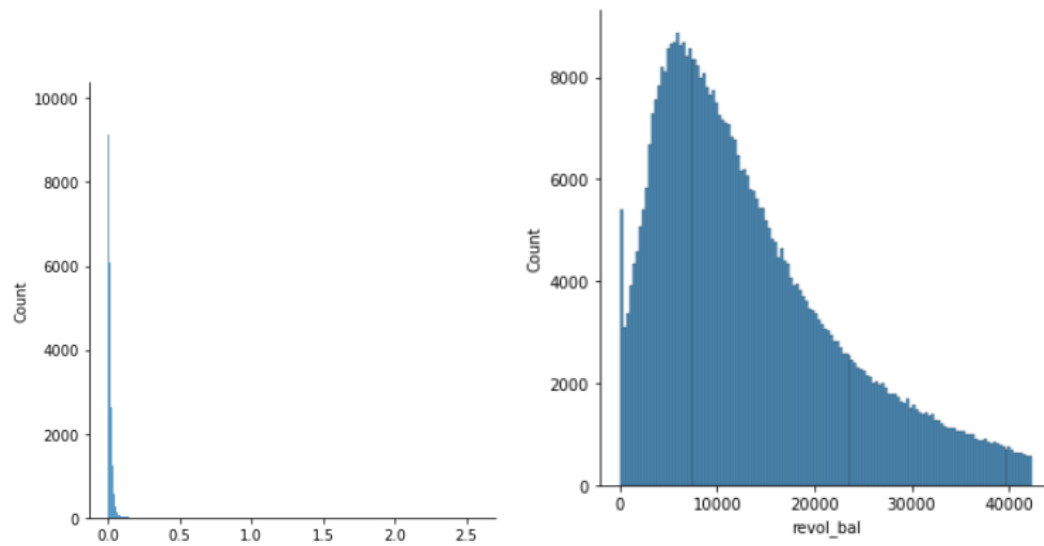
- Joint loan takers hold a minority of 0.1 per cent and the remaining i.e. individual loan are 99.9 per cent

Numerical data

except for the member id , loan amount,funded amount,funded amount invested,and months since last record the rest of the data has outliers as seen in the boxplots below

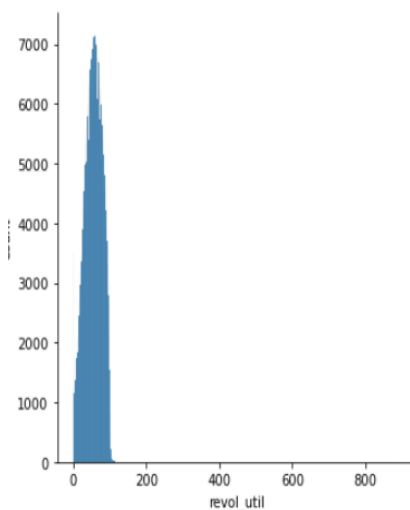


revol_bal: If you don't pay the balance on your revolving credit account in full every month, the unpaid portion carries over to the next month.



- The data has an extreme outlier of value of 2.5×10^5 which affect the mean, and other outliers
- The majority of the data has a total credit revolving balance in the range of 0 to 40000 and is positively skewed

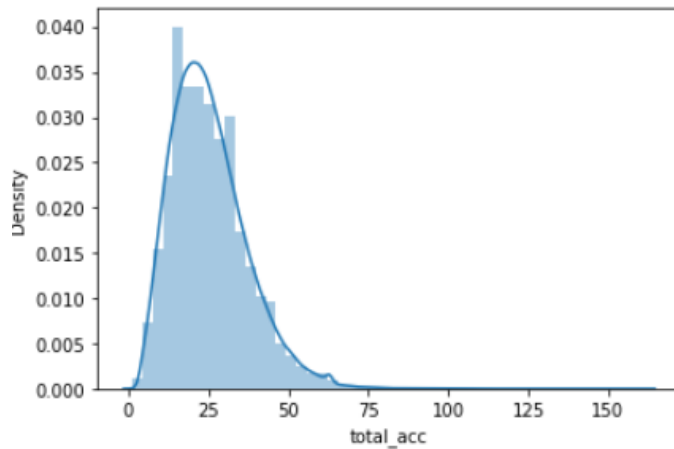
revol_util: (amount of credit a member is using relative to revol_bal)



After removal of the outliers the data is slightly skewed, but majority of the credit usages lies between 0 to 100k where the mean credit usage is 65k

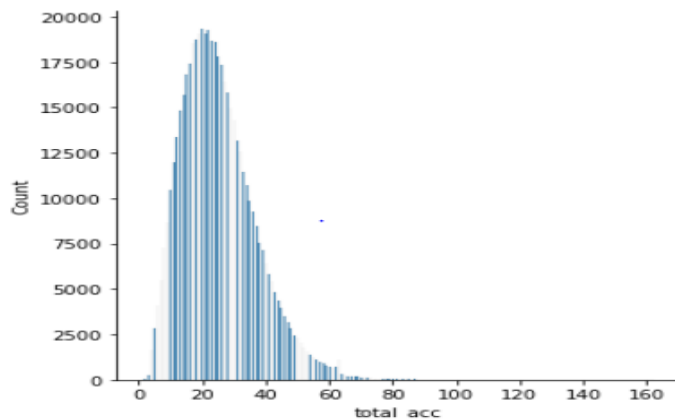
total_acc:

total number of credit lines available in members credit line



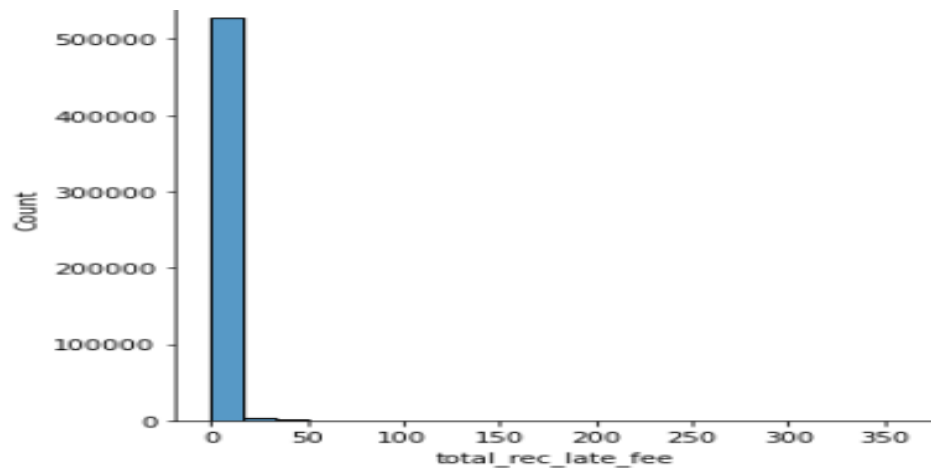
Since the data indicates the number of total credit lines available for in members credit line the range of credit lines available lie between 1 to 162 but the maximum number lies between 1-60 where the median is around 24, the data lies majorly between 0 to 50(IQR) the data also appears to be skewed to the right

total_rec_int: Interest received till date



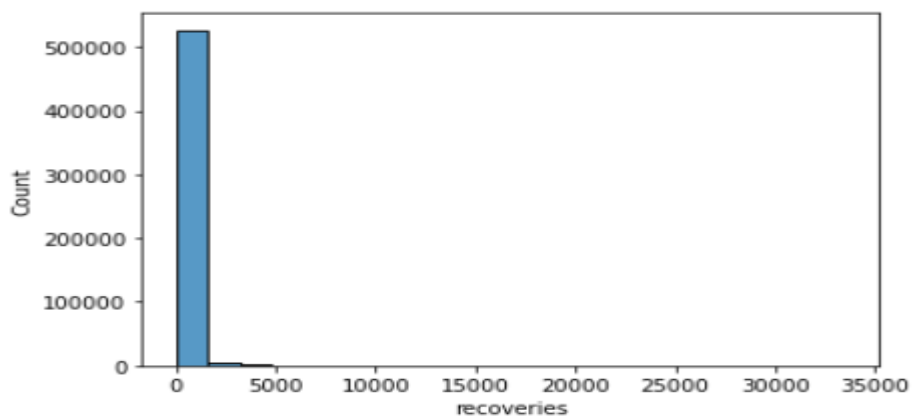
- The data also appears to be positively skewed,
- the interest received by the bank ranges between 0 to a maximum of 24205
- majorities of the data lies between 200-1900 after outlier removal
- the highest number of interest received lies between 0 to 2500

total_rec_late_fee: Late fee received till date



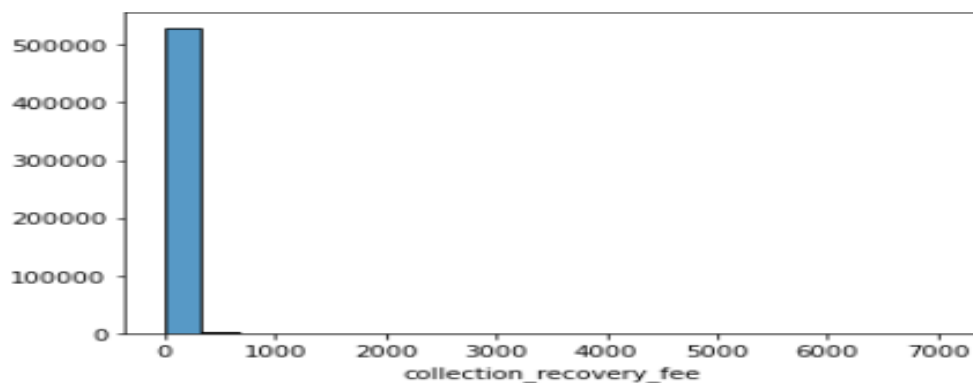
The range of late fee lies between 0 - 358, but the majority of the late fee is 0. One can conclude that most individuals pay on time. The outlier after zero to 350

recoveries:(post charge off gross recovery)



- The range of late fees lies between 0 - 35000, but the majority of the late fee is 0 ,most of the individuals pay on time and the recovery cost would not exist

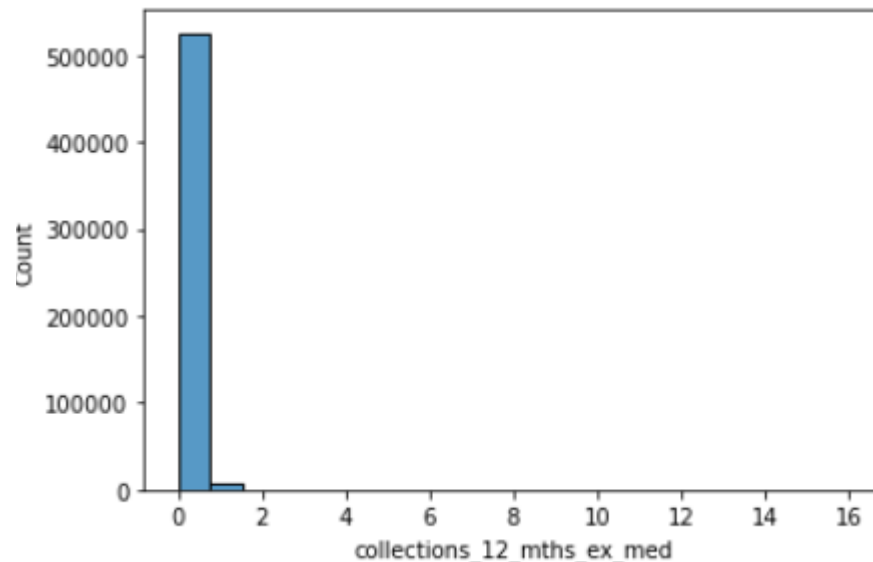
collection_recovery_fee: (post charge off collection fee)



The data has a lot of outliers but avg collection_recovery fee -0.4 to 0.4 after the removal of the outliers. There is a lot of multi correlation between collection recovery fee, recoveries, total_rec_late_fee

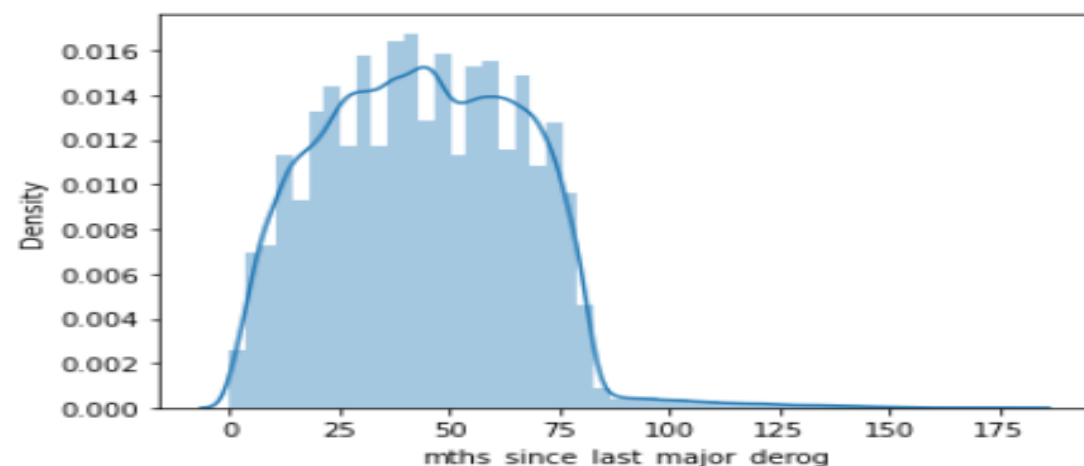
collections_12_mths_ex_med:

Number of collections in last 12 months excluding medical collections



- All the 4 data columns have right skew ,the collections_12_mths_ex_med multicollinear. Data majorly lies around 0

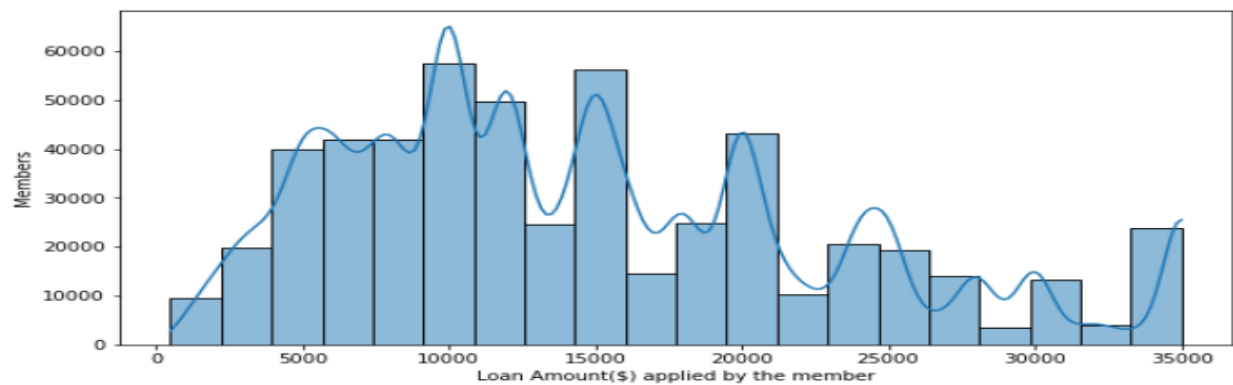
mths_since_last_major_derog: (months since the most recent 90-day or worse rating)



The iqr is between 25 -65, the mean appears to be 44, And has a slightly positive skew, outliers exist.

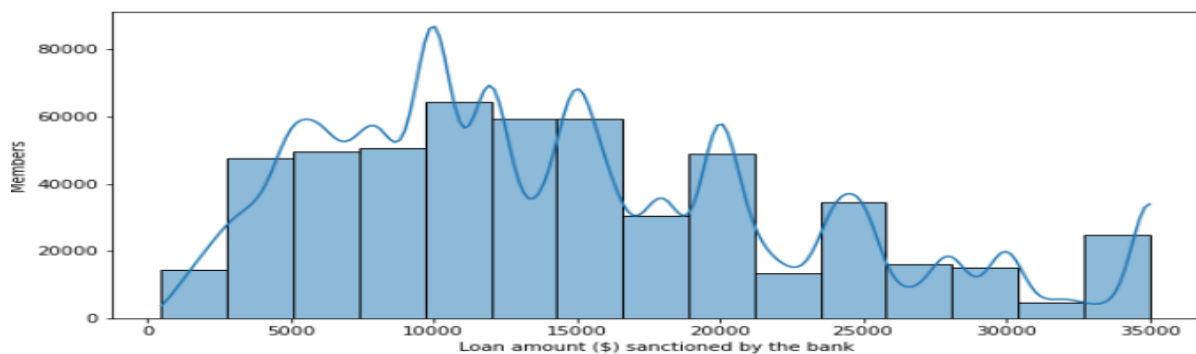
Loan Amount Attribute:

The maximum loan amount applied by the member is 35000\$ and approximately 58000 members applied the loan amount of 10000\$ and similarly 56000 members applied the loan amount of 15000\$.



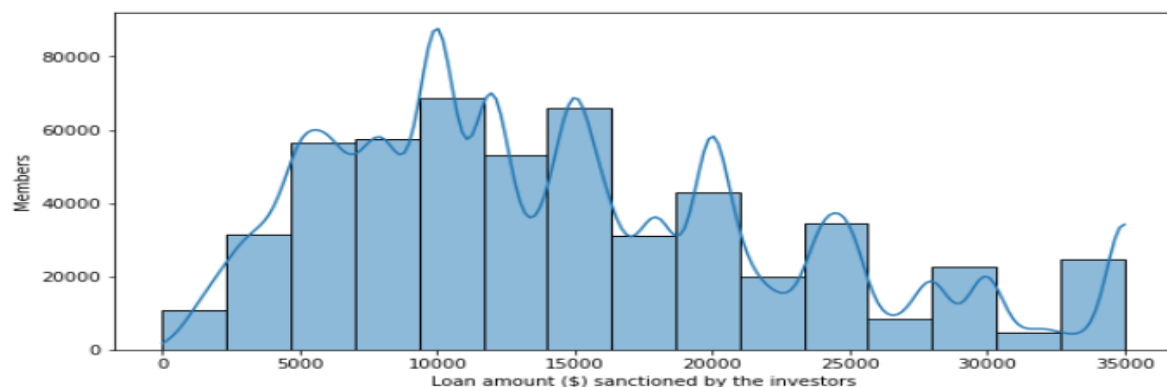
Funded Amount Bank Attribute:

The maximum loan amount sanctioned by the bank is 35000\$ and banks have also sanctioned the loan amount of 10000\$ for more than 60000 members.



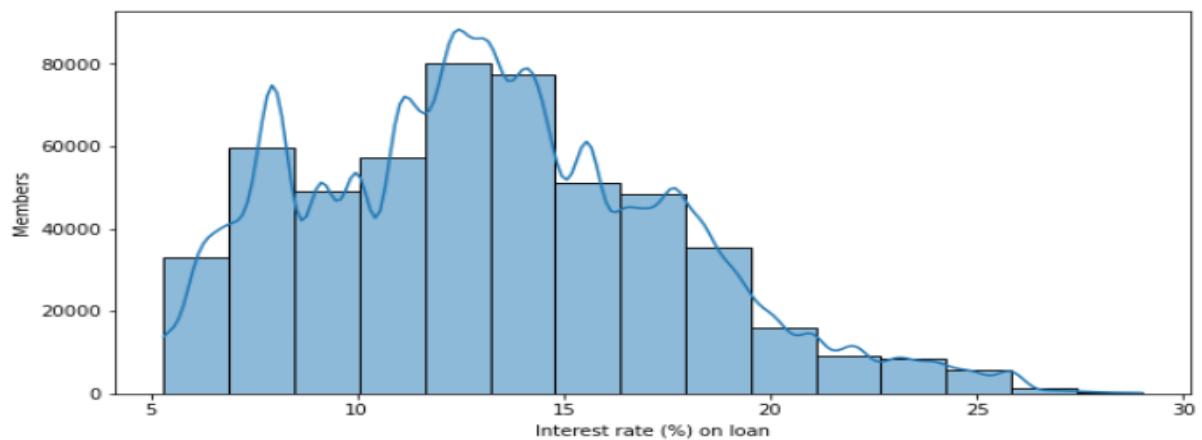
Funded Amount Investors Attribute:

From the above plot, we can observe that the maximum loan amount sanctioned by the investors is 35000\$ and investors have also sanctioned the loan amount of 10000\$ more than 68000 members

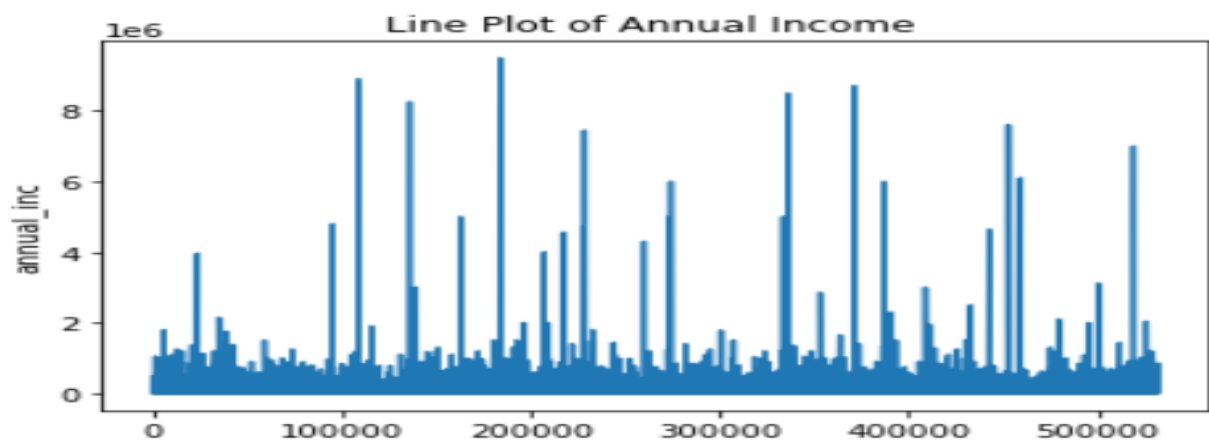


Interest Rate Attribute:

From the above plot, we can observe that the maximum interest rate given on a loan is around 26% and the minimum interest rate given on a loan is 6%. Approximately 80,000 members are paying the interest of around 13%.

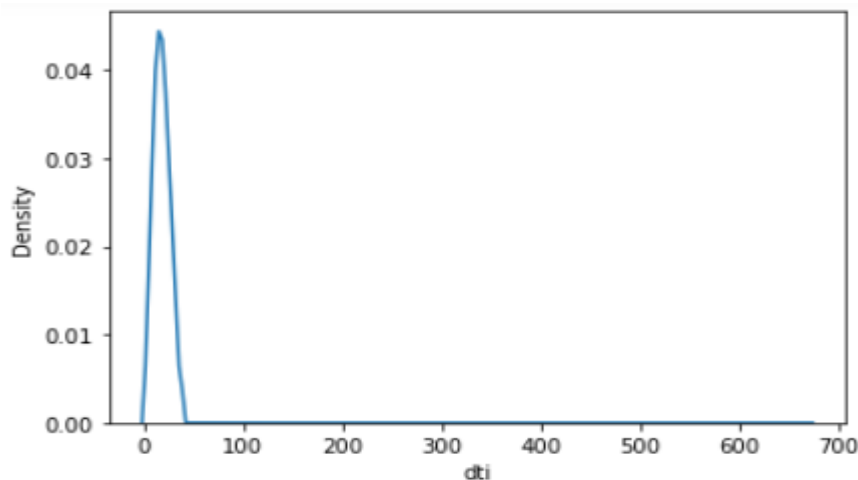


Annual Income Attribute:



Mean annual income reported by the members is approximately 75,000\$. Max annual income is 950,000\$. The range of the annual incomes reported by members lies between 45,000 \$ to 90,000\$.

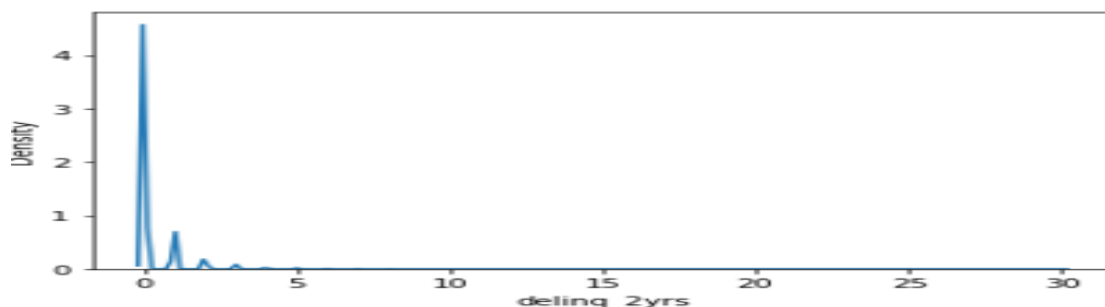
DTI :it is the ratio of member's total monthly debt repayment excluding mortgage divided by self reported monthly income



Lenders typically say the ideal front-end ratio should be no more than 28 per cent, and the back-end ratio, including all expenses, should be 36 per cent or lower. Lenders prefer borrowers with a lower DTI because that indicates less risk that you'll default on your loan. The majority of the data lies in between the range of 0 to 40.

delinq_2yrs: (number of 30+ days delinquency in past 2 years
delinquency:- minor crime)

If the payments are missed on payment in a single month across various credit accounts.

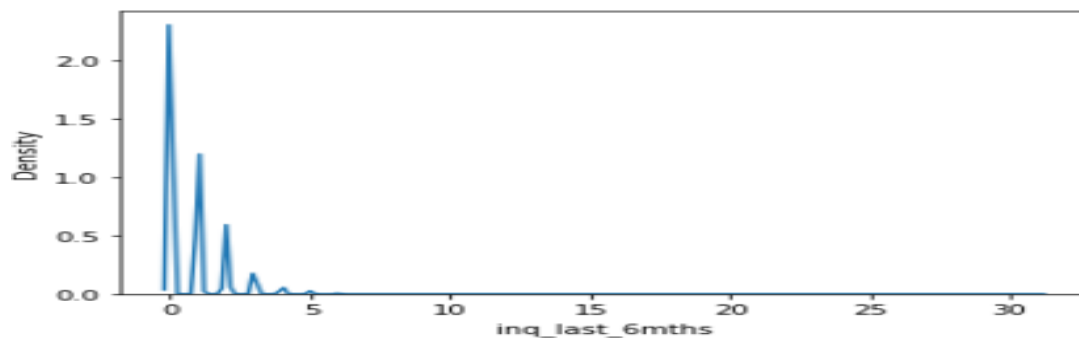


Most of the values lie at zero

inq_last_6mths: (number of inquiries in last 6 months)

- How many inquiries is too many in 6 months?
- For many lenders, six inquiries are too many to be approved for a loan or bank card. Even if you have multiple hard inquiries on your report in a short

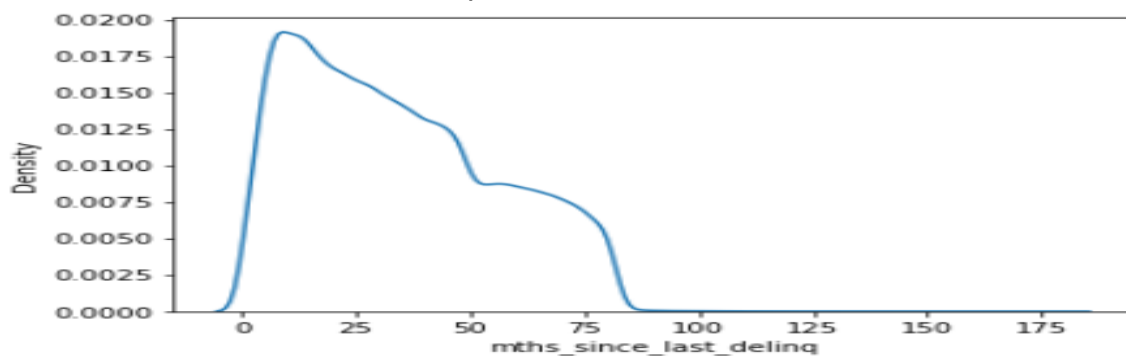
period of time, you may be spared negative consequences if you are shopping for a specific type of loan.



Most of the values lie at zero. There are outliers and there are a few null values, most of them are below 6 in our graph data. Increased the multiplying factor to 5 and reduced the outliers.

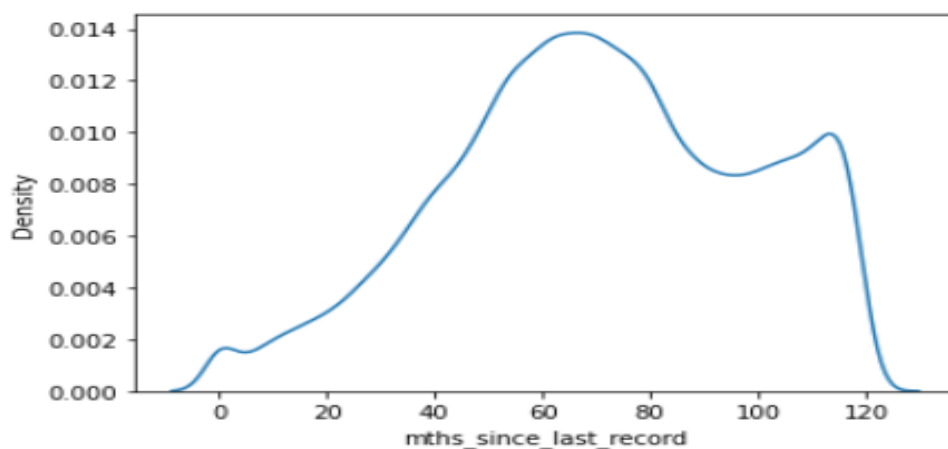
mths_since_last_delinq

Number of months since last delinq , few outliers and lots of null values



mths_since_last_record: (Number of months since last public record

- Public records may indicate you stopped paying your bills, it has two peaks

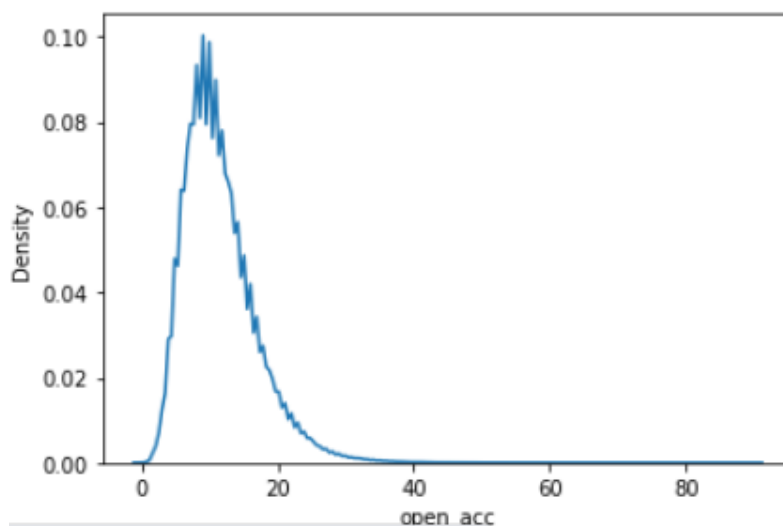


In the above graph, we could see that the record for defaulters has a high density towards the right side compared to the non-defaulters

open_acc:

- number of open credit lines in members' credit line
- Open credit is a pre-approved loan between a lender and a borrower.
- It allows the borrower to make repeated withdrawals up to a certain limit and then
- Make subsequent repayments before the payments become due

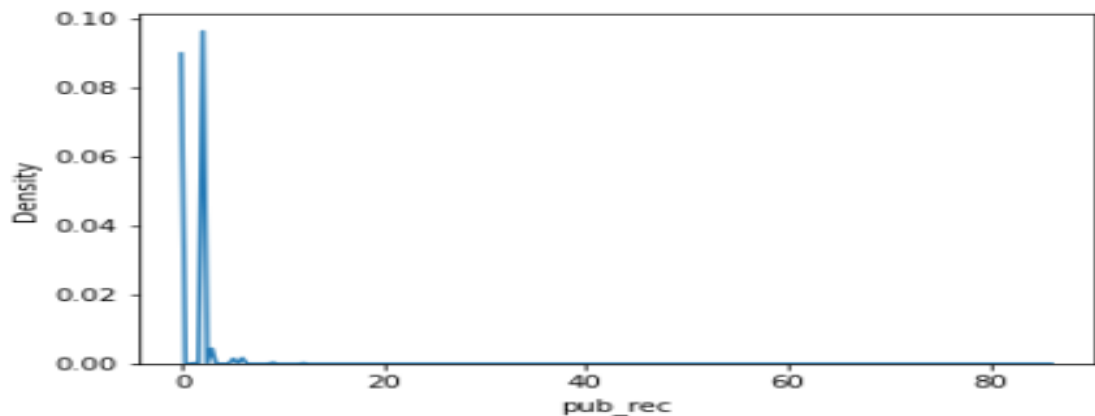
A new line of credit may improve your credit score. However, you should never take out an additional line of credit unless necessary. Applying for multiple lines of credit in a short period is not advised, and having too many lines of credit makes you look risky to lenders.



There is no evidence that the more the number of credit lines more is the chance of getting approved for a loan

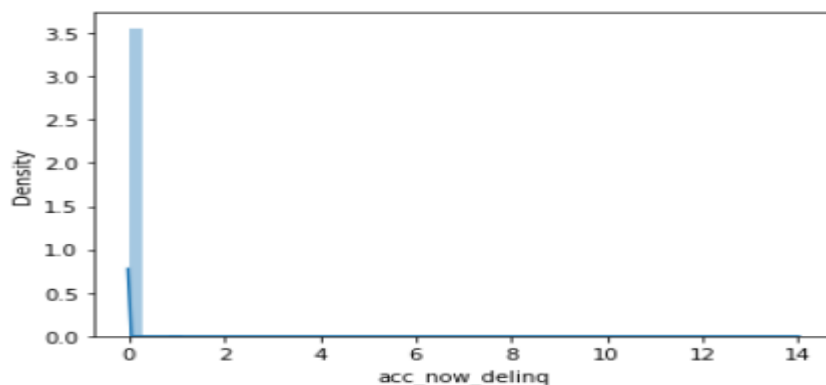
pub_rec: (number of derogatory public records)

A derogatory account is one that is seriously past due. Most commonly, the term derogatory refers to accounts that are 60 or 90 days past due or more. It also includes collection accounts, charge-offs, repossessions and foreclosures. The only type of public record information that would appear on your credit report is a bankruptcy filing.



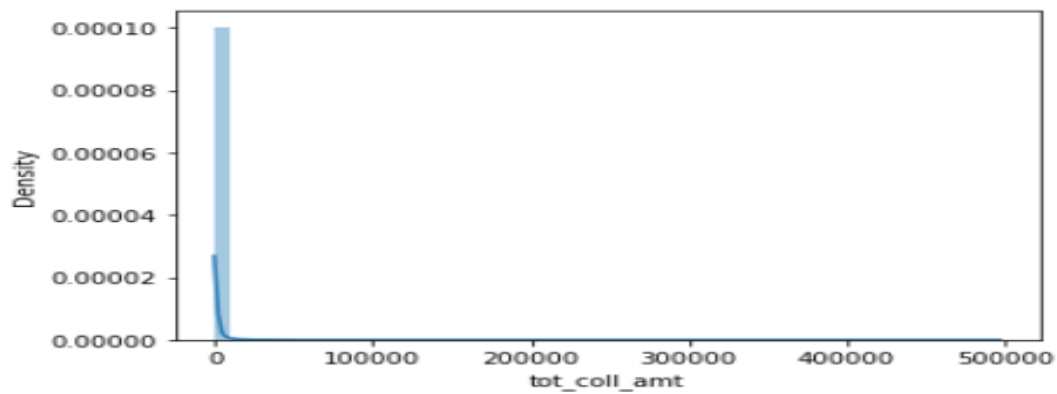
acc_now_delinq: number of accounts on which the member is delinquent

Most of the members pay on time and no default exists although there are a few outliers indicating that some people have a delay in their payment and their accounts have been marked



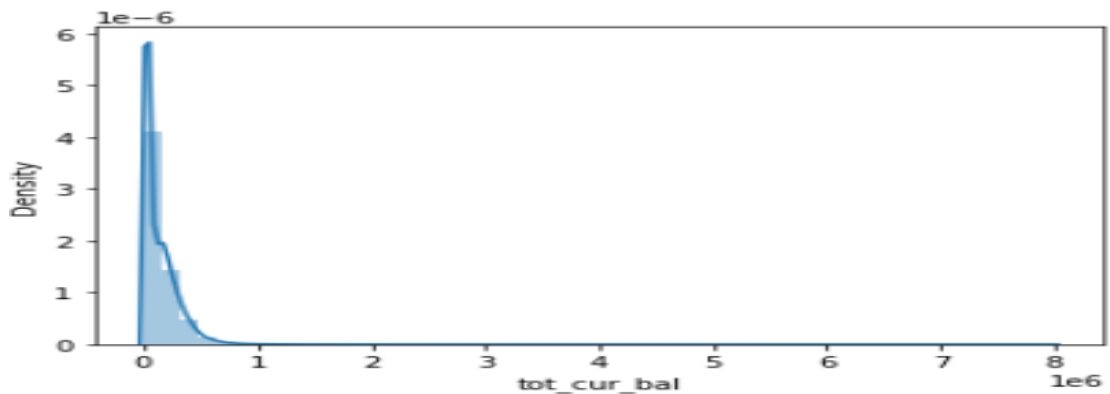
tot_coll_amt:

most of the data lies in zero which can be understood that people have yet to start to pay.



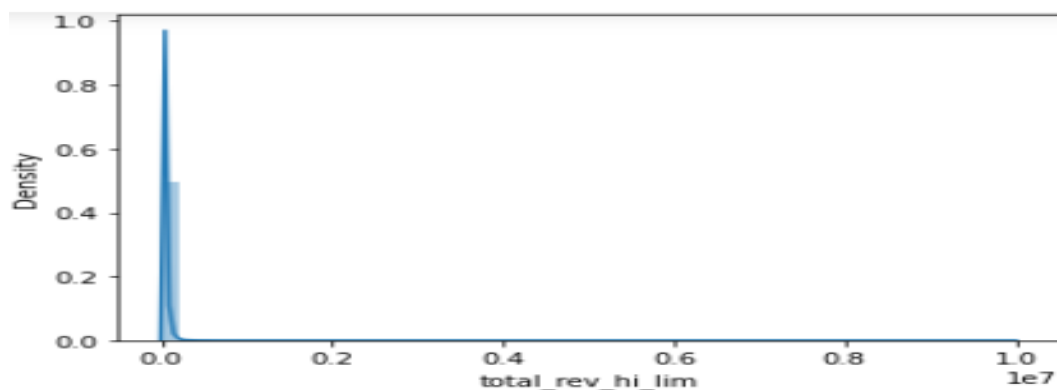
tot_cur_bal:

most of the data lie between 0 and 1, that people have yet to start to pay



total_rev_hi_lim:

most of the data lie between 0 and 0.2, that people have low credit revolving limit



Employee Title:

With more than 1lakh unique job/employer title of member, 'Teacher' title is the most number of titles given by an employer, almost 8000 members are with the Teacher employee title.

Desc:

This Attribute talks about the loan description given by the members. With more than 85% null values present in the said attribute. It is very insignificant . Hence can be dropped from the dataset.

Title:

This attribute is about the loan title given by the members. With almost 40000 unique records. 'Debt Consolidation' is the loan title provided by maximum members. Almost 2.48lakh members opted for Debt Consolidation as their loan title.

Batch Enrolled Attribute:

From the data, we can observe that Batch numbers allotted to members are mostly different.

Member ID Attribute:

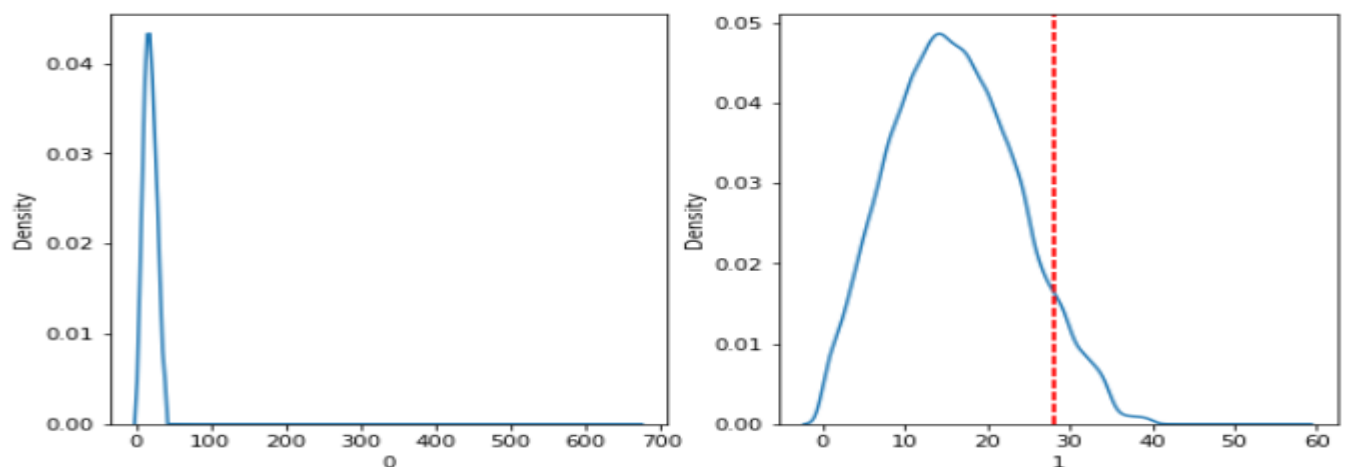
From the data, we can observe that the member id of each member is unique.

Outlier detection:

we have observed the presence of outliers in every numerical data type column which is mentioned during the univariate analysis

Bivariate data:

dti vs loan_status

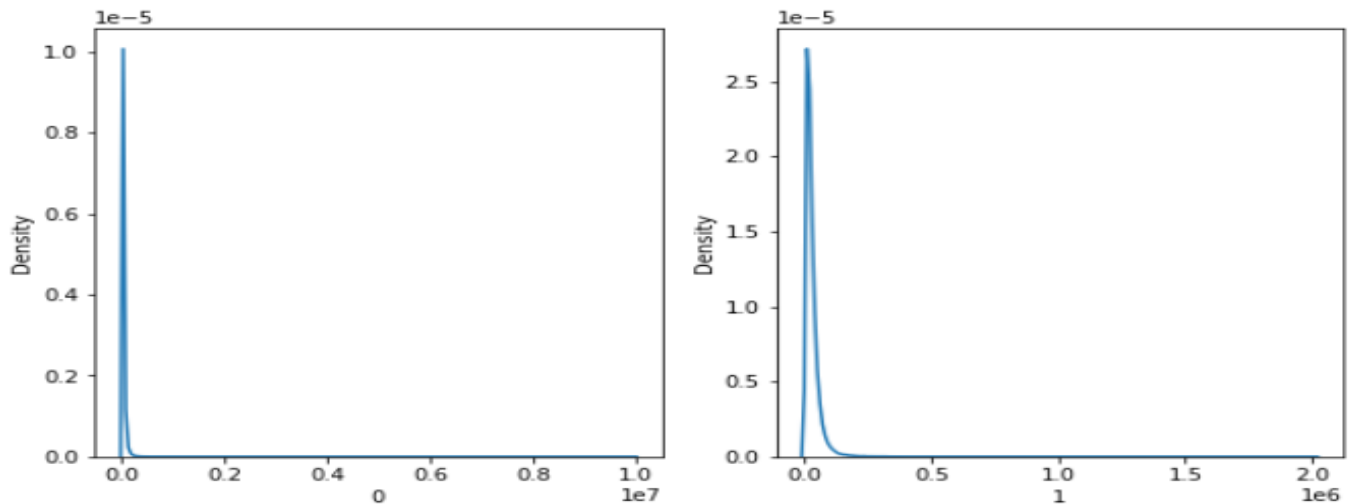


- Below graph shows two distributions of a non-defaulter and a defaulter

respectively

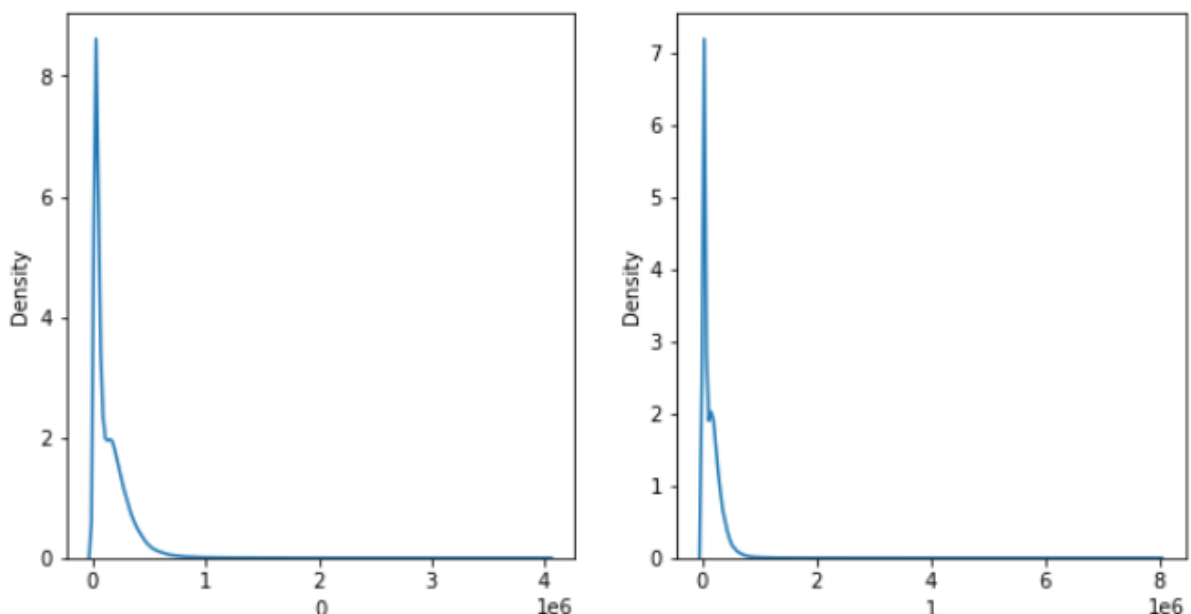
- In both cases we could see the majority of the data lies in between the range of 0 to 40
- In default graph we could see a lot of the member have ratios greater than 28 which could lead to rejection in loan

loan_status vs 'total_rev_hi_lim':



- Below graph shows two distribution of a non-defaulter and a defaulter respectively
- in non-default we see a most of the density of the data is situated in the zero value
- In default the area of density is more and has values greater than zero too.

loan_status vs tot_cur_bal

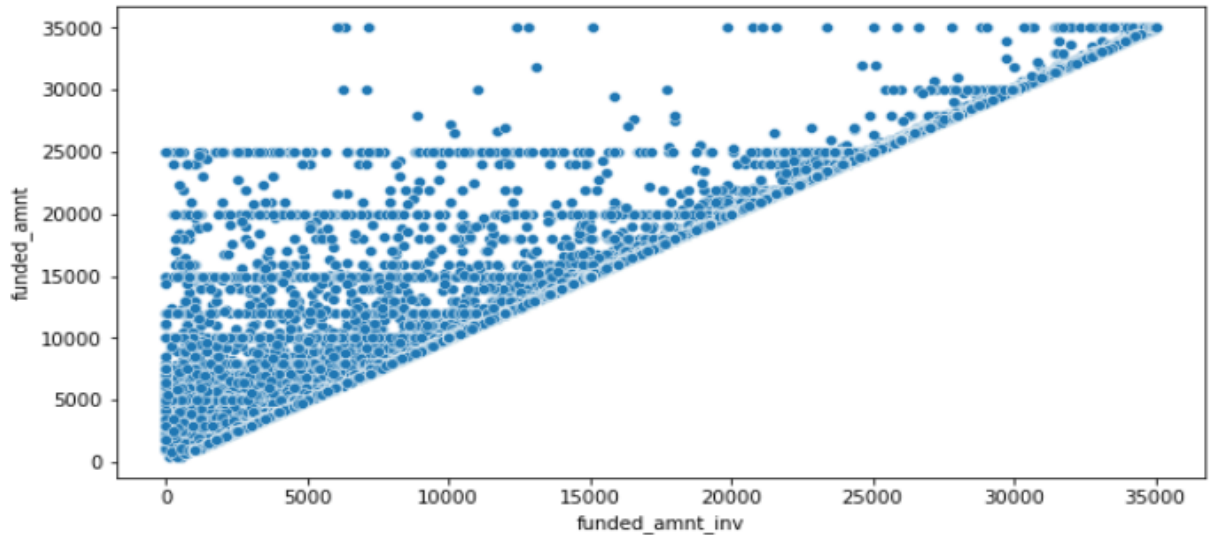


- Below graph shows two distributions of a non-defaulter and a defaulter

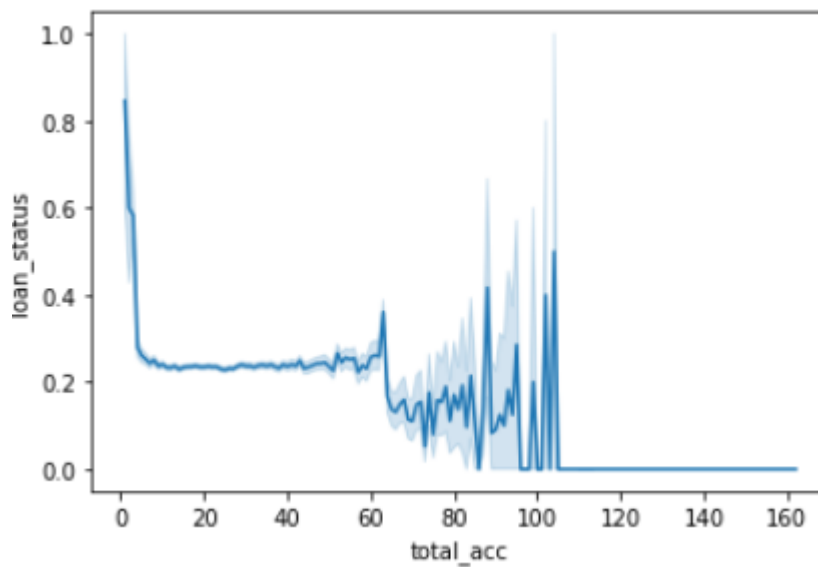
respectively

- In both the cases we can see the majority of the data lies in the non-default data
- In the default graph we can see a lot of the members have zero values

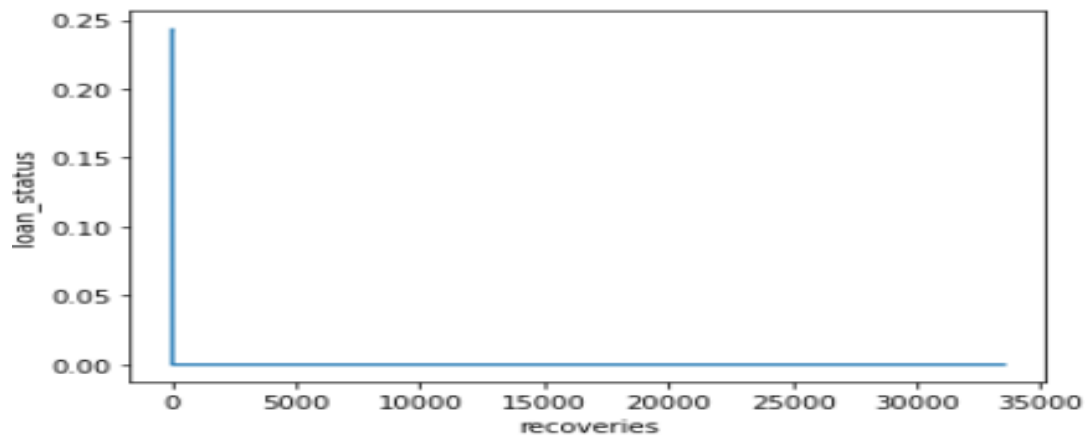
funded_amnt_inv vs funded_amnt: we can observe that Loan amount sanctioned by the bank is less or equal to the amount sanctioned by the investors.



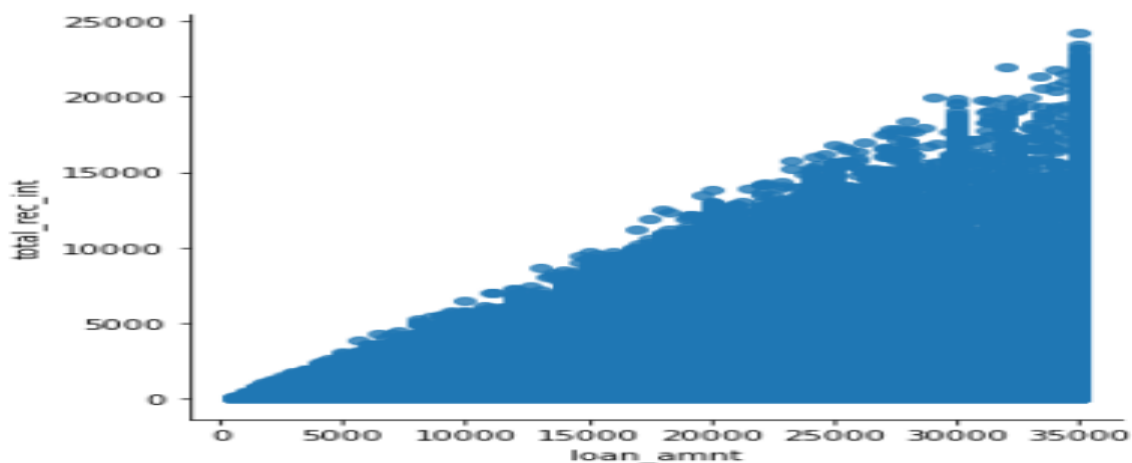
total_acc vs loan_status: there is a sudden drop which settles down and then we can observe and a smaller drop in loan status as total_acc increases



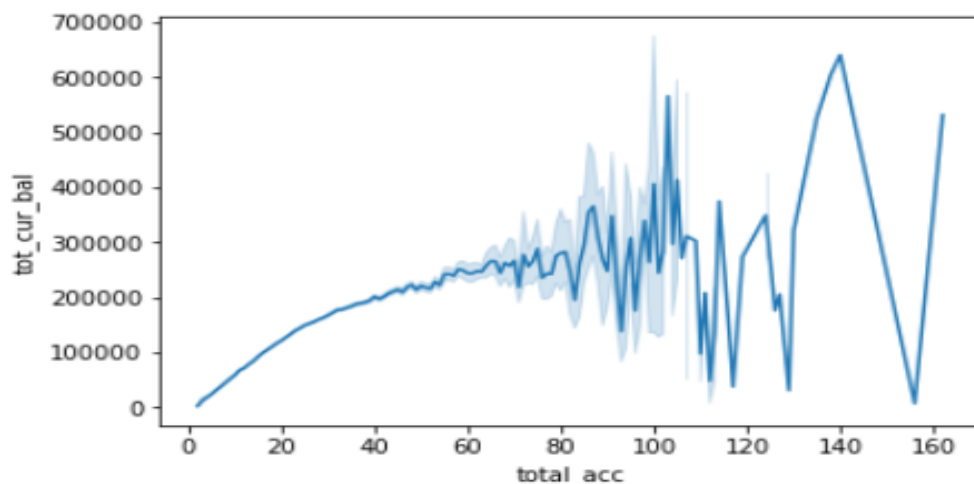
recoveries vs loan_status: L graph this indicates that most of the values lie on zero and there exists outliers.



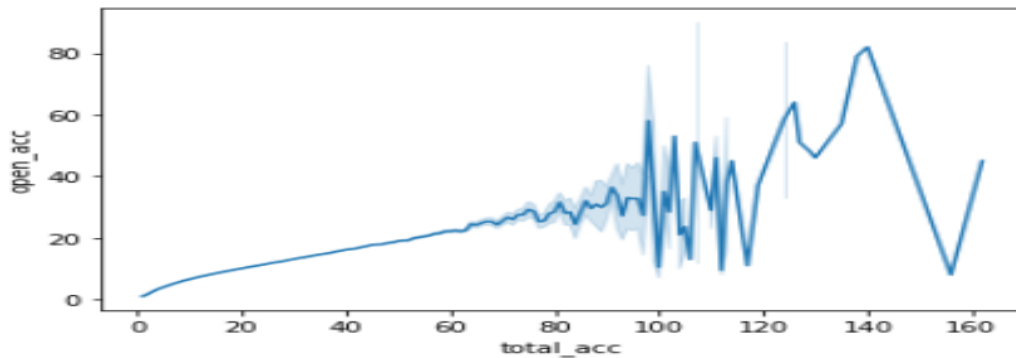
total_rec_int vs loan_amnt: Triangle-based filled graph indicating with total_rec_int increase there is a decrease in loan amount



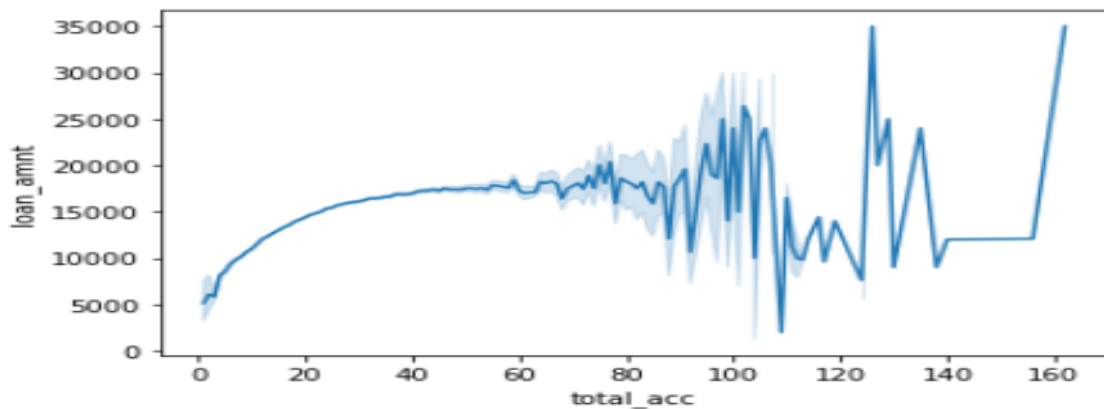
total_acc vs tot_cur_bal: With the increase in the total number of credit lines the loan status decreases



total_acc vs open_acc: The open_acc increases with an increase in total_acc and total_rec_int

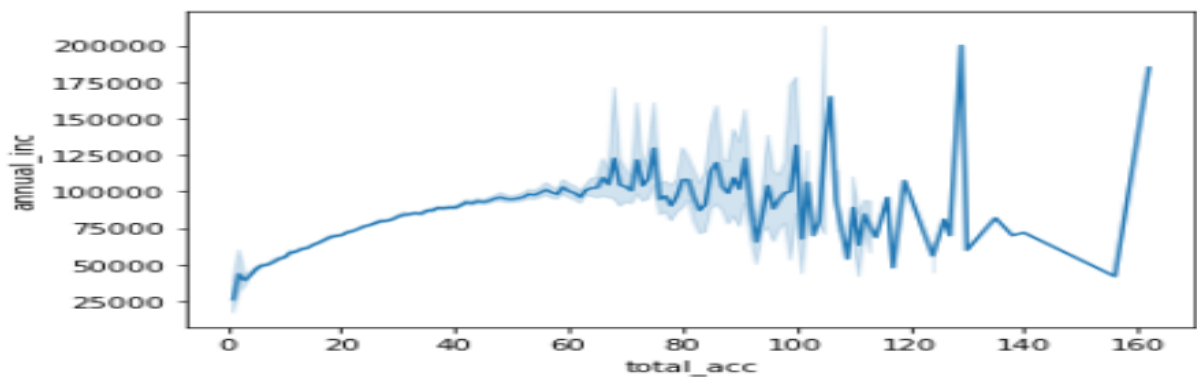


total_acc vs loan_amnt: Initially there is an exponential change and then it hits a plateau



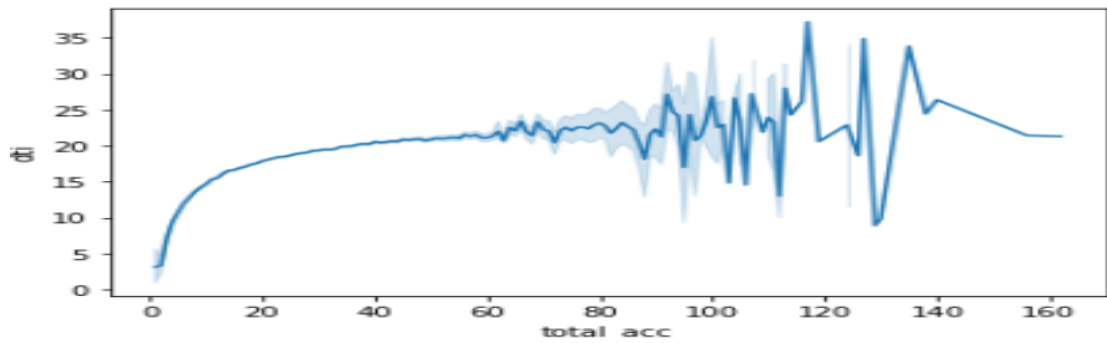
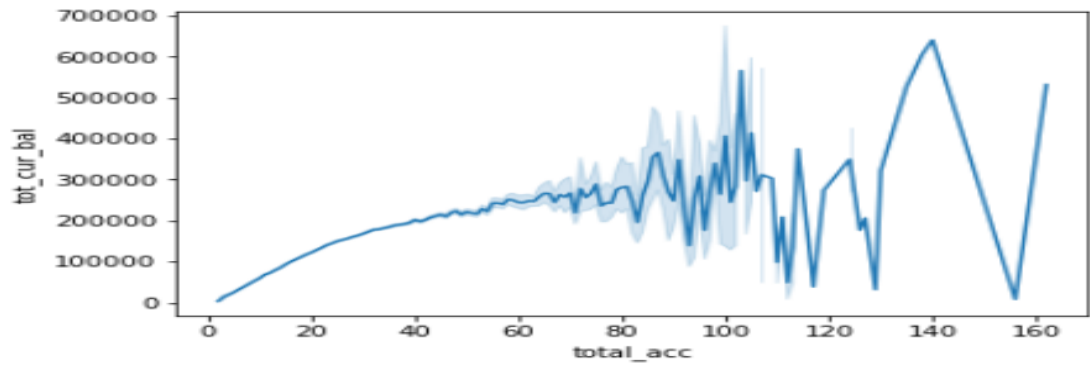
total_acc vs annual_inc:

annual_inc increases with increase in total_acc once total_acc reaches 80 it begins to drop



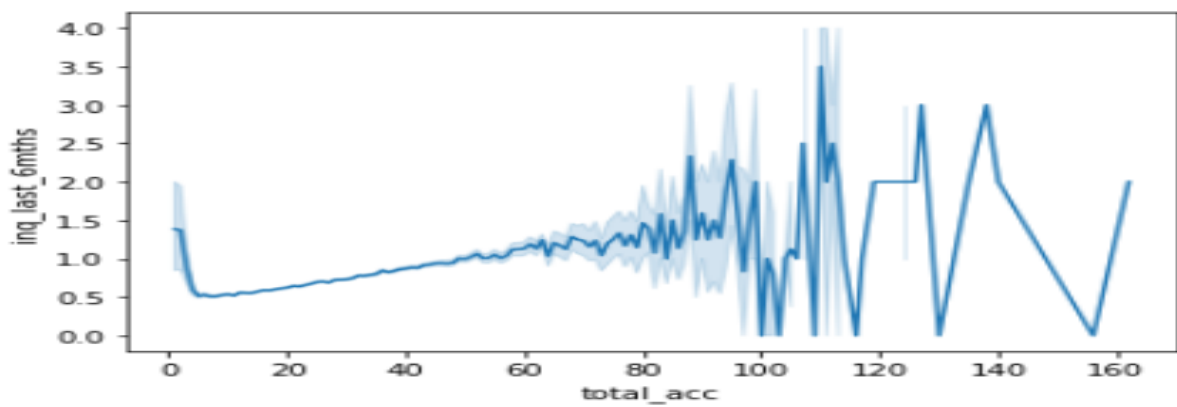
total_acc vs dti & ['tot_cur_bal']:

- initially there is a exponential change and then it hits a plateau at 40



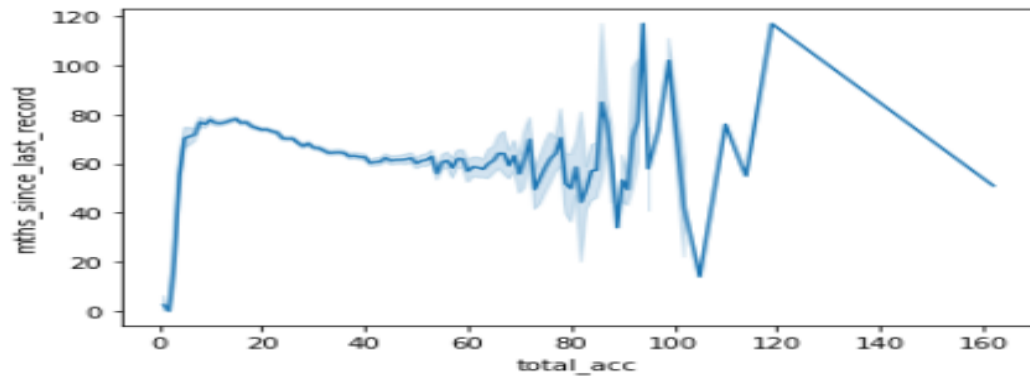
total_acc vs inq_last_6mths:

There was sudden drop and there was slow growth



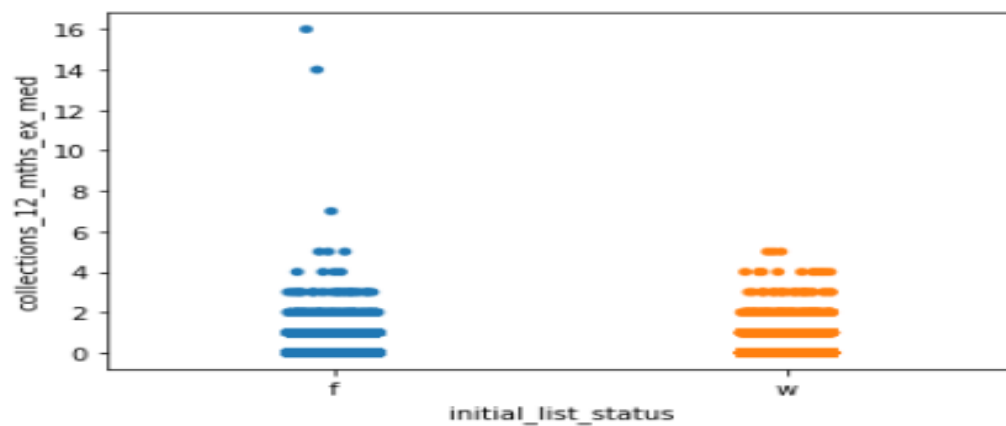
Total_acc vs mths_since_last_delinq and mths_since_last_record:

increases and then hits a plateau



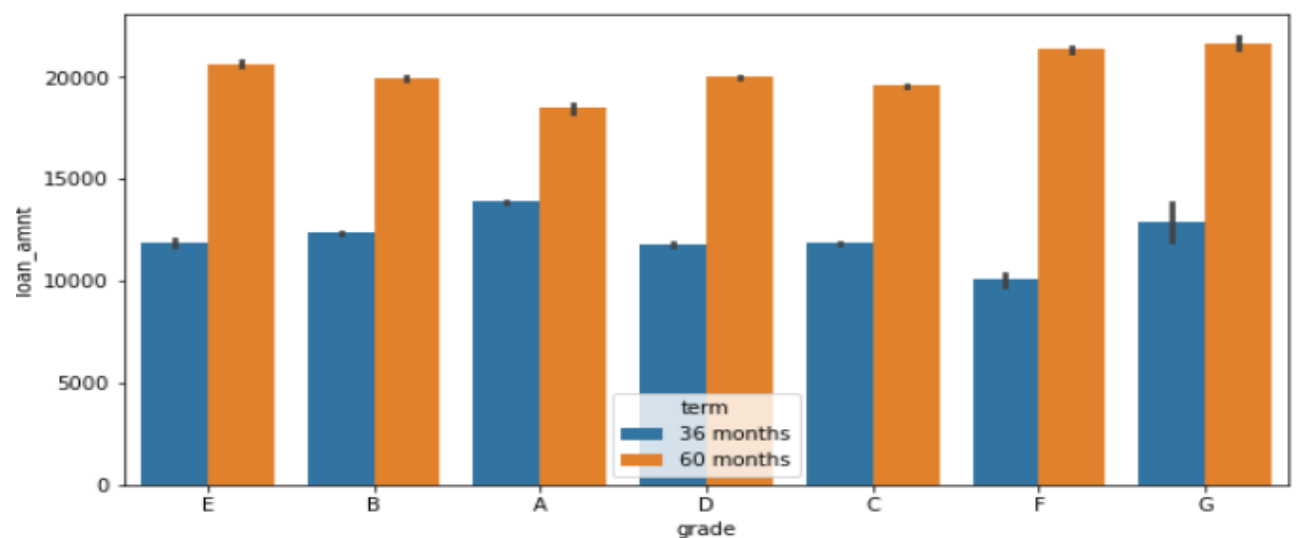
initial_list_status vs collections_12_mths_ex_med:

most of the values concentrate around 0-5



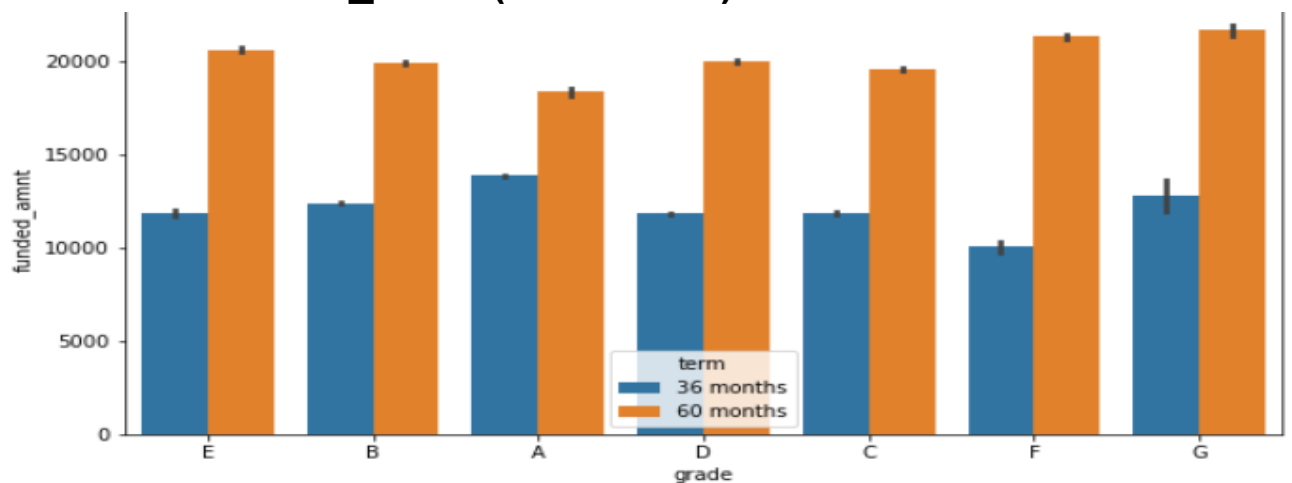
MULTIVARIATE

Grade vs Loan_amnt (Hue: Term)



- We can observe that the maximum loan amount applied to the bank is around 60 months and this is assigned under Grade G.

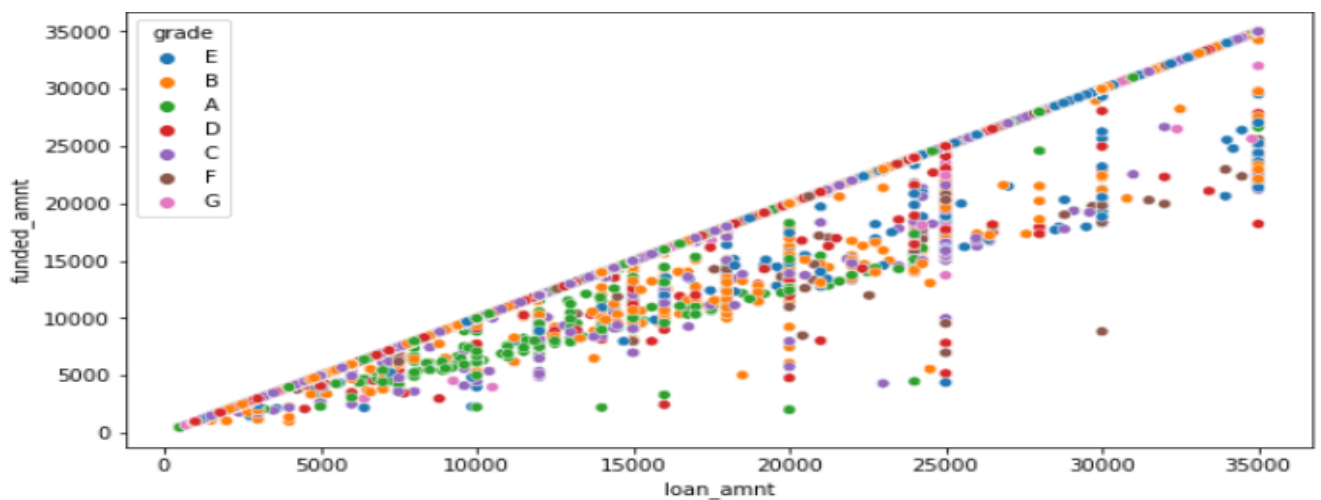
Grade vs funded_amnt (Hue: Term)



- we can observe that the maximum funded amount offered by the bank for the 60 months and it also comes under Grade G

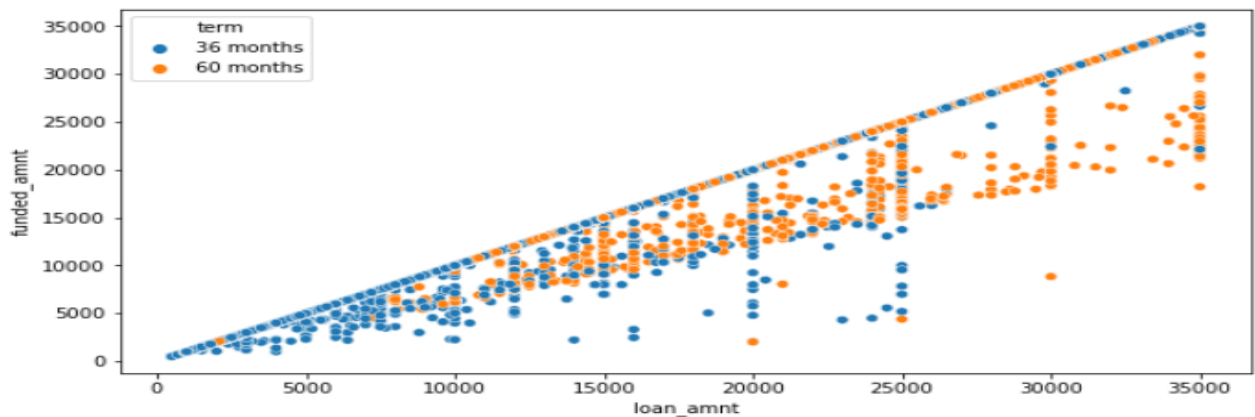
loan_amnt vs funded_amnt (Hue: grade)

- From the below plot, we can observe that as the loan amount is increasing the funded amount is also increasing.
- This is the directly proportional relationship between them.



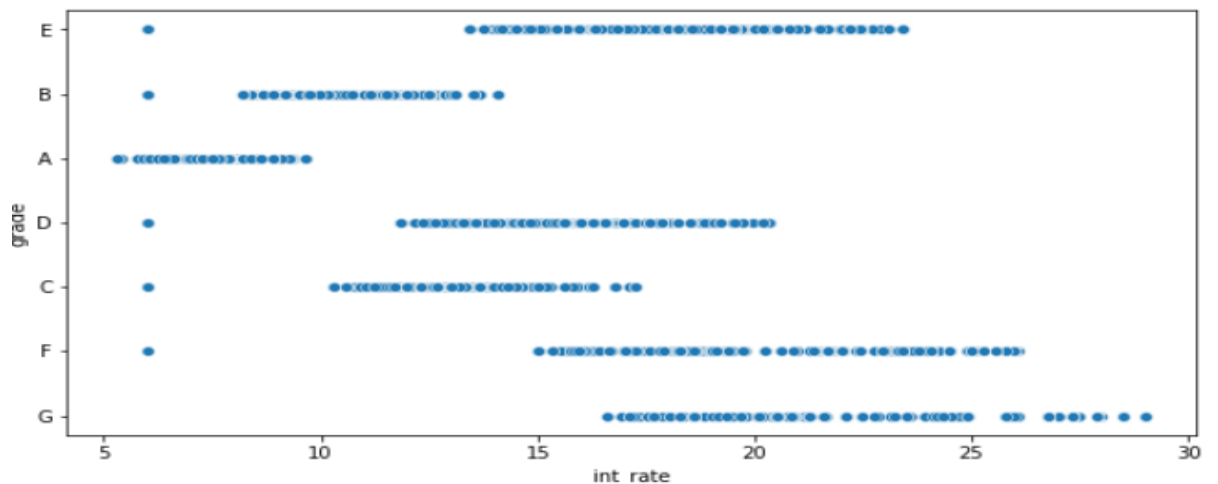
loan_amnt vs funded_amnt (Hue: term)

- From the below plot, we can observe that as the loan amount is increasing the funded amount is also increasing.
- This is the directly proportional relationship between them.
- when the loan_amount increases the number of default members also increases



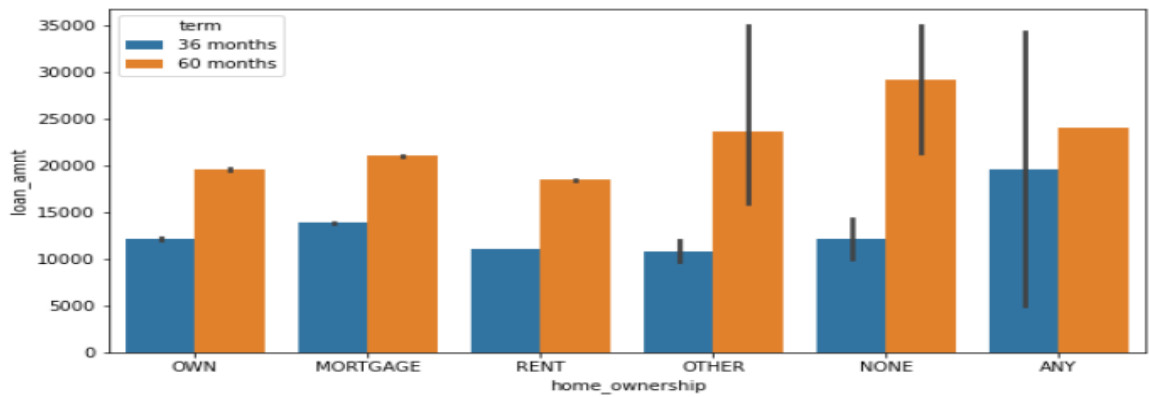
int_rate vs grade:

- From the below plot, we can observe that the maximum interest rate is applicable for the Grade G.



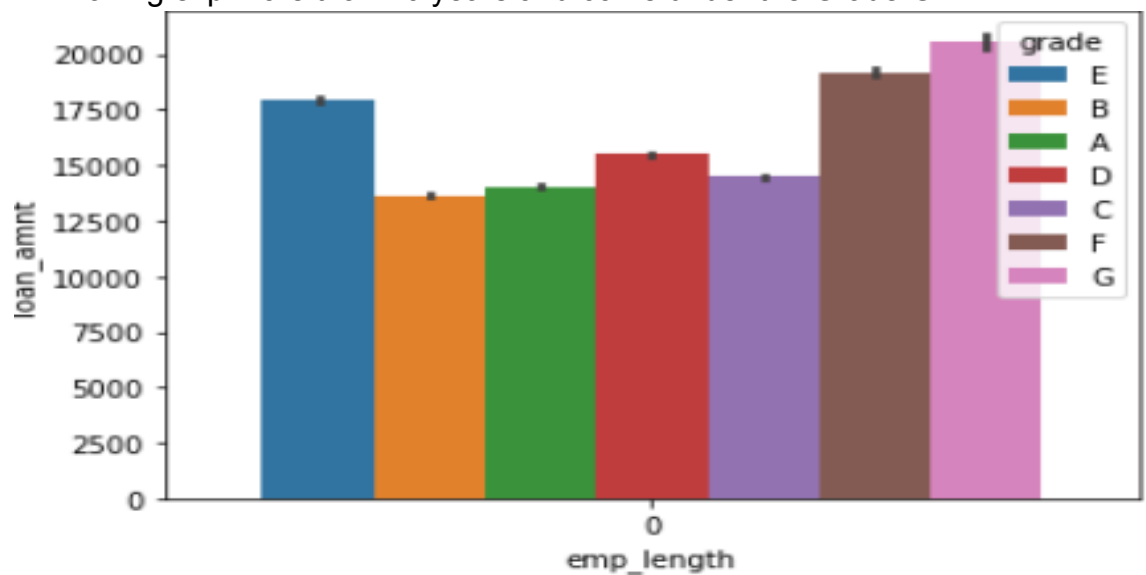
home_ownership vs loan_amnt (Hue: term)

- we can observe that maximum loan amount (30000\$) applied by the members whose ownership is categorised by 'None'.
- From the hist plot, the 'None' category home is very less.
- Maximum members (around 50%) come under the "Mortgage" category who applied the loan amount around 21000\$



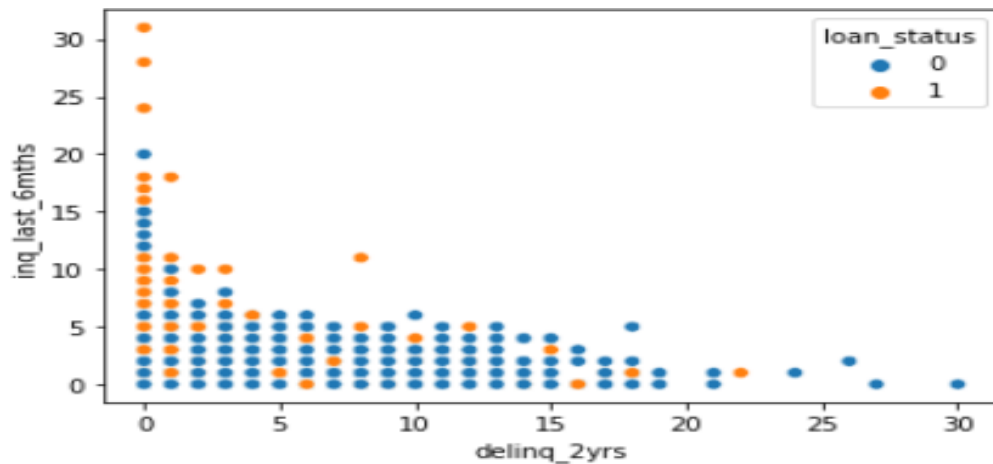
emp_length vs loan_amnt (Hue: grade)

- we can observe that maximum loan amount offered by the member having exp more than 10 years and come under the Grade G



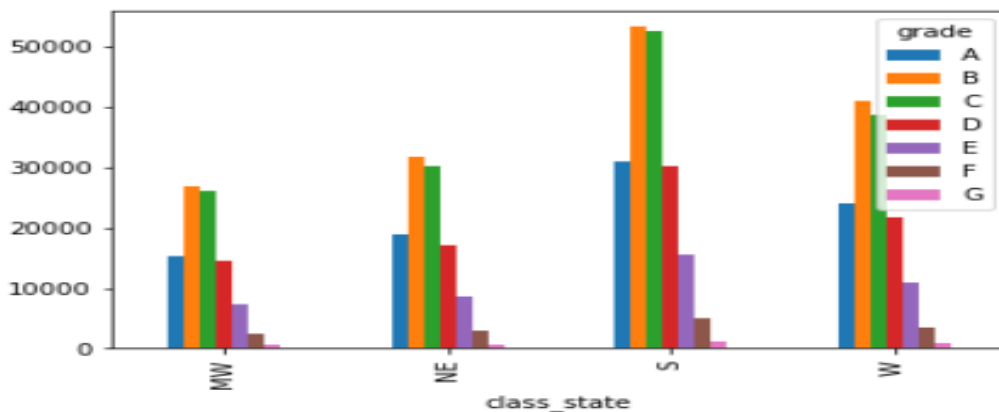
delinq_2yrs vs inq_last_6mths (hue: loan_status)

- From the plot we see that when the delinquency is zero the number of inquiries is more
- Although the delinquency is zero due to the low number in inquiries we see a lot of default loan status



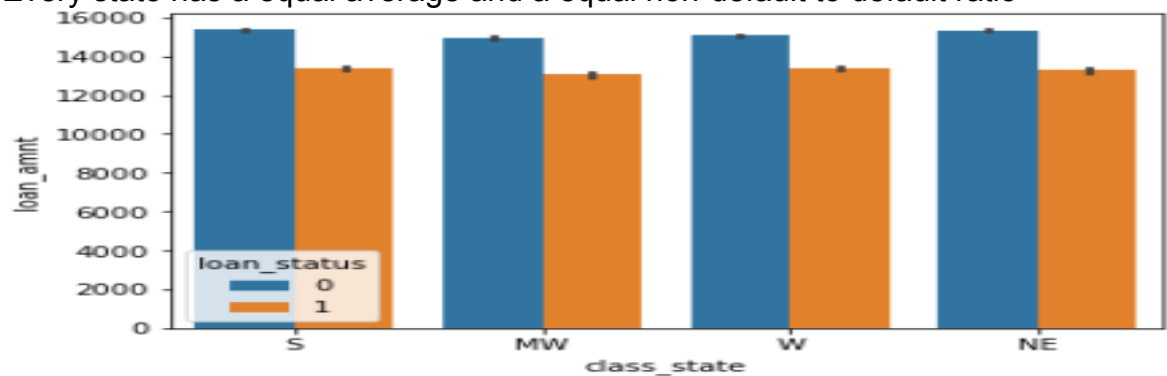
class_state vs grade

- The ratio of all the grades in every state is the same



class_state vs loan_amnt (hue: loan_status)

- Every state has a equal average and a equal non-default to default ratio



CORRELATION PLOT:

Inferences

- Member_id has a very weak negative correlation with loan_status.
- funded_amnt, funded_amnt_inv are perfectly positively correlated with loan_amnt.

- similarly funded_amnt,funded_amnt_inv,loan_amnt are perfectly positively correlated to each other.
- last_week_pay has a significant negative correlation with member_id.
- total_acc has a significant positive correlation with open_acc.
- total_rev_hi_lim has a significant positive correlation with revol_bal.
- collection_recovery_fee has a significant positive correlation with recoveries.
- total_rec_int has a weak positive correlation with loan_amnt,funded_amnt,funded_amnt_inv,last_week_pay columns.
- mths_since_last_dealing has a weak negative correlation with delinq_2yrs.
- member_id has a weak negative correlation with total_rec_int columns.
- mths_since_last_major_derog has weak positive relation with mths_since_last_delinq.
- tot_cur_bal has a fragile positive relation with annual_inc.

member_id	1	0.066	0.680	0.73	0.07	-0.140	0.120	0.35	0.13	0.05	-0.130	0.012	0.590	0.088	0.08	0.0360	0.5	0.031	-0.460	0.760	1.10	0.720	0.510	0.27	-0.790	0.180	0.370	0.510	0.667	-0.49
loan_amnt	-0.066	1	1	1	0.41	0.15	0.1	0.33	0.042	0.0048	0.340	0.28	0.0450	0.2	-0.081	0.33	0.12	0.22	0.53	0.0310	0.730	0.530	0.150	0.180	0.048	0.032	0.160	0.32	0.3	-0.096
funded_amnt	-0.068	1	1	1	0.41	0.15	0.1	0.33	0.048	0.0010	0.340	0.28	0.0450	0.2	-0.081	0.33	0.12	0.22	0.53	0.0310	0.730	0.530	0.150	0.180	0.048	0.032	0.160	0.32	0.3	-0.098
funded_amnt_inv	-0.073	1	1	1	0.41	0.15	0.1	0.33	0.048	0.0028	0.340	0.26	0.022	0.2	-0.08	0.33	0.12	0.22	0.53	0.0280	0.71	0.05	-0.0140	0.180	0.058	0.030	0.160	0.33	0.3	-0.1
term	-0.07	0.41	0.41	0.41	1	0.43	0.0670	0.59	0.10	0.00350	0.070	0.048	0.140	0.0850	0.020	0.920	0.86	0.1	0.380	0.053	0.580	0.340	0.038	0.035	0.038	0.050	0.078	0.11	0.065	-0.13
int_rate	-0.14	0.15	0.15	0.15	0.43	1	0.0670	0.720	0.16	0.057	0.230	0.0140	0.010	0.099	0.520	0.360	0.27	0.037	0.45	0.056	0.110	0.0710	0.0140	0.066	0.760	0.0270	0.040	0.0830	0.140	0.026
emp_length	-0.012	0.1	0.1	0.1	0.067	0.067	1	0.0630	0.0410	0.0230	0.018	0.0024	0.0210	0.0370	0.0310	0.0850	0.340	0.990	0.680	0.018	0.0048	0.038	0.00350	0.008	0.0840	0.0070	0.090	0.730	0.028	0.015
annual_inc	-0.035	0.33	0.33	0.33	0.0590	0.720	0.63	1	0.180	0.0480	0.330	0.330	0.340	0.130	0.007	0.03	0.037	0.18	0.13	0.011	0.0070	0.066	0.040	0.180	0.028	0.015	0.001	10.41	0.260	0.08
dti	-0.13	0.0420	0.0430	0.045	0.1	0.160	0.041	-0.18	1	0.0078	0.018	0.012	0.026	0.3	-0.0450	0.14	0.18	0.22	0.0180	0.18	0.0028	0.0039	0.011	0.015	-0.10	0.0078	0.012	0.0038	0.073	-0.13
delinq_2yrs	-0.050	0.0088	0.0190	0.0088	0.0035	0.570	0.230	0.440	0.07	1	0.022	0.480	0.370	0.530	0.130	0.320	0.180	0.120	0.028	0.18	0.0081	-0.09	0.63	-0.320	0.430	0.130	0.0340	0.07	0.330	0.046
inq_last_6mths	-0.130	0.340	0.340	0.340	0.037	0.230	0.0120	0.330	0.018	0.22	1	0.0130	0.390	0.11	0.0560	0.170	0.0870	0.140	0.0880	0.340	0.450	0.330	0.007	0.01	0.066	0.0320	0.01	0.020	0.000	0.087
mths_since_last_delinq	-0.0010	0.280	0.280	0.280	0.048	0.014	0.0024	0.330	0.00170	0.480	0.015	1	0.0160	0.0270	0.0750	0.026	0.0072	0.140	0.220	0.015	0.008	0.0028	0.02	0.47	0.0030	0.130	0.0330	0.0620	0.026	0.015
mths_since_last_record	-0.058	0.045	0.045	0.0022	0.014	0.01	0.0210	0.340	0.0260	0.370	0.030	0.1	1	0.00210	0.140	0.088	0.330	0.720	0.025	0.090	0.380	0.0340	0.010	0.170	0.046	0.078	0.320	0.008	0.0078	0.027
open_acc	-0.088	0.2	0.2	0.2	0.088	0.099	0.37	0.13	0.3	0.053	0.110	0.028	0.002	1	0.025	0.22	-0.14	0.7	0.062	0.099	0.00	0.0028	0.110	0.0040	0.770	0.018	0.00640	0.24	0.310	0.065
pub_rec	-0.08	-0.080	0.810	0.80	0.260	0.520	0.330	0.078	0.450	0.130	0.560	0.75	0.140	0.025	1	-0.1	-0.0790	0.130	0.060	0.110	0.0140	0.0078	0.0210	0.0690	0.780	0.0013	0.380	0.060	0.940	0.048
revol_bal	-0.036	0.33	0.33	0.33	0.0920	0.330	0.085	0.3	0.14	-0.0320	0.170	0.024	0.008	0.22	-0.1	1	0.22	0.19	0.140	0.00240	0.018	0.00820	0.028	0.008	0.018	0.00088	0.270	0.43	0.79	-0.04
revol_util	-0.0530	0.12	0.12	0.12	0.086	0.270	0.0340	0.037	0.18	-0.0180	0.088	0.0070	0.033	-0.140	0.790	0.22	1	-0.11	0.18	0.0210	0.0280	0.0190	0.036	0.00590	0.60	0.280	0.040	0.81	-0.110	0.047
total_acc	-0.031	0.22	0.22	0.22	0.1	-0.0370	0.099	0.18	0.22	0.12	0.140	0.0140	0.77	0.7	0.013	0.19	-0.11	1	0.0920	0.0380	0.0920	0.010	0.0098	0.220	0.048	0.260	0.29	0.3	0.250	0.002
total_rec_int	-0.46	0.53	0.53	0.53	0.38	0.450	0.068	0.13	0.018	0.0028	0.880	0.220	0.250	0.062	-0.06	0.14	0.18	0.092	1	0.098	0.680	0.530	0.0240	0.260	0.53	0.0012	0.210	0.11	0.066	0.038
total_rec_late_fee	-0.0780	0.310	0.310	0.310	0.005	0.560	0.110	0.130	0.110	0.180	0.340	0.015	0.090	0.099	0.12	0.028	0.028	0.003	0.09	1	0.0720	0.680	0.038	0.0038	0.740	0.00240	0.047	0.038	0.000	0.004
recoveries	-0.110	0.730	0.730	0.710	0.58	0.118	0.048	0.072	0.00250	0.00470	0.450	0.008	0.00380	0.089	0.01	-0.110	0.028	0.0092	0.680	0.72	1	0.8	0.0048	0.00540	0.018	0.018	0.00500	0.000	0.018	0.062
collection_recovery_fee	-0.0720	0.530	0.530	0.05	0.0360	0.710	0.0038	0.0068	0.00371e	0.5	0.03	0.00380	0.00028	0.00380	0.08	0.019	0.01	0.0520	0.65	0.8	1	0.0028	0.00380	0.0070	0.008	0.00380	0.000	0.00380	0.0018	0.043
collections_12_mths_ex_med	-0.0510	0.150	0.150	0.14	0.036	0.016	0.0015	0.04	0.0110	0.0630	0.0070	0.240	0.010	0.110	0.0210	0.230	0.036	0.0098	0.24	0.038	0.048	0.002	1	0.0790	0.430	0.520	0.50	0.0078	0.150	0.034
mths_since_last_major_derog	-0.0270	0.180	0.180	0.180	0.0038	0.0068	0.0028	0.180	0.015	-0.32	0.01	0.47	0.0170	0.0040	0.068	0.0083	0.0059	0.220	0.260	0.038	0.0058	0.033	0.07	1	0.0210	0.780	0.250	0.360	0.060	0.005
last_week_pay	-0.780	0.450	0.480	0.530	0.340	0.760	0.0088	0.28	-0.1	-0.0430	0.680	0.0038	0.440	0.770	0.780	0.190	0.06	-0.049	0.53	0.0740	0.018	0.0078	0.430	0.2	1	0.0170	0.330	-0.050	0.520	0.26
acc_now_delinq	-0.018	0.030	0.038	0.038	0.0058	0.028	0.0084	0.018	0.00730	0.130	0.00370	0.130	0.0078	0.018	0.000	0.0088	0.028	0.0240	0.00120	0.0028	0.0018	0.00038	0.520	0.780	0.01	1	0.0012	0.260	0.088	0.014
tot_coll_amnt	-0.0370	0.160	0.160	0.160	0.0078	0.0040	0.0070	0.00110	0.018	0.00340	0.010	0.0330	0.038	0.0064	0.380	0.270	0.040	0.0290	0.28	0.0047	0.0058	0.038	0.540	0.250	0.038	0.001	1	0.0030	0.220	0.023
tot_cur_bal	-0.051	0.32	0.32	0.33	0.110	0.088	0.099	0.48	0.00340	0.070	0.0220	0.0620	0.0990	0.24	0.066	0.43	0.081	0.3	0.110	0.00380	0.0098	0.018	0.0078	0.360	0.50	0.260	0.003	1	0.380	0.028
total_rev_hi_lim	-0.067	0.3	0.3	0.3	0.065	0.160	0.73	0.26	0.0730	0.0380	0.0078	0.26	0.0078	0.31	0.094	0.79	0.11	0.25	0.0660	0.008	0.0070	0.0018	0.019	0.066	0.58	0.008	0.220	0.38	1	0.041
loan_status	-0.490	0.990	0.98	0.1	-0.130	0.0020	0.028	0.0080	0.130	0.048	0.0870	0.150	0.270	0.050	0.490	0.040	0.48	0.0026	0.038	0.044	0.620	0.430	0.340	0.0050	0.260	0.0140	0.230	0.0280	0.4	-1

- before scaling the data it was observed that the data had a positive skew
- Also in the univariate analysis we noticed the presence of outliers and have chosen not to remove them since they can help in the identification of the possible defaulters

Model building:

Assumptions of logistic regression:

1. The logistic regression assumes that there is minimal or no multicollinearity among the independent variables.
2. The Logistic regression assumes that the independent variables are linearly related to the log of odds.
3. Logistic regression usually requires a large sample size to predict properly.

the model score is currently: 0.780

Logistic Regression Summary:

Logit Regression Results						
Dep. Variable:	loan_status	No. Observations:	425929			
Model:	Logit	Df Residuals:	425896			
Method:	MLE	Df Model:	32			
Date:	Fri, 22 Jul 2022	Pseudo R-squ.:	0.1820			
Time:	13:01:14	Log-Likelihood:	-1.9051e+05			
converged:	False	LL-Null:	-2.3290e+05			
Covariance Type:	nonrobust	LLR p-value:	0.000			
	coef	std err	z	P> z	[0.025	0.975]
const	-16.6707	0.147	-113.609	0.000	-16.958	-16.383
loan_amnt	0.0003	1.76e-05	19.453	0.000	0.000	0.000
funded_amnt	0.0003	2.58e-05	10.043	0.000	0.000	0.000
funded_amnt_inv	-0.0006	1.82e-05	-32.410	0.000	-0.001	-0.001
term	-0.4677	0.012	-39.576	0.000	-0.491	-0.444
int_rate	0.5464	0.005	117.194	0.000	0.537	0.556
emp_length	0.0008	0.001	0.729	0.466	-0.001	0.003
annual_inc	-1.065e-07	7.76e-08	-1.372	0.170	-2.59e-07	4.56e-08
verification_status	-0.1506	0.009	-16.320	0.000	-0.169	-0.133
dti	-0.0341	0.001	-58.722	0.000	-0.035	-0.033
delinq_2yrs	-0.1909	0.006	-32.992	0.000	-0.202	-0.180
inq_last_6mths	0.1356	0.004	32.636	0.000	0.127	0.144
mths_since_last_delinq	-0.0037	0.000	-20.036	0.000	-0.004	-0.003
open_acc	-0.0439	0.001	-37.725	0.000	-0.046	-0.042
pub_rec	-0.3223	0.009	-34.555	0.000	-0.341	-0.304
revol_bal	4.176e-06	4.11e-07	10.161	0.000	3.37e-06	4.98e-06
revol_util	-0.0075	0.000	-34.288	0.000	-0.008	-0.007
total_acc	0.0274	0.000	55.965	0.000	0.026	0.028
initial_list_status	-0.7697	0.009	-87.336	0.000	-0.787	-0.752
total_rec_int	-0.0002	3.91e-06	-48.012	0.000	-0.000	-0.000
total_rec_late_fee	-0.0087	0.001	-7.450	0.000	-0.011	-0.006
recoveries	-1.9192	115.970	-0.017	0.987	-229.216	225.378
collection_recovery_fee	-231.3463	4330.291	-0.053	0.957	-8718.560	8255.868
collections_12_mths_ex_med	-0.6024	0.043	-14.022	0.000	-0.687	-0.518
last_week_pay	0.0092	0.000	70.476	0.000	0.009	0.009
acc_now_delinq	-0.2480	0.064	-3.892	0.000	-0.373	-0.123
tot_coll_amt	-5.73e-05	4.05e-06	-14.157	0.000	-6.52e-05	-4.94e-05
tot_cur_bal	1.171e-07	3.33e-08	3.511	0.000	5.17e-08	1.82e-07
total_rev_hi_lim	-3.494e-06	3.06e-07	-11.436	0.000	-4.09e-06	-2.9e-06
grade_num	0.1660	0.002	105.933	0.000	0.163	0.169
addr_state_NE	0.0880	0.013	6.531	0.000	0.062	0.114
addr_state_S	0.1132	0.012	9.311	0.000	0.089	0.137
addr_state_W	0.2736	0.013	21.708	0.000	0.249	0.298

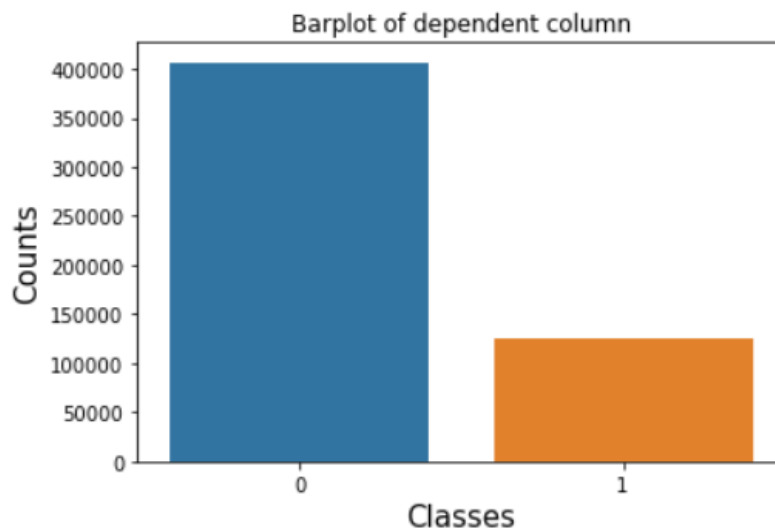
Goals as issues:

- There were 400000 rows and 118000 rows under class '0' and class '1' of loan_status respectively which resulted in Class Imbalance in the target data:
- our main goal is to predict the defaulters i.e. the 1 - in the target which is the loan status as defaulters and in equal proportions predict the not defaulters as well
- This has been attempted in many ways using different parameters eg: using value added columns, based on pvalues and on the basis of standard deviation.
- Since the target column is imbalanced we have tried to synthetically sample them using SMOTE and random sampling for both upsample and down

sample and reduce the number of false negatives and false positives in the target column while predicting in the test data

Model Tuning:

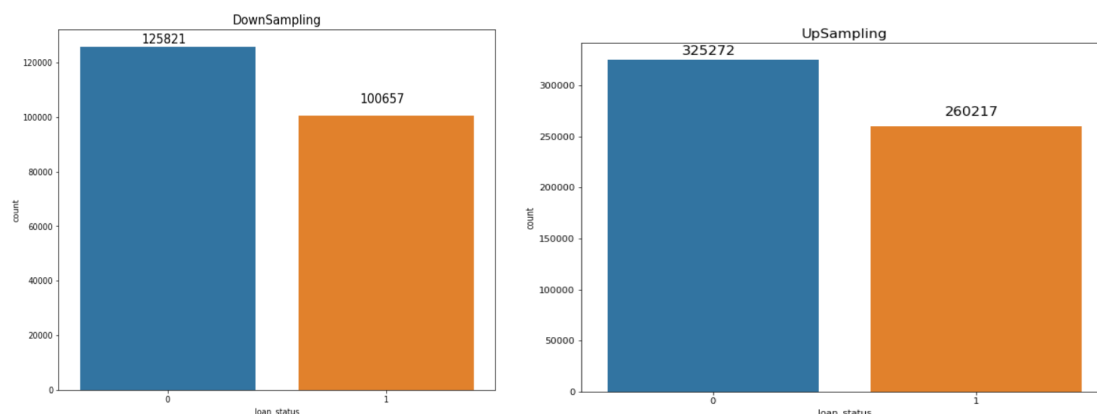
Supervised Learning Classification & Ensemble



Interpretations on the count plot:

→ We can see that the count of defaulters(1) has been represented in much less than the non defaulters(0). So building a model on this target column will give a very biased results in the favour of non defaulters that will in fact taint our objective. So we will try with upsampling and downsampling techniques using SMOTE .

SMOTE Sampling



Inference For Upsampled Target :-

- As we can see we have upsampled the defaulters in the loan_status by 80% .We have tried building a model with the upsampled target .The results are explained below using the ensembling techniques.

- After UpSampling:- 325272 Non default data
260217 Default data

Inference For Downsampled Target :-

- The Downsampling was also done with a sampling strategy of 80% .The Downsampled countplot also shows the same.Results are again tabulated further .
- After DownSampling :- 12821 Non Default data
100657 Default data

1) Model building using Value Addition Attributes:-

VA columns	Description	Relevance
Loan_to_Income_ratio	How big the loan a member has taken with respect to his income.	High loan ratio means likely to default.
Interest paid	Interest paid so far	Less interest paid implies a good credit score.
Bad state	gives a magnitude of how much the repayment has gone off course in terms of ratios/ Convert to boolean for convenience .	Bad state 0 means a person is likely to get a loan accepted.
emi_paid_progress_perc	Calculating EMIs paid (in terms of percent)	Enables a person to pay in smaller amounts than initially anticipated
Available lines	Total number of available 'Credit Lines	More available credit lines mean a person's credit score is good.
Total repayment progress	Calculating total repayments received so far in terms of EMI or recoveries after charge off	identification if a person might default in future

Formulas for the value addition columns:

- How big the loan a member has taken with respect to his income
 $\text{loan_to_income} = \text{annual_inc} / \text{funded_amnt}$
- Magnitude of how much the repayment has gone off course in terms of ratios/
Convert to boolean for convenience.

$$\begin{aligned} \text{bad_state} = & \text{acc_now_delinq} + (\text{total_rec_late_fee} / \text{funded_amnt}) + \\ & (\text{recoveries} / \text{funded_amnt}) + \\ & (\text{collection_recovery_fee} / \text{funded_amnt}) + \\ & (\text{collections_12_mths_ex_med} / \text{funded_amnt}) \end{aligned}$$

- Total number of available Credit Lines: $\text{avl_lines} = \text{total_acc} - \text{open_acc}$
- Interest paid so far: $\text{int_paid} = \text{total_rec_int} + \text{total_rec_late_fee}$
- Calculating EMIs paid (in terms of percent): $\text{emi_paid_progress_perc}] = (\text{last_week_pay} / (\text{term} / 12 * 52 + 1)) * 100$

Process For the Value added Model building :-

- A full model with all the above mentioned attributes was considered for Value added modeling.
- For efficiently predicting the target variable 'loan_status', it was downsampled using RandomOverSampler using 80% sampling strategy. (0 class and 1 class respectively). There was 125821 rows and 100657 rows under class '0' and class '1' of loan_status respectively achieved with RandomOverSampler.

Classification Report Analysis:

	Logistic_reg_up	Decision_Tree_Base	Random_forest	Adaboost_upsample	XGB_up_samp
Accuracy	0.650000	0.650000	0.680000	0.720000	0.690000
Precision	0.620000	0.580000	0.810000	0.730000	0.710000
recall	0.520000	0.730000	0.360000	0.600000	0.520000
f1_score	0.580000	0.650000	0.500000	0.660000	0.600000

Results:

- We built the above mentioned models i.e Logistic Regression, Decision Tree, XG Boost, ADA Boost, Random Forest.
- We dropped the columns having the highest VIF. (open_acc, total_rec_int,

funded_amnt, funded_amnt_inv ,int_rate ,total_acc ,last_week_pay) sequentially and checked the Accuracy,Precision,Recall and F1 Score For every model.

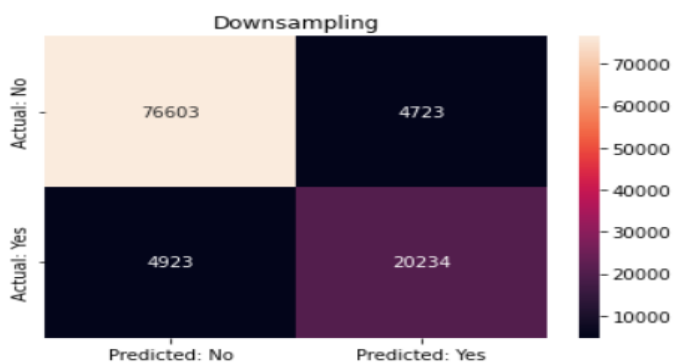
- **The best result** obtained was after dropping the columns (open_acc, total_rec_int, funded_amnt,funded_amnt_inv,int_rate,total_acc) was for the model ADA BOOST with accuracy of **72%**.

2) Model building using UpSampling and DownSampling :-

Classification Report :-

	Logistic_reg_up_samp	Decision_Tree_Base	Random_forest_down_samp	Adaboost_Down_samp	Gradient_Down_samp	XGB_Down_samp
Accuracy	70.431	84.215	84.36	77.653	81.191	90.941
Precision	0.390	0.660	0.67	0.520	0.590	0.810
recall	0.470	0.680	0.66	0.600	0.640	0.800
f1_score	0.430	0.670	0.67	0.560	0.620	0.810

Confusion Matrix:-



Summary :-

- We used different classification models
- In the ensemble techniques the best model was obtained by the XGBoost technique using downsampling. which had an F1 score of 0.81 for the default.

The best result: XGBoost Model has an accuracy of **90.94%** and an F1 score of 0.81.

3) Model Building using STD condition :-

- Since delinq_2yrs,inq_last_6mths,pub_rec,collections_12_mths_ex_med , acc_now_delinq columns had low std i.e. below 10 -i.e - 5 columns they have been dropped.
- 21 columns have been retained '
- This issue was corrected using SMOTE upsampling in order not to lose data in the train dataset.

Classification Report Analysis:

	Logistic_reg_up	Decision_Tree_Base	Random_forest	Adaboost_upsample	Gradient_up_samp	XGB_up_samp
Accuracy	0.78	0.85	0.84	0.80	0.87	0.92
Precision	0.58	0.90	0.75	0.57	0.91	0.90
recall	0.25	0.43	0.49	0.62	0.75	0.75
f1_score	0.35	0.58	0.60	0.59	0.80	0.82

The Best Hyperparameters for the Model Buliding:

ADABOOST :

learning_rate: 1.5, n_estimators: 100

GradientBoostingClassifier:

n_estimators = 100, learning_rate=0.6,max_depth = 7

XGBClassifier:

(max_depth = 10, gamma = 1)-this has the best model score.

Classification Report for the best model:

this report belongs to the *XGBoost classification* for the best results in classification so far obtained in the criterion.

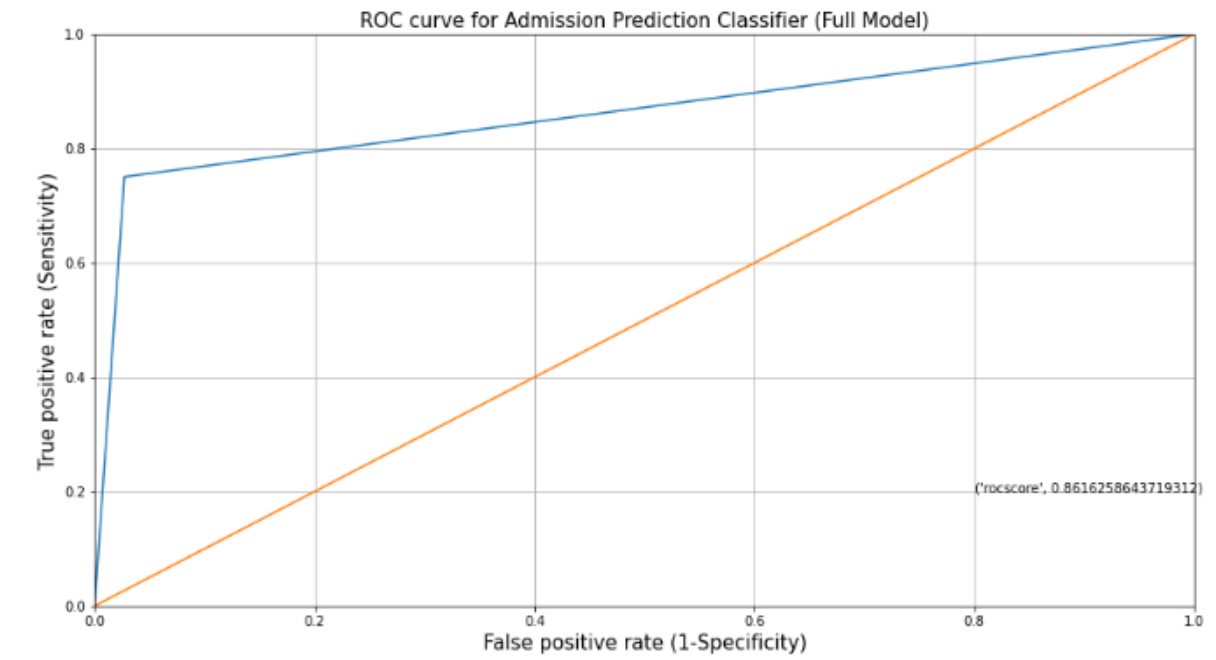
```
test classification report
              precision    recall  f1-score   support

     0           0.93       0.97       0.95         81326
     1           0.90       0.75       0.82         25157

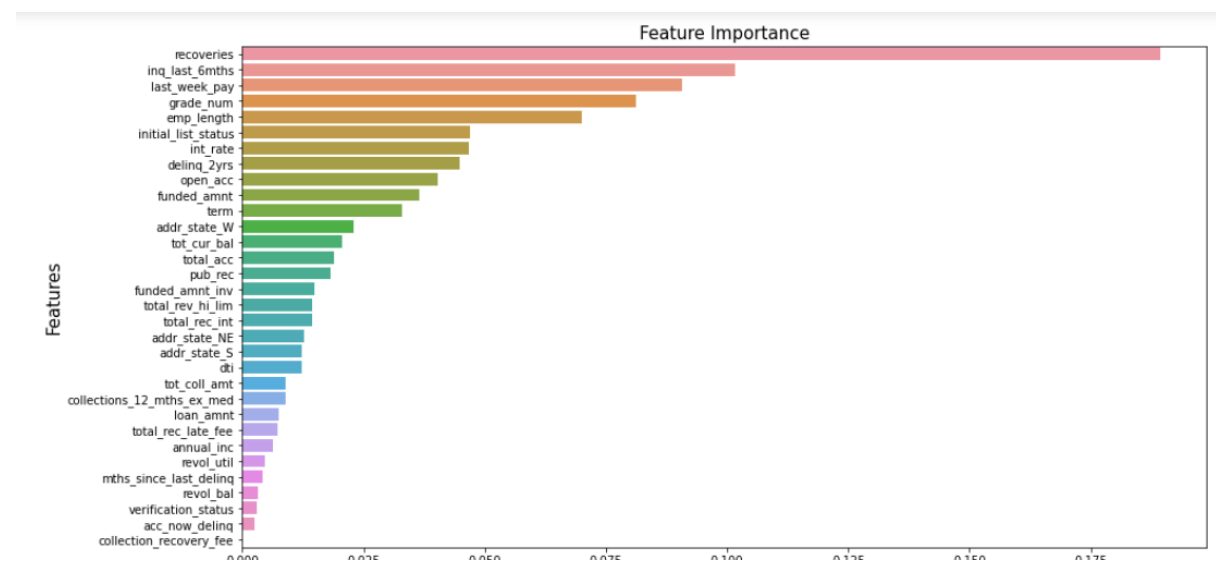
 accuracy          0.92         106483
 macro avg          0.91       0.86       0.88         106483
 weighted avg          0.92       0.92       0.92         106483
```

Here we can see the effect of the SMOTE which let us improve the scores in predicting the model by reducing the number of False negatives and false positive predictions thereby improving the overall model performance.

ROC AUC Curve and confusion matrix:



Feature Importances Plot



Inference:

The major feature that has the highest influence on the xgboost model is the recoveries column followed by inq_last_6months, last week pay and so on as seen in the above graph

Summary:

- For this situations we used different hyperparameters in different classifier models and recursively and found that the XGBoost model performed the best based on the recall,precision and f1 in the classification report
- For the ADBOost classifier with the hyperparameters learning_rate: 1.5, n_estimators: 100 the f1 score was 0.59
- For the gradient boost classifier n_estimators = 100, learning_rate=0.6,max_depth = 7 the f1 score was 0.8
- XGBClassifier (max_depth = 10, gamma = 1)-with these hyperparameters the best model score achieved was 0.82 of F1 score.

Result:-

As seen above, results from all the iterations XGB model with upsampling and considering the columns based on the standard deviation is giving the best result as compared to all the other models in the prediction of the defaulters in the test data.

So the final XGB model has overall accuracy=92%, Precision=0.93, Recall= 0.97, and F1 score=0.82

Unsupervised Learning

Problem Statement:

To apply clustering techniques and come up with business Insights for the bank so as to analyze,promote ,enhance better facilities to customers and account holders keeping in mind the welfare of the bank and its shareholders.

DBSCAN:

We have also tried the DBSCAN method of clustering. Here are the inferences we could figure from the clusters formed. After separating the defaulters .

- we can see that the largest clusters that the average employee length of of the defaulters can lie between 4- 6 years.
- annual income being in the range of 68 K - 47 K.
- their dti can be between 13-16 (ratio of member's total monthly debt repayment excluding mortgage divided by self reported monthly income) so the general standard of it being below 26-28 cannot be used.
- Also the term taken for 36 months is common among most defaulters.
- Based on the combination of the features in the below table it lets us know about the verification status.
- It is common to note that the grades of loan is quite low for the defaulters.

cluster	loan_amount mean	emp_length	annual income	dti	term	verification status	grade
A0	12772.24	5.5 years	68247	16	36 months	not verified	B3
A1	6985.394	4 years	46959.53	14.19745	36 months	not verified	B3
A2	7127.857	6 years	63362.86	13.57029	36 months	verified	C1
Outliers	18677.994	5.8 years	63362.857	15.59	36 months	verified	C3

- we tried to visualise the clusters using the scatter plot



clustering for dti vs total_rec_int

Implications :

This allows us to determine typical features and patterns of behaviour that lead to a future inability to make debt repayments. This helps the bank predict the probability of a customer to default on their loan and helps the bank to reduce its NPA(Non performing Assets) and reduce liability .

The model may fail to identify the default member who might take larger loans. Loans which can potentially risk the bank once again. If the data for default members are increased the model could predict the pattern more accurately for the default members. Also obtaining data from different banks would have an impact in the predictions.

Limitations :

Extreme cases can harm the model which could be a risk . New relevant features can add value to the prediction. Our model could struggle with datasets having small data points. As the market changes the banks NPAs could increase in the long term. Every year the model should get fresh data for training which could solve the problem.

Conclusion :

From the initial logistic regression we got the results of recall of 0.25 and f1 score of 0.35 we have improved the overall scores.

If 1 class (class '1' i.e. True Positive or defaulters) gets predicted as '0' class (class '0' i.e. False Negative), it will increase the False Negatives and will fail to prioritize the appropriate parameters which needs more turn-around-time to resolve. Since the cost of the False Negatives is high in this case, Recall score is used. Higher Recall score of '1' class (class '1' i.e. True Positive or defaulters) implies a better Model.

We have also used the F1 score to evaluate the models

It can be inferred that XGBoost Classification Model for columns removed based on standard deviation upsampled and XGBoost Classification full Model with downsampling with high Recall score of 0.75 and 0.80 they also showed F1 scores of 0.82 and 0.81 respectively are the best Model Predictors of target variable loan_status ,that showed better performance wrt to the earliest model build.