

## 9. Bank Loan Defaulter Prediction

### Writing a Data Description Report

#### Data Quantity

- What is the format of the data?

The total 3 of the given files have a common file format of Comma-separated values(.CSV)

FILE NAME	FILE FORMAT
train	.csv(Comma-separated values)
test	.csv(Comma-separated values)
submission	.csv(Comma-separated values)

- Identify the method used to capture the data--for example, ODBC.

We can't identify how the data is collected but on visualizing the data we can tell the data collection should have been taken by the banks itself from all the data they have gathered from giving out all the previous loans to the borrowers.

- How large is the database (in numbers of rows and columns)?

FILE NAME	NO. OF ROWS	NO. OF COLUMNS
train	67463	35
test	28913	34
submission	28913	1

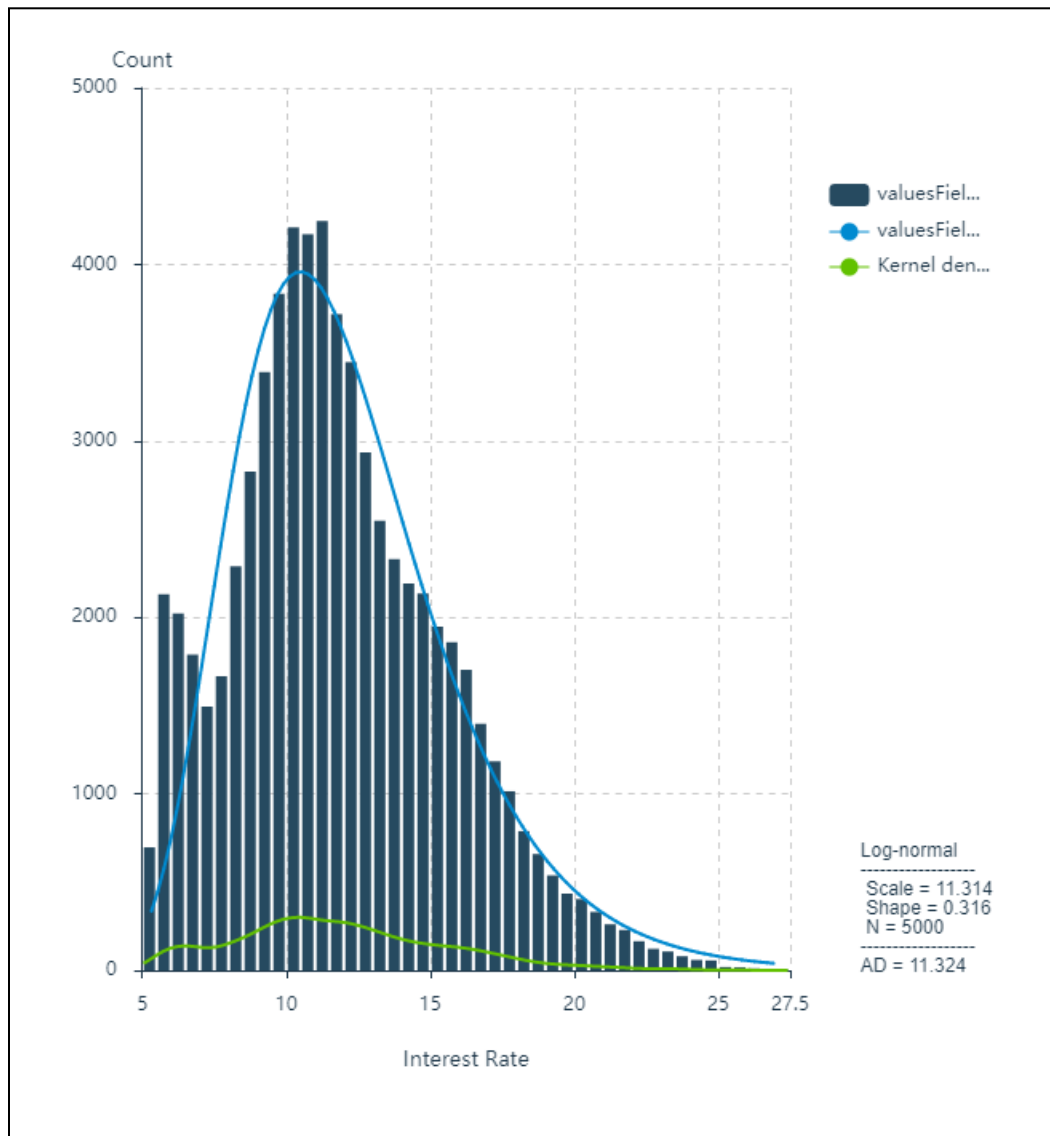
## Data Quality

- Does the data include characteristics relevant to the business question?

### Accuracy

We can check accuracy by checking given data with real life scenario:

➤ Interest Rates Range given in dataset:



➤ Current Interest Rate Range in Real Life by banks:

Bank	Interest Rate(p.a.)
HDFC Bank	10.5% p.a. - 21.00% p.a.
ICICI Bank	10.75% p.a. - 19.00% p.a.
Yes Bank	10% p.a. onwards - 24% p.a.
Citibank	10.50% p.a. - 16.49% p.a.
Kotak Mahindra Bank	10.99% and above
Axis Bank	12% p.a.- 21% p.a.
IndusInd Bank	10.49% p.a. - 26.5% p.a.
Home Credit Cash Loan	24% p.a. - 49.5% p.a.
Aditya Birla Capital	14% p.a. -26% p.a.
State Bank of India	10.65% p.a. - 13.65% p.a.
Federal Bank	10.49% p.a. - 17.99% p.a.
IIFL	11.75% p.a. - 34% p.a.

### Completeness

- **id**: Unique ID of the loan application.
- **grade**: LC assigned loan grade.
- **annual\_inc**: The self-reported annual income provided by the borrower during registration.
- **loan\_amount**: Total loan amount given to the borrower.
- **interest\_rate**: Interest rate at which a particular loan is given.
- **home\_ownership**: Type of home ownership.
- **dti (Debt-To-Income Ratio)**: A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage

and the requested LC loan, divided by the borrower's self-reported monthly income.

- **loan\_title:** A category provided by the borrower for the loan request.
- **term:** The number of payments on the loan. Values are in months and can be either 36 or 60.
- **delinquency:** 1 when the borrower had at least one event of delinquency in two years.
- **revolving\_utilities:** Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.
- **total\_received\_late\_fee:** Late fees received to date.
- **total\_current\_balance:** Total balance in the bank of a particular borrower.
- **recoveries:** The recovery rate is the estimated percent of a loan or an obligation that will still be repaid to creditors in the event of a default or bankruptcy.

### Reliability

In the realm of data quality characteristics, reliability means that a piece of information doesn't contradict another piece of information in a different source or system.

The data provided to us doesn't have contradicted pieces of information.

### Relevance

ATTRIBUTES	DESCRIPTION
ID	Unique ID of the loan application.
Loan Amount	Total loan amount given to the borrower.

## Predictive Modelling Project

---

Funded Amount	The aggregate amount of Purchase Prices paid by the Banks
Funded Amount Investor	Aggregate amount
Term	Period of time for which loan is taken
Batch Enrolled	Loan enrolled in a batch
Interest Rate	The amount charged over and above the principal amount by the lender from the borrower
Grade	A particular level of rank, quality, proficiency, or value
Sub Grade	Subset for grades
Employment Duration	Home ownership type
Home Ownership	Amount for home the borrower is paying
Verification Status	Status for verification of a loan
Payment Plan	Plan for repayment of loan
Loan Title	A category provided by the borrower for the loan request.
Debt to Income	A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower's self-reported monthly income
Delinquency - two years	1 when the borrower had at least one event of delinquency in two years
Inquires - six months	If borrower has looked into getting any loans in last six months
Open Account	Total accounts opened for a particular borrower
Public Record	Documents or pieces of information that are not considered confidential
Revolving Balance	A line of credit they can keep using and repaying over and over
Revolving Utilities	How much of credit balance has been used
Total Accounts	Total no. of accounts
Initial List Status	A request by a Borrower to borrow money on the terms of a Loan
Total Received Interest	Principal loan amount x interest rate x loan term
Total Received Late Fee	A late fee is a charge imposed on a consumer who fails to make the payment on a debt or other financial obligation by the due date

## Predictive Modelling Project

---

Recoveries	The recovery rate is the estimated percent of a loan or an obligation that will still be repaid to creditors in the event of a default or bankruptcy
Collection Recovery Fee	Covers costs incurred to collect tax and fee liabilities that are unpaid for more than 90 days
Collection 12 months Medical	Paying the debt in an agreed-on timeframe, no interest will be charged
Application Type	Whether a borrower is identified as individual or asking for loan in joint status
Last week Pay	Payment received in last week
Accounts Delinquent	A delinquency is when your account is late by 30 days or more.
Total Collection Amount	Total amount collected from a particular borrower.
Total Current Balance	Total current balance in a borrower's account.
Total Revolving Credit Limit	A line of credit that remains available over time, even if you pay the full balance.
Loan Status	Indicates where your loan is in the process.

- What data types are present (symbolic, numeric, etc.)?

ATTRIBUTES	DATA TYPES
ID	Integer
Loan Amount	Integer
Funded Amount	Integer
Funded Amount Investor	Real
Term	Integer
Batch Enrolled	String
Interest Rate	Real
Grade	String
Sub Grade	String
Employment Duration	String
Home Ownership	Real
Verification Status	String
Payment Plan	String

## Predictive Modelling Project

---

Loan Title	String
Debt to Income	Real
Delinquency - two years	Integer
Inquires - six months	Integer
Open Account	Integer
Public Record	Integer
Revolving Balance	Integer
Revolving Utilities	Real
Total Accounts	Integer
Initial List Status	String
Total Received Interest	Real
Total Received Late Fee	Real
Recoveries	Real
Collection Recovery Fee	Real
Collection 12 months Medical	Integer
Application Type	String
Last week Pay	Integer
Accounts Delinquent	Integer
Total Collection Amount	Integer
Total Current Balance	Integer
Total Revolving Credit Limit	Integer
Loan Status	Integer

## Predictive Modelling Project

- Did you compute basic statistics for the key attributes? What insight did this provide into the business question?

View data from Filter Node

<



# Predictive Modelling Project

View data from Filter Node

	Audit	Quality	Statistics	Pearson Correlations							
	MEAN	STD. DEV	VARIANCE	MEAN STD. ERR.	SKEWNESS	SKEWNESS STD. ERR.	KURTOSIS	KURTOSIS STD. ERR.	UNIQUE	VALID	
346	25,627,607.746	21,091,554.024	444,853,651,136,038.06	81,203.688	0.557	0.009	-1.09	0.019	--	67,463	
	--	--	--	--	--	--	--	--	7	67,463	
	15,770.599	8,150.993	66,438,681.38	31.382	0.673	0.009	-0.617	0.019	--	67,463	
	16,848.903	8,367.866	70,021,176.815	32.217	0.288	0.009	-0.798	0.019	--	67,463	
	11.846	3.719	13.828	0.014	0.563	0.009	0.149	0.019	--	67,463	
	--	--	--	--	--	--	--	--	3	67,463	
	23.299	8.452	71.433	0.033	0.081	0.009	-0.905	0.019	--	67,463	
	--	--	--	--	--	--	--	--	109	67,463	
	58.174	3.327	11.072	0.013	-6.381	0.009	39.593	0.019	--	67,463	
	0.327	0.801	0.641	0.003	4.635	0.009	30.673	0.019	--	67,463	
	52.889	22.539	508.027	0.087	-0.237	0.009	-0.545	0.019	--	67,463	
	1.144	5.244	27.503	0.02	5.084	0.009	25.99	0.019	--	67,463	
5	159,573.934	139,033.246	19,330,243,372.283	535.286	1.512	0.009	3.125	0.019	--	67,463	
67	59.692	357.026	127,467.812	1.375	7.371	0.009	58.177	0.019	--	67,463	

- Are you able to prioritize relevant attributes? If not, are business analysts available to provide further insight?

**id:** Unique ID of the loan application.

**grade:** LC assigned loan grade.

**annual\_inc:** The self-reported annual income provided by the borrower during registration.

**loan\_amount:** Total loan amount given to the borrower.

<b>interest_rate:</b> Interest rate at which a particular loan is given.
<b>home_ownership:</b> Type of home ownership.
<b>dti (Debt-To-Income Ratio):</b> A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower's self-reported monthly income.
<b>loan_title:</b> A category provided by the borrower for the loan request.
<b>term:</b> The number of payments on the loan. Values are in months and can be either 36 or 60.
<b>delinquency:</b> 1 when the borrower had at least one event of delinquency in two years.
<b>revolving_utilities:</b> Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.
<b>total_received_late_fee:</b> Late fees received to date.
<b>total_current_balance:</b> Total balance in the bank of a particular borrower.
<b>recoveries:</b> The recovery rate is the estimated percent of a loan or an obligation that will still be repaid to creditors in the event of a default or bankruptcy.

## Data Exploration Report

- What sort of hypotheses have you formed about the data?
  1. The people who are not verified are more likely to default on their loan amounts.
  2. The people who have high debt income are more likely to default on their loan amounts.
  3. Loan amount is a significant predictor of loan default. Customers with higher loan amounts are more likely to default on their loans compared to those with lower loan amounts.
  4. Loan term is a significant predictor of loan default. Customers with longer loan terms are more likely to default on their loans compared to those with shorter loan terms.
  5. Employment status is a significant predictor of loan default. Customers who are unemployed or self-employed are more likely to default on their loans compared to those who are employed or in other employment status categories.
- Which attributes seem promising for further analysis?

<b>id:</b> Unique ID of the loan application.
<b>loan_amount:</b> Total loan amount given to the borrower.
<b>interest_rate:</b> Interest rate at which a particular loan is given.
<b>home_ownership:</b> Type of home ownership.
<b>dti (Debt-To-Income Ratio):</b> A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower's self-reported monthly income.

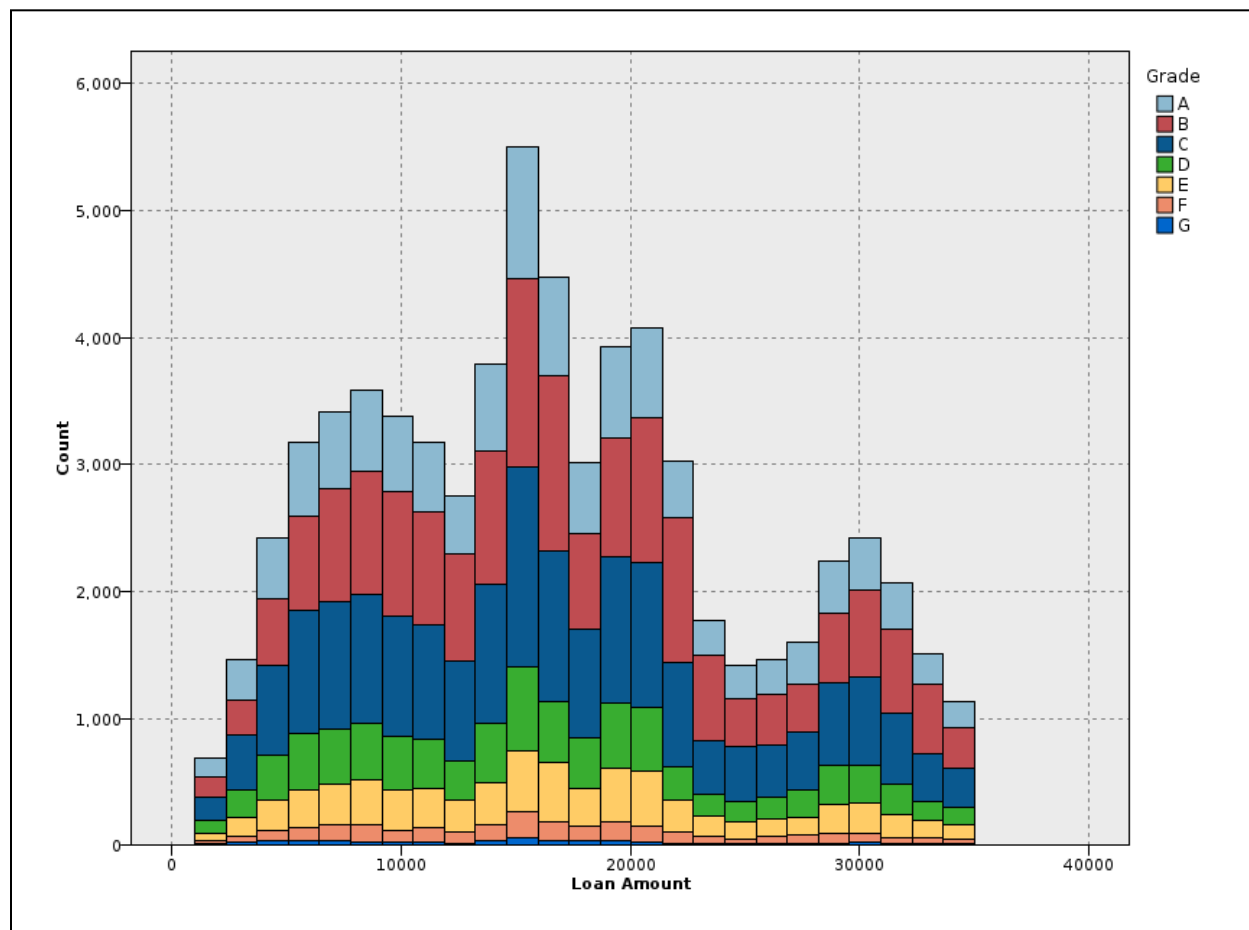
<b>loan_title:</b> A category provided by the borrower for the loan request.
<b>verification_status:</b> Status for verification of a loan.
<b>term:</b> The number of payments on the loan. Values are in months and can be either 36 or 60.
<b>delinquency:</b> 1 when the borrower had at least one event of delinquency in two years.
<b>total_received_late_fee:</b> Late fees received to date.
<b>total_current_balance:</b> Total balance in the bank of a particular borrower.
<b>recoveries:</b> The recovery rate is the estimated percent of a loan or an obligation that will still be repaid to creditors in the event of a default or bankruptcy.

- Have your explorations revealed new characteristics about the data?

### Univariate Analysis:

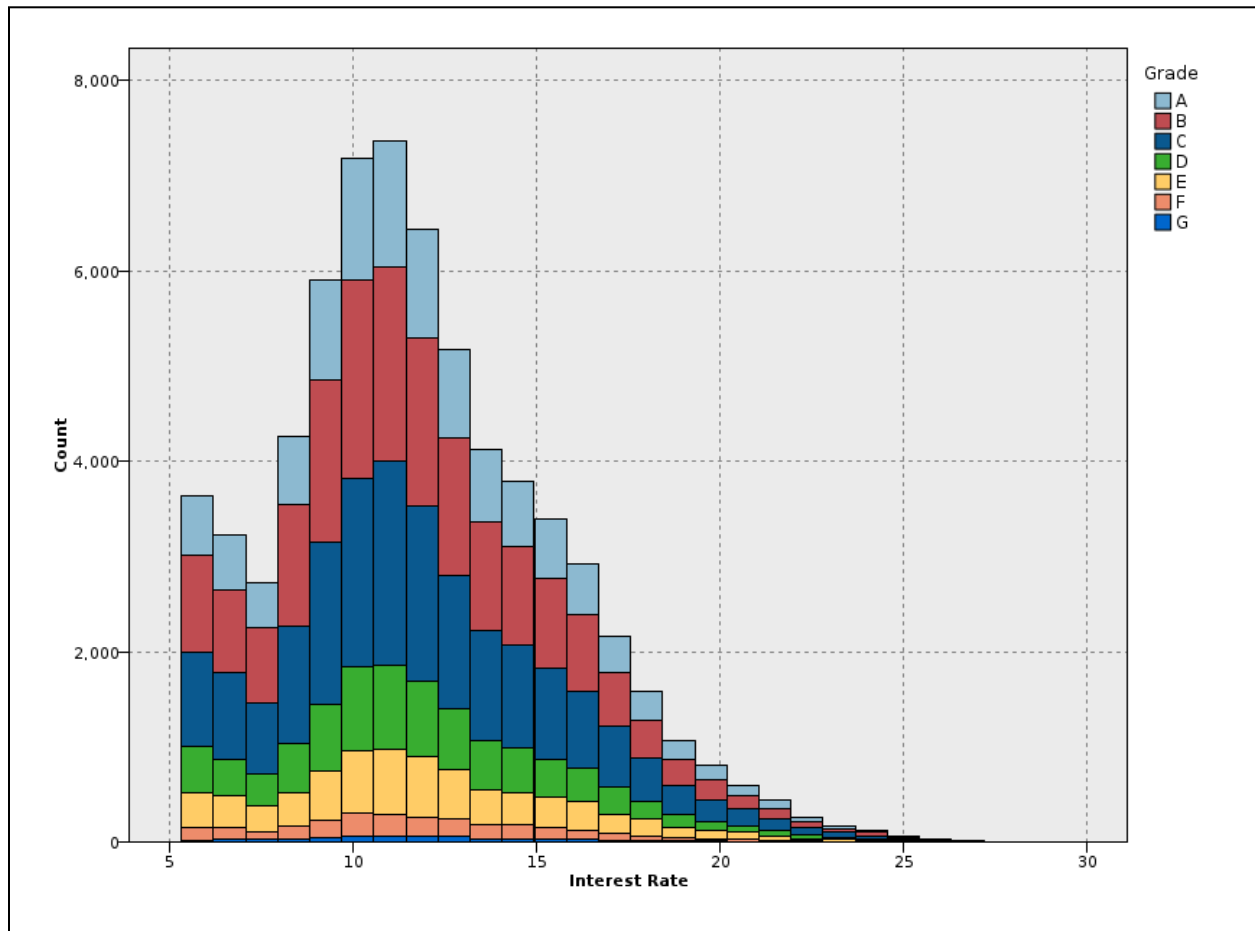
#### 1. Loan Amount:

*Insights: Most of the loan amounts are distributed between 10000 to 20000.*



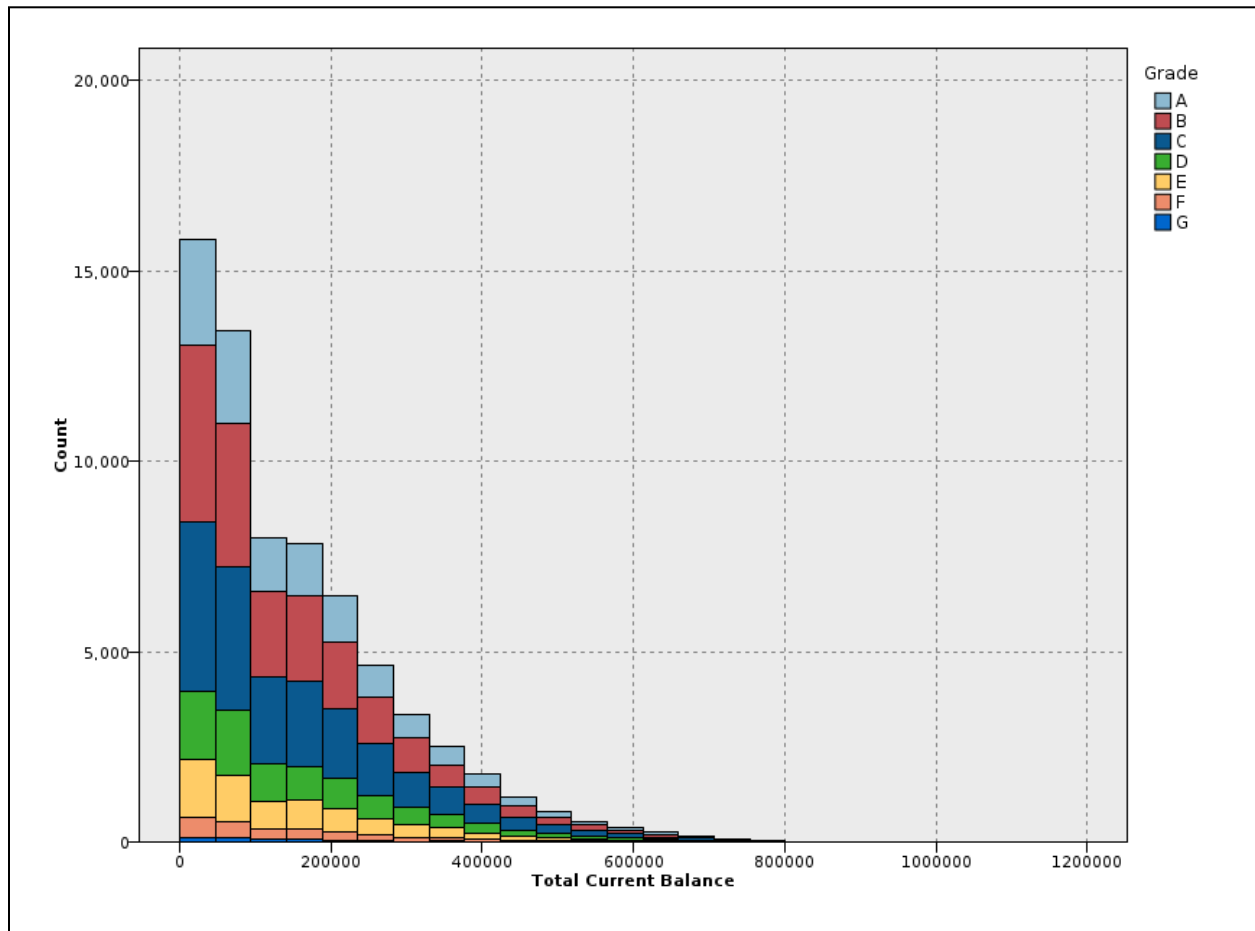
## 2. Interest Rate:

*Insights: Most of the loans interest rates are distributed between 10% to 13%.*



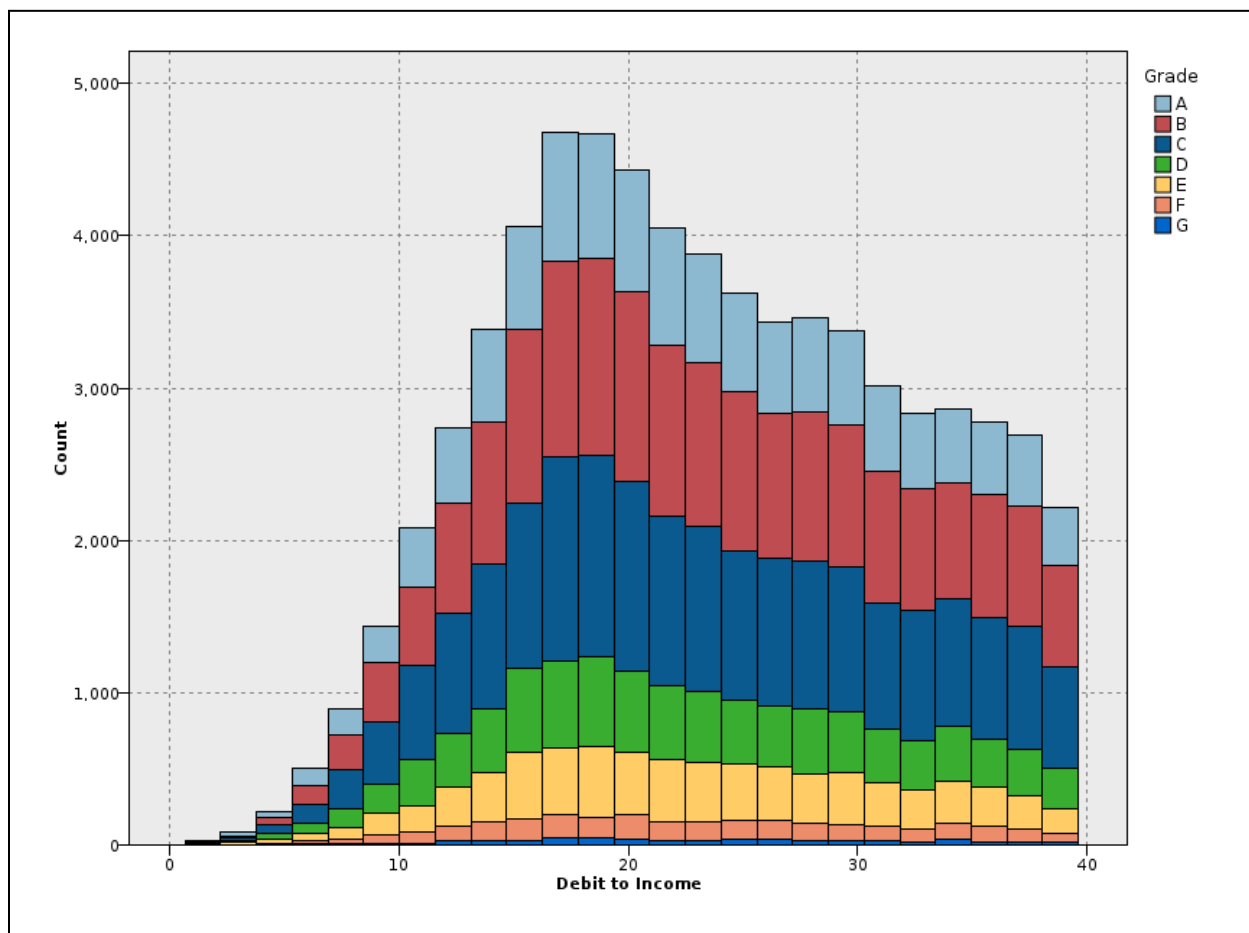
## 3. Total Current Balance:

*Insights: Most of the applicants have total balance less than 200000.*



## 4. dti (Debt to Income Ratio):

*Insights: Most of the applicants have a debt-to-income (dti) ratio between 15 to 20.*

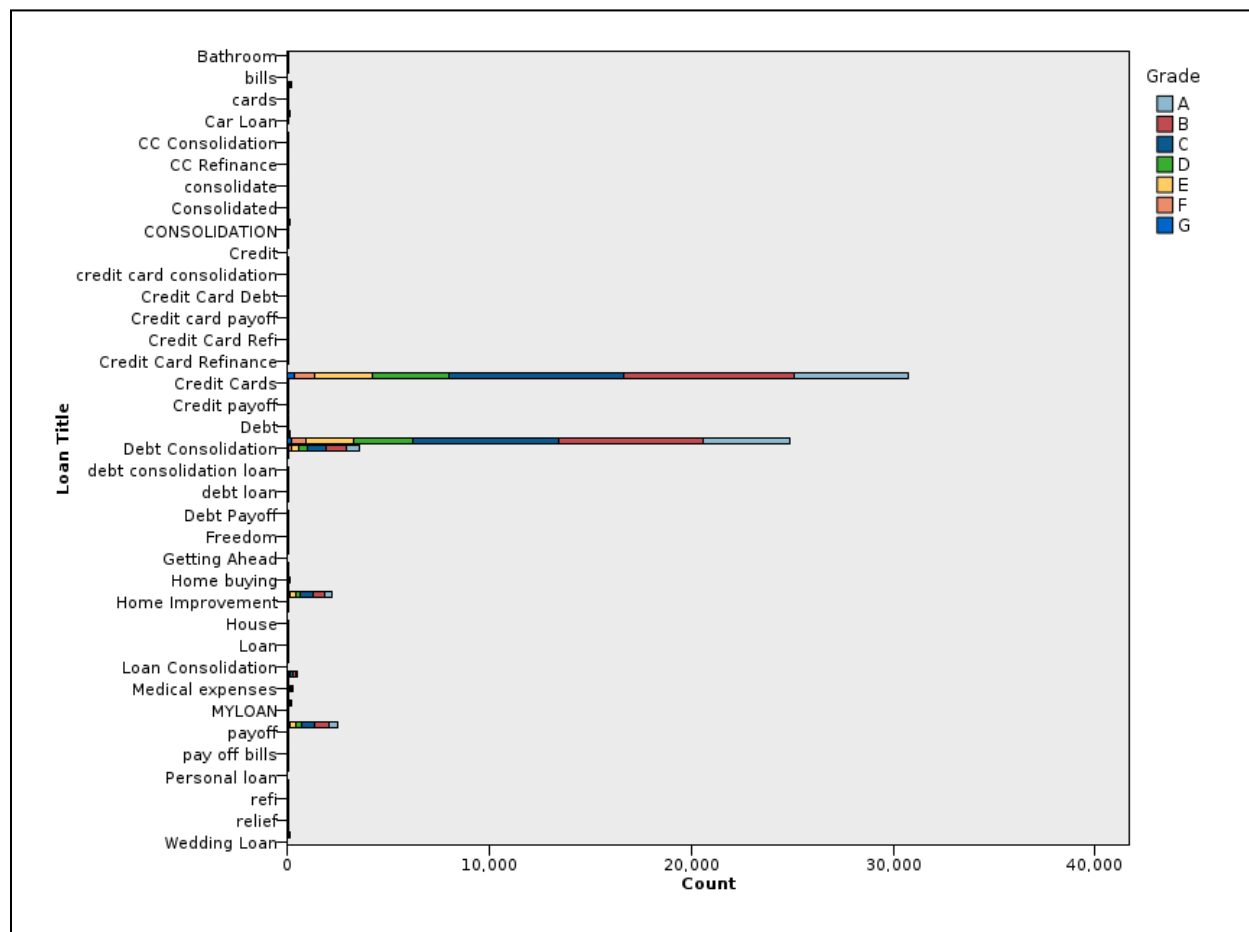




## Categorical Variables:

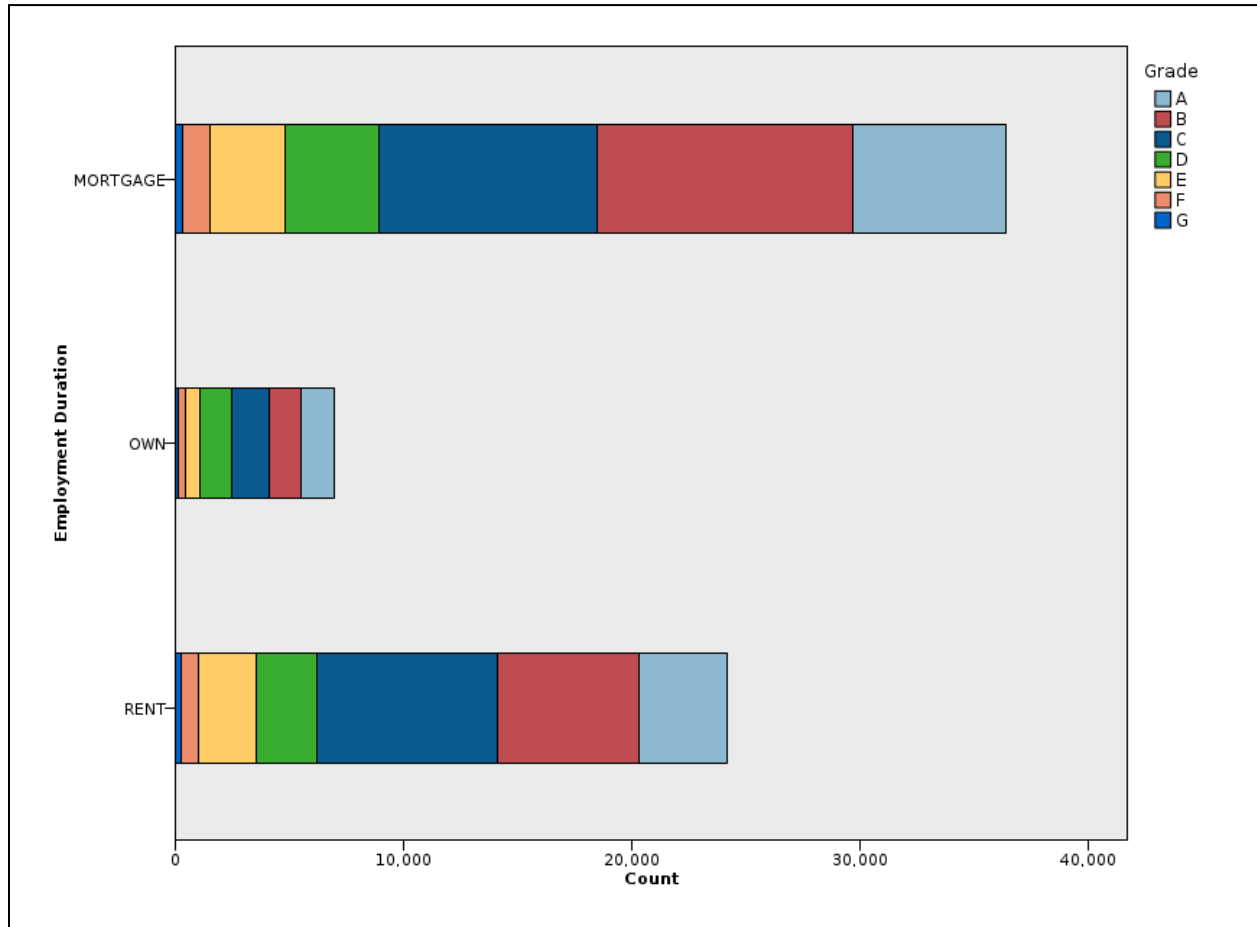
### 5. Purpose of Loan:

*Insights: Insights: Approx. 60% of the applicants applied for a loan for paying their other loans and credit card bills (Debt Consolidation).*



### 6. Home Ownership:

*Insights: 30% of applicants are living in rented homes whereas 60% applicants were mortgaged their home.*



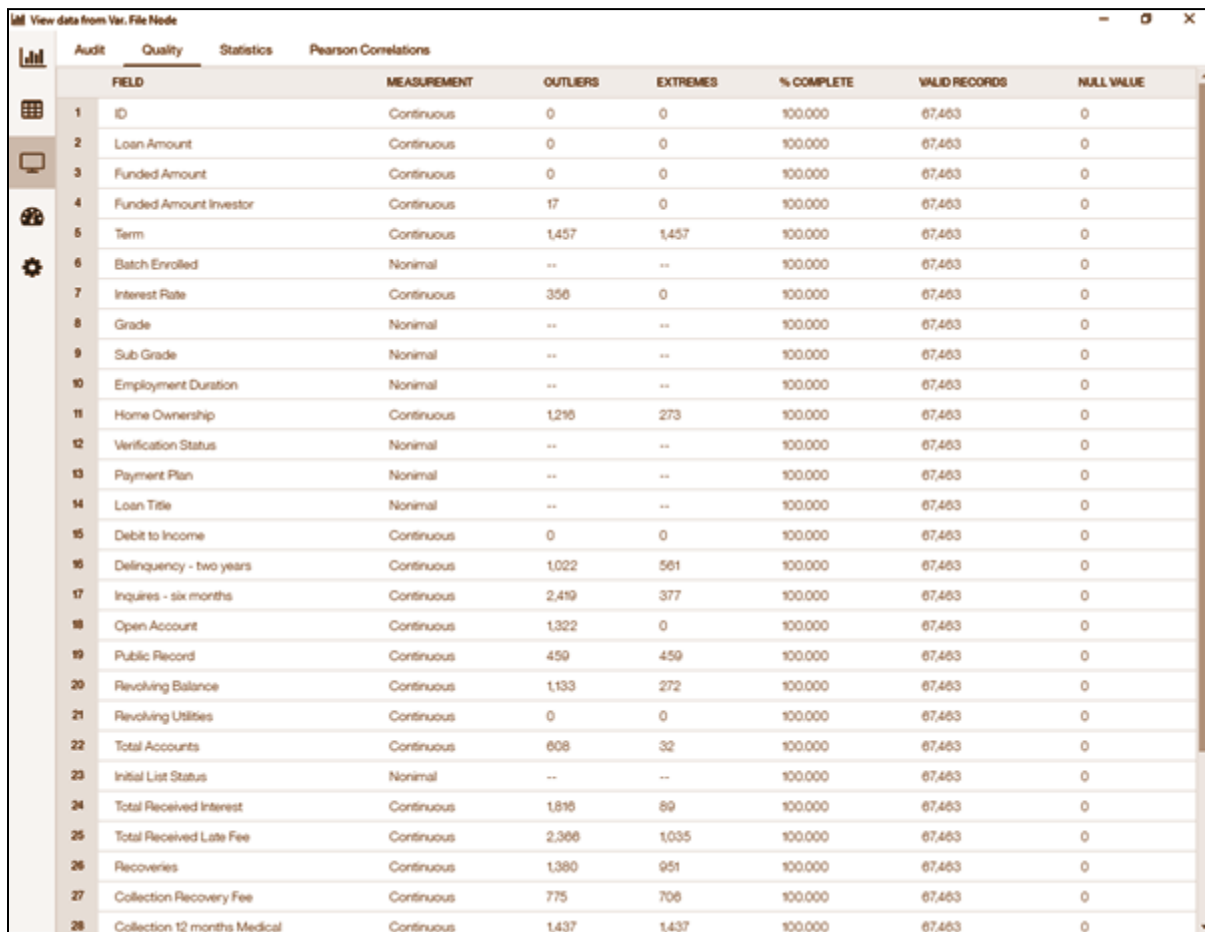
- Can you identify particular subsets of data for later use?
  - Subsets of people with higher & lower income levels.
  - Subset of data in which a person has its own home ownership.
  - Subset of person with verification & non-verification status.
  - Subset of inquiries of the last six months.
  - Subsets of whether recoveries are high are not.

- Take another look at your data mining goals. Has this exploration altered the goals?
  - No, because we have to predict that a person can pay the loan back or not & after visualizing the data we came to the conclusion that it doesn't alter the main goal.

## Writing a Data Quality Report

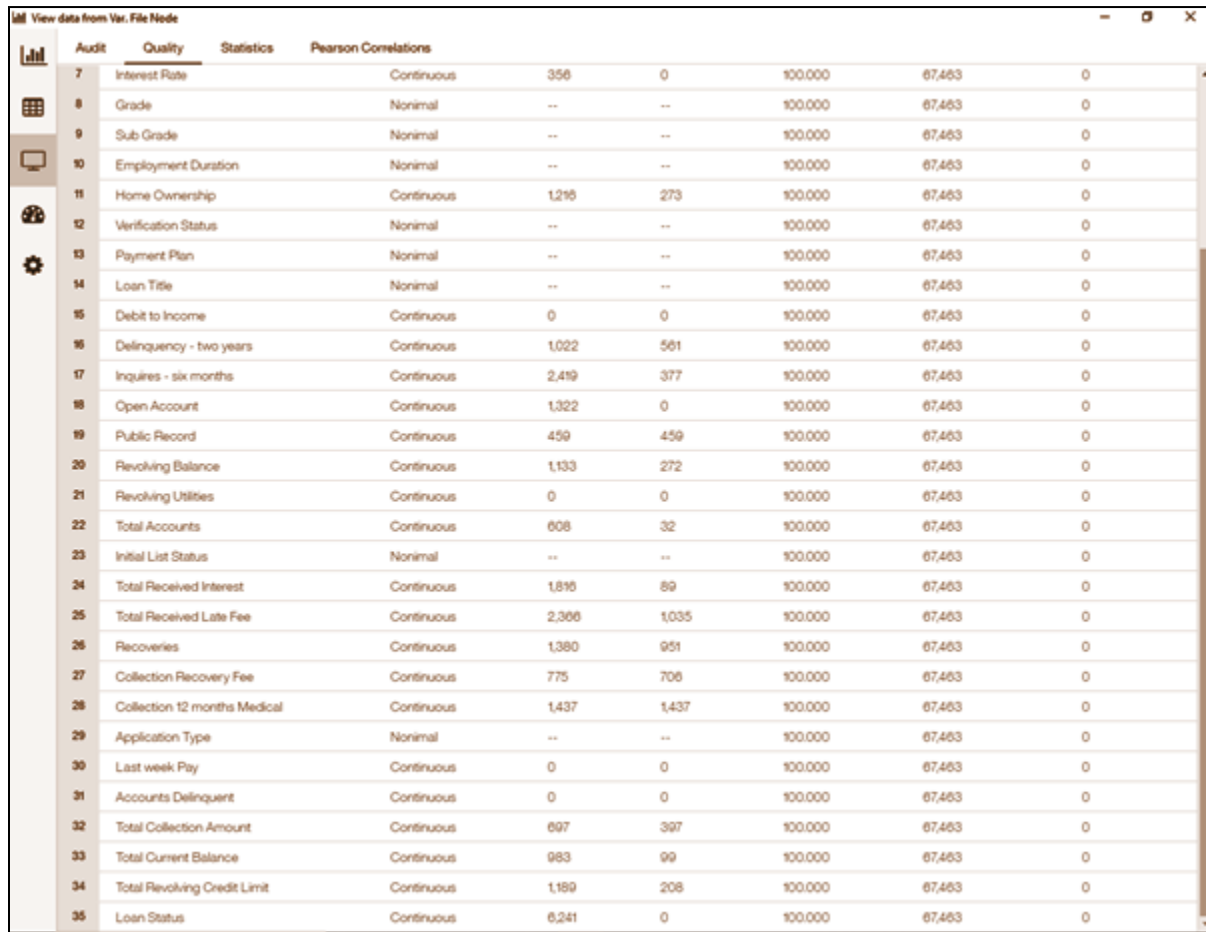
- Have you identified missing attributes and blank fields? If so, is there meaning behind such missing values?

No we haven't identified any missing values or blank fields in any attributes



	FIELD	MEASUREMENT	OUTLIERS	EXTREMES	% COMPLETE	VALID RECORDS	NULL VALUE
1	ID	Continuous	0	0	100.000	67,463	0
2	Loan Amount	Continuous	0	0	100.000	67,463	0
3	Funded Amount	Continuous	0	0	100.000	67,463	0
4	Funded Amount Investor	Continuous	17	0	100.000	67,463	0
5	Term	Continuous	1,457	1,457	100.000	67,463	0
6	Batch Enrolled	Nonimal	--	--	100.000	67,463	0
7	Interest Rate	Continuous	356	0	100.000	67,463	0
8	Grade	Nonimal	--	--	100.000	67,463	0
9	Sub Grade	Nonimal	--	--	100.000	67,463	0
10	Employment Duration	Nonimal	--	--	100.000	67,463	0
11	Home Ownership	Continuous	1,216	273	100.000	67,463	0
12	Verification Status	Nonimal	--	--	100.000	67,463	0
13	Payment Plan	Nonimal	--	--	100.000	67,463	0
14	Loan Title	Nonimal	--	--	100.000	67,463	0
15	Debit to Income	Continuous	0	0	100.000	67,463	0
16	Delinquency - two years	Continuous	1,022	561	100.000	67,463	0
17	Inquires - six months	Continuous	2,419	377	100.000	67,463	0
18	Open Account	Continuous	1,322	0	100.000	67,463	0
19	Public Record	Continuous	459	459	100.000	67,463	0
20	Revolving Balance	Continuous	1,133	272	100.000	67,463	0
21	Revolving Utilities	Continuous	0	0	100.000	67,463	0
22	Total Accounts	Continuous	608	32	100.000	67,463	0
23	Initial List Status	Nonimal	--	--	100.000	67,463	0
24	Total Received Interest	Continuous	1,816	89	100.000	67,463	0
25	Total Received Late Fee	Continuous	2,366	1,035	100.000	67,463	0
26	Recoveries	Continuous	1,380	951	100.000	67,463	0
27	Collection Recovery Fee	Continuous	775	706	100.000	67,463	0
28	Collection 12 months Medical	Continuous	1,437	1,437	100.000	67,463	0

## Predictive Modelling Project



The screenshot shows a window titled "View data from Var. File Node". On the left is a sidebar with icons for a bar chart, a grid, a monitor, a network, and a gear. The main area displays a table with the following columns: "Audit", "Quality", "Statistics", and "Pearson Correlations". The table lists 35 variables, each with a row number, a name, a data type, and several numerical values. The "Statistics" column contains two values for each variable, and the "Pearson Correlations" column contains two values. The values for "Statistics" and "Pearson Correlations" are consistently 0 or 100.000 across all rows. The "Quality" column contains values like "Continuous" or "Nonimal" (note the spelling).

	Audit	Quality	Statistics	Pearson Correlations
7	Interest Rate	Continuous	356 0	100.000 67,463 0
8	Grade	Nonimal	-- --	100.000 67,463 0
9	Sub Grade	Nonimal	-- --	100.000 67,463 0
10	Employment Duration	Nonimal	-- --	100.000 67,463 0
11	Home Ownership	Continuous	1,216 273	100.000 67,463 0
12	Verification Status	Nonimal	-- --	100.000 67,463 0
13	Payment Plan	Nonimal	-- --	100.000 67,463 0
14	Loan Title	Nonimal	-- --	100.000 67,463 0
15	Debit to Income	Continuous	0 0	100.000 67,463 0
16	Delinquency - two years	Continuous	1,022 561	100.000 67,463 0
17	Inquiries - six months	Continuous	2,419 377	100.000 67,463 0
18	Open Account	Continuous	1,322 0	100.000 67,463 0
19	Public Record	Continuous	459 459	100.000 67,463 0
20	Revolving Balance	Continuous	1,133 272	100.000 67,463 0
21	Revolving Utilities	Continuous	0 0	100.000 67,463 0
22	Total Accounts	Continuous	608 32	100.000 67,463 0
23	Initial List Status	Nonimal	-- --	100.000 67,463 0
24	Total Received Interest	Continuous	1,816 89	100.000 67,463 0
25	Total Received Late Fee	Continuous	2,366 1,035	100.000 67,463 0
26	Recoveries	Continuous	1,380 951	100.000 67,463 0
27	Collection Recovery Fee	Continuous	775 706	100.000 67,463 0
28	Collection 12 months Medical	Continuous	1,437 1,437	100.000 67,463 0
29	Application Type	Nonimal	-- --	100.000 67,463 0
30	Last week Pay	Continuous	0 0	100.000 67,463 0
31	Accounts Delinquent	Continuous	0 0	100.000 67,463 0
32	Total Collection Amount	Continuous	697 397	100.000 67,463 0
33	Total Current Balance	Continuous	983 99	100.000 67,463 0
34	Total Revolving Credit Limit	Continuous	1,189 208	100.000 67,463 0
35	Loan Status	Continuous	6,241 0	100.000 67,463 0

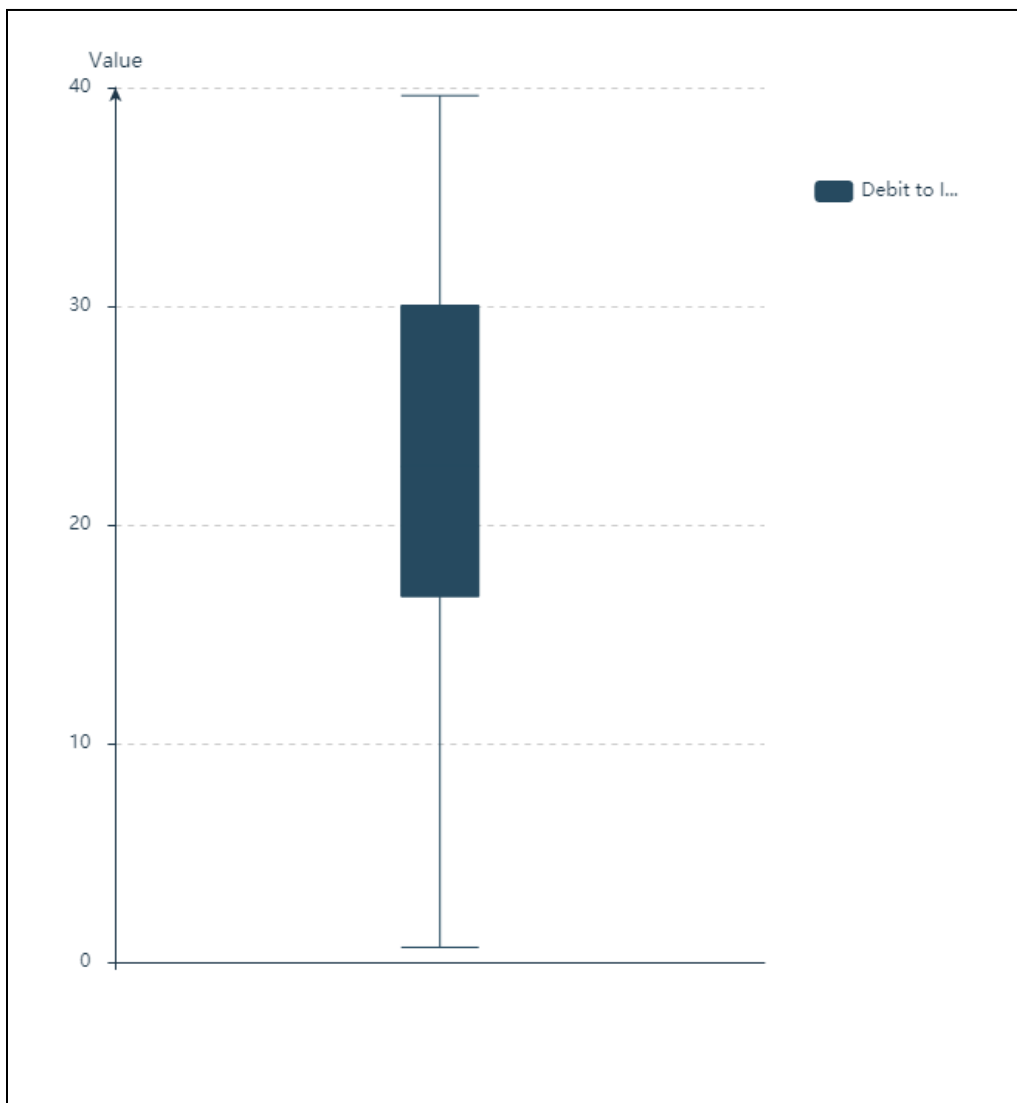
- Are there spelling inconsistencies that may cause problems in later merges or transformations?

There are no spelling inconsistencies in the attributes so it will not cause any problems in later merger or transformation.

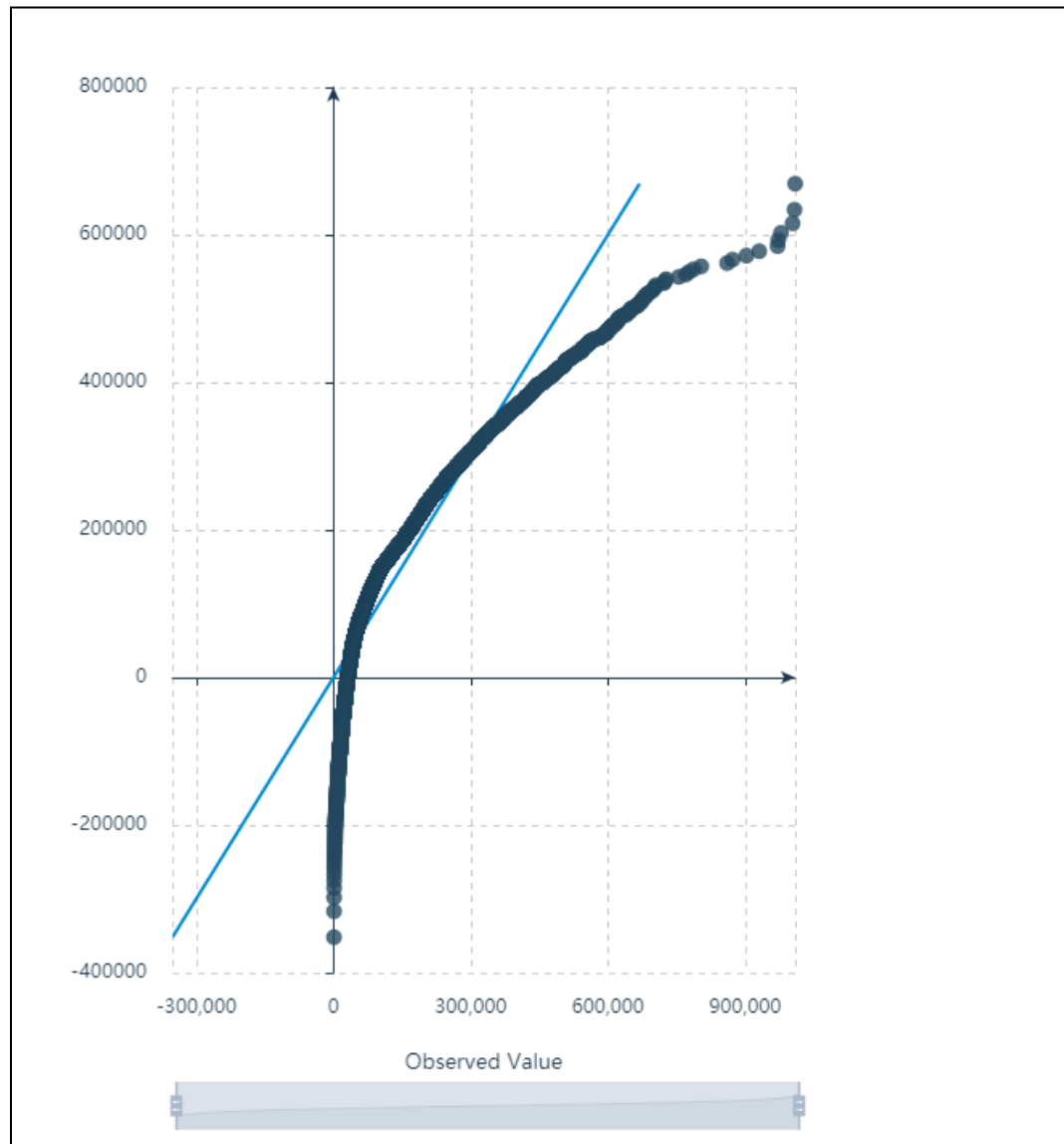
- Have you explored deviations to determine whether they are "noise" or phenomena worth analyzing further?

Yes, we explored datasets where there are outliers and extremes are present in some columns but they are valid outliers, all of these outliers and extremes didn't affect the further analysis.

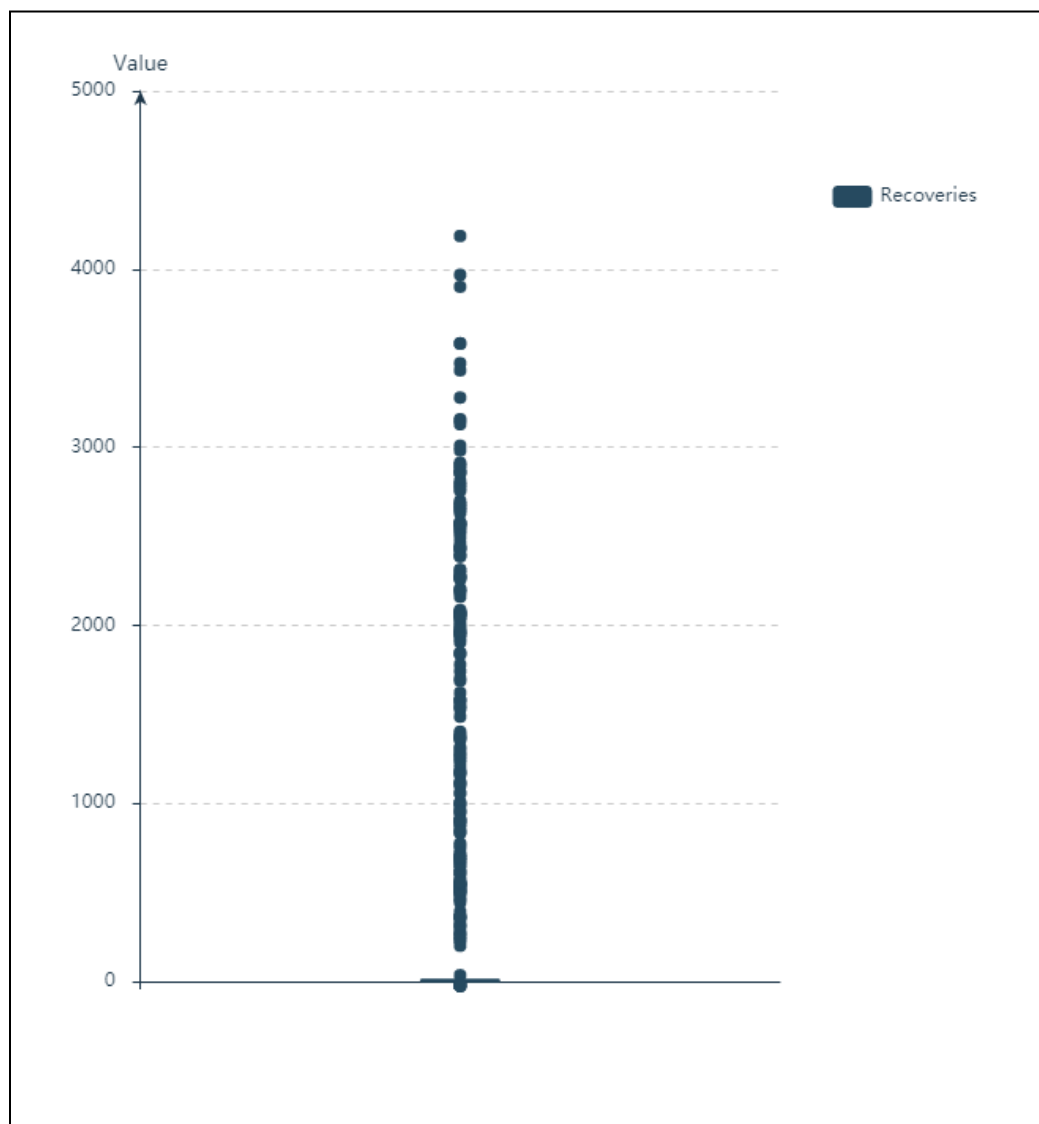
➤ Box plot for debt to income



- Q-Q plot of total\_current\_balance



➤ Box plot for recoveries attribute



- Have you conducted a plausibility check for values? Take notes on any apparent conflicts (such as teenagers with high income levels).

Yes, there is plausibility in our data.



# Predictive Modelling Project

- Have you considered excluding data that has no impact on your hypotheses?

Yes , we have considered excluding the data that has no impact on this hypothesis.

- Is the data stored in flat files? If so, are the delimiters consistent among files? Does each record contain the same number of fields?

Yes , data is stored in a flat file and comma delimiter is consistent in the file.

Yes , each record contains the same number of fields.

File: C:\Usershp\Downloads\train.csv

ID	Loan Amount	Funded Amount	Funded Amount	Investor Term	Batch Enrolled	Interest Rate	Grade	Sub Grade	Employment Duration	Home Ownership	Verification Status	Payment Plan	Plan L
45087372	10000	12236	12329	34256	59	BAT2522922	11	13500484	B,C4	MORTGAGE	174346	4267	Not Verified,n,Debt Consolidation,16.28475781,1,0,13,0,24246,74.93255103,7,w,2929.646315,0.1205519800000001,2.495290942,0.79372376,0,INDIVIDUAL,49,0,31,311301,6
1450153	3609	11940	12191	99492	59	BAT1564599	12	23756263	C,D3	RENT	39033	921	Source Verified,n,Debt Consolidation,15.41240945,0,0,12,0,812,78.29718401,13,f,772.7693549999999,0.03618117,2.377214756,0.974821103,0,INDIVIDUAL,109,0,53,182610,20885,0
1949101	28276	9311	21603	22455	59	BAT2134391	12	54508365	F,D4	MORTGAGE	91504	69105	Source Verified,n,Debt Consolidation,28.13761842,0,0,14,0,1843,2.073039673,20,w,863.3243
4651430	11170	4954	17877	15555	59	BAT2428731	16	73120146	C,C3	MORTGAGE	108284	5759	Source Verified,n,Debt Consolidation,18.04373003,1,0,7,0,13919,67.46795121,12,w,288.1731
14354449	14890	13226	13539	924469999999	59	BAT5341619	15	00830022	C,D4	MORTGAGE	44234	82545	Source Verified,n,Credit card refinancing,17.20988404,1,3,13,1,1544,65.25076114
50509044	34631	30203	5635	931612999999	36	BAT4694572	17	24498602	B,G5	RENT	95957	47541000001	Not Verified,n,Credit card refinancing,7.914332742000001,3,2,16,0,2277,51.5644
32737431	30544	19773	15777	511830000001	59	BAT4808022	10	73143231	C,C5	RENT	102391	8243	Verified,n,Home improvement,15.08391059,0,0,11,0,14501,46.80880397,37,w,525.7381086
43151450	20744	10409	7645	014802	55	BAT2558388	13	99368784	A,A5	OWN	61723	52014	Not Verified,n,Debt Consolidation,29.82971503,0,0,14,0,13067,23.93462444,33,w,1350.24521199
4279642	9299	11235	13429	45641	59	BAT5341619	11	17645655	G,C2	MORTGAGE	63205	09072000001	Verified,n,Credit card refinancing,26.24471043,0,0,6,0,549,15.94738559,17,w,4140.1
4431034	19232	8942	7004	097481	55	BAT2078974	5	520412677	C,B5	RENT	42015	465840000004	Source Verified,n,Credit card refinancing,10.04854906,1,0,11,0,1361,35.07334484,30,f,2
22010590	5640	6425	8169	779401000001	59	BAT2252229	8	627334502999999	B,A5	RENT	41858	88652	Verified,n,Debt Consolidation,14.26647349,0,0,17,0,4005,61.27955642,46,f,107.890
3975542	16581	8767	10637	04903	59	BAT2333412	9	535987839	A,D4	MORTGAGE	39405	50605	Source Verified,n,Credit Consolidation,32.5409807,1,0,10,2,8266,25.62485547,21,w,1276.25
1793329	31235	11317	12595	097459999999	55	BAT5341619	8	009151305	E,D3	MORTGAGE	32108	74657	Not Verified,n,Credit card refinancing,30.57773072,1,0,8,0,148,83.59974136,20,w,4196824
41949024	34631	15985	4917	125714	55	BAT4694572	10	59015342	F,B1	MORTGAGE	60883	55632999999	Source Verified,n,Green loan,11.16411216,0,0,11,0,19097,28.50511393,33,w,6344.62
3922183	31157	9855	31585	98335	59	BAT5849876	9	014251987999999	C,B4	OWN	83900	04312999999	Not Verified,n,Debt Consolidation,29.23281105,0,0,10,0,542,46.39051225,12,f,3157
1313088	27859	33502	16545	20307	36	BAT2833642	14	84819001	B,A4	OWN	60042	5995	Source Verified,n,Credit card refinancing,17.71530168,0,0,9,0,5904,44.38899747,13,f,150.37666
48621970	25721	25880	12035	60189	55	BAT2033411	15	90466542	A,B2	MORTGAGE	114223	9547	Verified,n,Credit card refinancing,24.82525634,0,1,8,1,1232,56.01538182,19,f,1340.1122
44295964	16442	9405	10727	56492	55	BAT5525466	11	18268547	B,A4	MORTGAGE	67166	66525	Not Verified,n,Credit card refinancing,19.15539315,0,0,14,0,1889,46.45458905,10,w,1363
30700045	6373	12341	10452	66161	59	BAT5714674	7	031038523999999	E,D2	OWN	33028	2143	Not Verified,n,Debt Consolidation,8.080329016,0,0,12,0,6475,56.87128113,13,f,1494.958
9030338	14058	24523	9422	342364	59	BAT2003845	14	52337296	A,C1	MORTGAGE	48423	25599	Verified,n,Credit card refinancing,32.48889253,6,0,26,0,18181,46.83020379,16,f,1752.799

File: C:\Usershp\Downloads\train.csv

Interest Rate	Grade	Sub Grade	Employment Duration	Home Ownership	Verification Status	Payment Plan	Loan Title	Debit to Income	Delinquency - two years	Inquires - six months	Open
346	4267	Not Verified,n,Debt Consolidation,16.28475781,1,0,13,0,24246,74.93255103,7,w,2929.646315,0.1205519800000001,2.495290942,0.79372376,0,INDIVIDUAL,49,0,31,311301,6									
1.49105	Source Verified,n,Debt Consolidation,15.41240945,0,0,12,0,812,78.29718401,13,f,772.7693549999999,0.03618117,2.377214756,0.974821103,0,INDIVIDUAL,109,0,53,182610,20885,0										
14.5759	Source Verified,n,Debt Consolidation,28.13761842,0,0,14,0,1843,2.073039673,20,w,863.3243956,18.77846007,4.3162773439999999,1.020074954,0,INDIVIDUAL,66,0,34,39801,2										
AGE,44234.82545	Source Verified,n,Credit card refinancing,17.20988404,1,3,13,1,1544,65.25076114,22,w,129.2395533,19.30664639,1294.818751,0.34895291799999994,0,INDIVIDUAL,39,0,40,9189,602										
95957.47541000001	Not Verified,n,Credit card refinancing,7.914332742000001,3,2,16,0,2277,51.56447406,20,w,464.818124,0.08858434800000001,0.043575438,0.581687716,0,INDIVIDUAL,102391.8243										
102391.8243	Verified,n,Home improvement,15.08391059,0,0,11,0,14501,46.80880397,37,w,525.7381086,0.08382828400000001,3.167937127,0.553076428,0,INDIVIDUAL,71,0,3388,42049,3										
1014	Not Verified,n,Debt Consolidation,29.82971503,0,0,14,0,13067,23.93462444,33,w,1350.2452119999999,0.044964815,0.098448003,0.047589103,0,INDIVIDUAL,87,0,48,184909,43303										
5.0907200001	Verified,n,Credit card refinancing,26.24471043,0,0,6,0,549,15.94738559,17,w,4140.198978,0.01710584,0.53021355799999999,0.216985337,0,INDIVIDUAL,144,0,26,48126										
8.06000004	Source Verified,n,Credit card refinancing,10.04854906,1,0,11,0,1361,35.07334484,30,f,2149.4669429999999,0.008337931,2.912214614,0.856464352,0,INDIVIDUAL,9,0,35										
MT,41858.88652	Verified,n,Debt Consolidation,14.26647349,0,0,17,0,4005,61.27955642,46,f,107.8900429,0.04420826,3.502267497,0.769926275,0,INDIVIDUAL,34,0,26,141492,18025,0										
5.50605	Source Verified,n,Credit Consolidation,32.5409807,1,0,10,2,8266,25.62485547,21,w,1276.258125,0.029185139,7.742979139,1.120056184,0,INDIVIDUAL,55,0,30,33294,11356,0										
AGE,32108.74657	Not Verified,n,Credit card refinancing,30.57773072,1,0,8,0,148,83.59974136,20,w,415.0105247000001,0.032341435,7.90234658,0.9157525990000001,0,INDIVIDUAL,11										
183	35632999999	Source Verified,n,Green loan,11.16411216,0,0,11,0,19097,28.50511393,33,w,6344.628302,0.073004435,4.483083989999999,0.359227761,0,INDIVIDUAL,72,0,1572,30417									
20.04312999999	Not Verified,n,Debt Consolidation,29.23281105,0,0,10,0,542,46.39051225,12,f,3157.741904,0.024345594,6.603832364,1.605330334,0,INDIVIDUAL,83,0,58,207613,1282										
55	Source Verified,n,Credit card refinancing,17.71530168,0,0,9,0,5904,44.38899747,13,f,150.37666299999999,0.039649602,3.345746622,0.627521787,0,INDIVIDUAL,26,0,26,212092										
1223.9547	Verified,n,Credit card refinancing,24.82525634,0,1,8,1,1232,56.01538182,19,f,1340.112263,0.041917396,5.86918474,0.301490787,0,INDIVIDUAL,128,0,23,42122,4280,0										
6.46525	Not Verified,n,Credit card refinancing,19.15539315,0,0,14,0,1889,46.45458905,10,w,1363.05653,0.055918373,5.61049765,0.194200373,0,INDIVIDUAL,66,0,16,128988,15802										
3026.2143	Not Verified,n,Debt Consolidation,8.080329016,0,0,12,0,6475,56.87128113,13,f,1494.969294,0.026837542000000002,1.75119024800000002,1.14841005,0,INDIVIDUAL,109,0,7										
13.25599	Verified,n,Credit card refinancing,32.48889253,6,0,26,0,18181,46.83020379,16,f,1752.7990120000002,32.28935624,2.12352138,1.0540051000000000,0,INDIVIDUAL,74,0,47,1										