



ANALYSIS ON LOAN APPLICANT

Data Exploration Project Report

ABSTRACT:

Applying loan for immediate/essential financial needs is almost done by everyone. The only factor that varies is the reason and the amount required.

There are many factors that are involved in determining if the applicant is going to repay the loan or not.

Through this analysis process, I have tried to find if there are any correlations which have an impact on the loan defaulters.

AKSHAYA KUMAR

CHANDRASEKARAN

ID: 31021301

Contents:

1.INTRODUCTION	3
1.1 PROBLEM SUMMARY:	3
1.2 MOTIVATION:	3
1.3 QUESTIONS:	3
2.DATA WRANGLING:	3
2.1 DATASET REFERENCE:	3
2.2 DATA WRANGLING AND EXPLORATION:	4
2.3 DATA WRANGLING:	5
3. DATA CHECKING AND HANDLING ANOMALIES:	5
3.1 MISSING VALUES:	5
3.2 INCONSISTENCY:	5
3.3 CORRELATION:	6
3.4 DATA TYPE ERROR:	6
4.DATA EXPLORATION:	6
4.1 ANNUAL INCOME:	6
4.2 APPLICATION TYPE	7
4.3 HOME OWNERSHIP:	8
4.4 EMPLOYMENT LENGTH:	8
4.5 OTHER INSIGHTS AND ANALYSIS PERFORMED WITH THE DATASET:	8
4.5.1 COUNTRY:	8
4.5.2 ATTRIBUTE "ANNUAL INCOME" VS "LOAN AMOUNT"	9
4.5.3 PURPOSE	10
4.5.4 EMPLOYEE_TITLE:	11
4.5.5 STATISTICAL TESTS:	11
5.CONCLUSION:	12
6.REFLECTIONS:	12
7. BIBLIOGRAPHY	13

FACTORS INFLUENCING LOAN REPAYMENT.

1.INTRODUCTION:

Almost all of us apply for loan for some of the essential needs and requirements in an intention to repay to the loan within the stipulated date and time. However, few of them are unfortunately not able to repay the loan as planned and end up being a defaulter.

So what loan providers generally do is collect required basic information about the person to assess the application and then decide to proceed with the application or not.

1.1 PROBLEM SUMMARY:

The more the number of defaulters, the more the loss is for the loan providers. Hence before processing and granting any type loan for the applicants, the loan providers do some ground works and analyse if they can proceed. Analysis and explorations are done to identify any such factors which will influence the granting on loan.

1.2 MOTIVATION:

Primary focus of this analysis is on some important factors like Annual Income, Employment length, House ownership, type of application, grade of the applicant, loan amount requested and other few factors.

1.3 QUESTIONS:

1. Does the **type of application** have impact on loan repayment?
2. Does the applicant's **annual income** have any influence towards loan repayment?
3. How does **employment length, home ownership** have influence on loan repayment individually?

2.DATA WRANGLING:

As part of data wrangling process, columns that were insignificant for the analysis were removed and the analysis were done after data cleansing was done.

2.1 DATASET REFERENCE:

The dataset used for the analysis is taken from the site Kaggle. One large data set is used for analysis.

1. **Dataset** → Tabular data consisting of 855696 records and 74 variables. (All kinds of variables such as **location, continuous, categorical** and **simple text with punctuations and numbers** are present in the dataset)

Source Link : <https://www.kaggle.com/sonujha090/xyzcorp-lendingdata>

This data source will help me to answer the above mentioned 3 questions.

Dataset Description: Attributes mentioned below are **used in analysis** from the dataset. Other columns were removed as part of the data wrangling process as their contribution in the analysis were found insignificant.

Attribute Name	Description
annual_inc	The self-reported annual income provided by the borrower during registration.
Application_type	Indicates whether the loan is an individual application or a joint application with two co-borrowers
loan_amnt	The listed amount of the loan applied for by the borrower.
Emp_length	Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years.
Home_ownership	The home ownership status provided by the borrower during registration. Our values are: RENT, OWN, MORTGAGE, OTHER.
Sub_grade	XYZ assigned loan subgrade
verification_status	Was the income source verified
Purpose	A category provided by the borrower for the loan request.
Title	The loan title provided by the borrower.
Emp_title	The job title supplied by the Borrower when applying for the loan. (Text Data)
addr_state	The state provided by the borrower in the loan application (Geo Location)
default_ind	If he is a defaulter or not. 0 is Non defaulter and 1 is defaulter.

Table 2.1 Major Attributes used for analysis

2.2 DATA WRANGLING AND EXPLORATION:

To start with, few packages were imported. Additional packages were added when new functions were introduced. Setting up working directory was one of the important tasks that was carried out that ensured the code and the environment variables were saved and available in the IDE instance.

Data auditing was carried out as a first step to have some understanding on the data using basic commands like *describe()*, *summary()* function in R and observed data anomalies like missing values, inconsistencies, correlation and data type error.

Below is **the sample summary output** and the highlighted columns consisting of the above-mentioned anomalies.

```
> summary(dataSet)
```

id	member_id	loan_amnt	funded_amnt	funded_amnt_inv	term	int_rate	installment
Min. : 54734	Min. : 70699	Min. : 500	Min. : 500	Min. : 0	36 months:600221	Min. : 5.32	Min. : 15.69
1st Qu.: 9067986	1st Qu.:10792732	1st Qu.: 8000	1st Qu.: 8000	1st Qu.: 8000	60 months:255748	1st Qu.: 9.99	1st Qu.: 260.55
Median :34313546	Median :36975319	Median :13000	Median :13000	Median :13000		Median :12.99	Median : 382.55
Mean :32240726	Mean :34762690	Mean :14746	Mean :14732	Mean :14700		Mean :13.19	Mean : 436.24
3rd Qu.:54463114	3rd Qu.:58035586	3rd Qu.:20000	3rd Qu.:20000	3rd Qu.:20000		3rd Qu.:15.99	3rd Qu.: 571.56
Max. :68616867	Max. :73519693	Max. :35000	Max. :35000	Max. :35000		Max. :28.99	Max. :1445.46

grade	sub_grade	emp_title	emp_length	home_ownership	annual_inc	verification_status
A:145665	B3 : 54958	: 49439	10+ years:282090	ANY : 3	Min. : 0	Not Verified :257742
B:247998	B4 : 54116	Teacher : 12965	2 years : 75986	MORTGAGE:429106	1st Qu.: 45000	Source Verified:318178
C:236855	C1 : 51588	Manager : 10821	< 1 year : 67597	NONE : 45	Median : 65000	Verified :280049
D:132802	C2 : 50457	Registered Nurse: 5341	3 years : 67392	OTHER : 144	Mean : 75071	
E: 66448	C3 : 48337	RN : 5182	1 year : 54855	OWN : 84136	3rd Qu.: 90000	
F: 21328	B2 : 47589	Owner : 5157	5 years : 53812	RENT :342535	Max. :9500000	
G: 4873	(Other):548924	(Other) :767064	(Other) :254237			

issue_d	pymnt_plan	desc	purpose
Oct-15 : 48212	n:855964	:734156	debt_consolidation:505392
Jul-15 : 44906	y: 5	: 231	credit_card :200144
Oct-14 : 37442		Borrower added on 03/10/14 > Debt consolidation : 10	home_improvement : 49956
Nov-15 : 37211		Borrower added on 03/17/14 > Debt consolidation : 10	other : 40949
Dec-15 : 35638		Debt Consolidation : 10	major_purchase : 16587
Aug-15 : 35267		Borrower added on 02/19/14 > Debt consolidation : 9	small_business : 9785
(Other):617293	(Other)	:121543	(Other) : 33156

title	zip_code	addr_state	dti	delinq_2yrs	earliest_cr_line	inq_last_6mths
Debt consolidation :398089	945xx : 9466	CA :125172	Min. : 0.00	Min. : 0.0000	Aug-01 : 6433	Min. :0.0000
Credit card refinancing:159228	750xx : 9111	NY : 71114	1st Qu.: 11.88	1st Qu.: 0.0000	Aug-00 : 6322	1st Qu.:0.0000
Home improvement : 38633	112xx : 8894	TX : 68708	Median : 17.61	Median : 0.0000	Oct-00 : 6117	Median :0.0000
Other : 30522	606xx : 8370	FL : 58639	Mean : 18.12	Mean : 0.3116	Oct-01 : 5924	Mean :0.6809
Debt Consolidation : 15469	300xx : 7820	IL : 34379	3rd Qu.: 23.90	3rd Qu.: 0.0000	Aug-02 : 5858	3rd Qu.:1.0000
Major purchase : 11519	100xx : 7348	NJ : 32061	Max. :9999.00	Max. :39.0000	Sep-00 : 5712	Max. :8.0000
(Other) :202509	(Other):804960	(Other):465896			(Other):819603	

mths_since_last_major_derog	policy_code	application_type	annual_inc_joint	dti_joint	verification_status_joint	acc_now_delinq
Min. : 0.0	Min. :1	INDIVIDUAL:855527	Min. : 17950	Min. : 3.0	:855527	Min. : 0.000000
1st Qu.: 27.0	1st Qu.:1	JOINT : 442	1st Qu.: 75000	1st Qu.:13.2	Not Verified : 252	1st Qu.: 0.000000
Median : 44.0	Median :1		Median :100000	Median :17.7	Source Verified: 51	Median : 0.000000
Mean : 44.1	Mean :1		Mean :107412	Mean :18.3	Verified : 139	Mean : 0.004944
3rd Qu.: 61.0	3rd Qu.:1		3rd Qu.:130750	3rd Qu.:22.6		3rd Qu.: 0.000000
Max. :188.0	Max. :1		Max. :410000	Max. :43.9		Max. :14.000000
NA's :642830			NA's :855527	NA's :855529		

tot_coll_amt	tot_cur_bal	open_acc_6m	open_il_6m	open_il_12m	open_il_24m	mths_since_rcnt_il	total_bal_il
Min. : 0	Min. : 0	Min. : 0.0	Min. : 0.0	Min. : 0.0	Min. : 0.0	Min. : 0.0	Min. : 0
1st Qu.: 0	1st Qu.: 29870	1st Qu.: 0.0	1st Qu.: 1.0	1st Qu.: 0.0	1st Qu.: 0.0	1st Qu.: 6.0	1st Qu.: 10390
Median : 0	Median : 81008	Median : 1.0	Median : 2.0	Median : 0.0	Median : 1.0	Median : 12.0	Median : 24960
Mean : 225	Mean : 139766	Mean : 1.1	Mean : 2.9	Mean : 0.7	Mean : 1.7	Mean : 20.8	Mean : 36512
3rd Qu.: 0	3rd Qu.: 208703	3rd Qu.: 2.0	3rd Qu.: 4.0	3rd Qu.: 1.0	3rd Qu.: 2.0	3rd Qu.: 23.0	3rd Qu.: 47493
Max. :9152545	Max. :8000078	Max. :12.0	Max. :40.0	Max. :12.0	Max. :15.0	Max. :300.0	Max. :634217
NA's :67313	NA's :67313	NA's :842681	NA's :842681	NA's :842681	NA's :842681	NA's :843035	NA's :842681

inq_last_12m	default_ind
Min. : -4.0	Min. :0.00000
1st Qu.: 0.0	1st Qu.:0.00000
Median : 2.0	Median :0.00000
Mean : 1.8	Mean :0.05429
3rd Qu.: 3.0	3rd Qu.:0.00000
Max. :32.0	Max. :1.00000
NA's :842681	

2.3 DATA WRANGLING:

From the above output, the following key features were found.

1. Every grade has 5 different sub-grades, so if we know the sub-grade, the applicant's grade can be found.
2. The attributes "Title" and "Purpose" were almost equal. They depict the same meaning in different kinds.
3. There were few attributes with missing values more than 95% in them.

Attributes that exhibited the above-mentioned properties were identified and removed as part of wrangling process for exploration.

3. DATA CHECKING AND HANDLING ANOMALIES:

3.1 MISSING VALUES:

Few attributes were found to have more number of missing values as shown in the below figure 3.1. While analysing I found the **reason** for very high number of missing values are, a **particular application type "joint"** were found to **have all records**. When the application type is **"JOINT"** those columns had values and, if the **application type is individual** the attribute's **value will be NA**. Those columns were dropped using *dplyr* and *plyr* libraries in R as their contribution were insignificant for the analysis. Other attributes such as "emp_length" that had few missing values were handled and the missing values were imputed with the method "measure of central tendency" (**median** for **continuous** data type and **mode** for **categorical** data type)

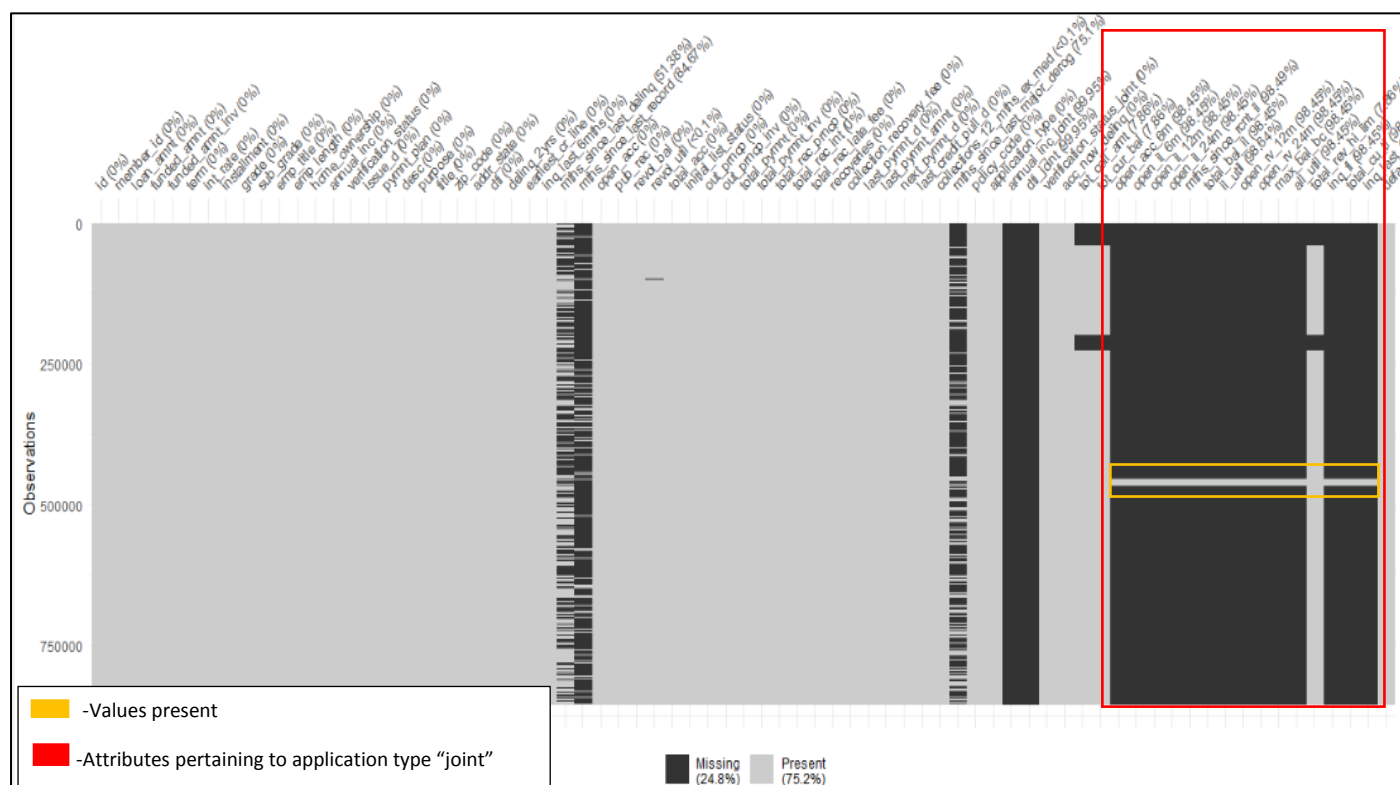


Figure 3.1 Image showing Missing values

3.2 INCONSISTENCY:

Attributes namely "home_ownership", "verification_status" and "emp_title" had inconsistency in them. For eg: attribute "verification_status" had **Source Verified, verified and Not verified** unique values as shown in the figure 3.2. Here, both **source verified and verified means the same**. Home ownership has **Any** and **None** as extra values other than mentioned in business requirement (in table 2.1). Anomalies like these were handled as well as shown in the figure 3.3.

```
> unique(df$home_ownership)
[1] RENT      OWN      MORTGAGE OTHER    NONE     ANY
Levels: ANY MORTGAGE NONE OTHER OWN RENT
> unique(df$verification_status)
[1] Verified      Source Verified Not Verified
Levels: Not Verified Source Verified verified
> |
```

Figure 3.2 Attributes having inconsistency

```
> unique(df$home_ownership)
[1] RENT      OWN      MORTGAGE OTHER
Levels: MORTGAGE OTHER OWN RENT
> unique(df$verification_status)
[1] Source Verified Not Verified
Levels: Not Verified Source Verified
> |
```

Figure 3.3 Inconsistency eliminated

3.3 CORRELATION:

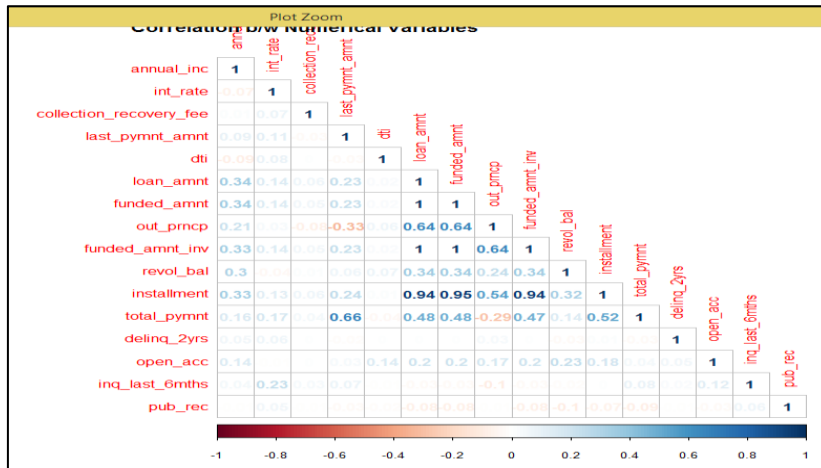


Figure 3.4 Correlation plot of numerical variables

Few attributes were exhibiting correlation. Correlation plot was drawn using `cor()` function to find the correlation between the attributes as shown in the figure 3.4. Attributes having high correlation value greater than 0.8 was removed using `dplyr` and `plyr` libraries in R.

There was one other interesting property found during the analysis, *Total_payment* attribute was having a correlation value of 1 with the sum of (*Principal + Interest + Total Recoveries Late Fee + Recoveries*) attributes. This was also handled

3.4 DATA TYPE ERROR:

```
> str(df$default_ind)
int [1:855969] 0 1 0 0 0 0 0 1 1 ...
> df$default_ind <- as.factor(df$default_ind)
> str(df$default_ind)
Factor w/ 2 levels "0","1": 1 2 1 1 1 1 1 2 2 ...
> |
```

Figure 3.5 Imaging showing data type handled

Few attribute's data type was misinterpreted by R. For example, "default_ind" attribute is of **type categorical** with unique values 0 and 1 but R **misinterpreted** it as a **numerical** attribute. Attributes that I considered for the analysis were all checked and reformatted it to correct data type as shown in the figure 3.5.

All the anomalies were handled, and cleaning was done making the data ready for the exploration and analysis.

4.DATA EXPLORATION:

During the process of exploration, analysis was also done parallelly. 4 Major attributes namely,

1. Annual Income
2. Application Type
3. Home Ownership
4. Employment length , are analysed in detail on their influence on the loan payment to answer the question. Other attributes were also analysed additionally.

4.1 ANNUAL INCOME:

It is evident from the figure 4.1 that the **applicants with less than 1000K as their annual income are subject to be defaulters**. On drilling down further, only those whose salary is less than 300k and 1000k was analysed as shown in the figure 4.2(a) and (b) based on their home_ownership. It is understood from figure 4.2(b) that applicants who are salaried less than **750000 and who are either living in a rented house or has mortgaged their house are likely to be defaulters**. Also, people who **own the house but are getting less than 500000** are also likely to be defaulters.

Since the box-plot was not clearly visible and the stats of a box plot could not be interpreted from figure 4.1, subset of the data was taken with annual salary less than 300k. It was found that the annual salary's 25 percentile (Q1), 50 percentile/median(Q2), 75 percentile (Q3) were almost the same for both defaulters and non-defaulters shown in the figure 4.2(a)

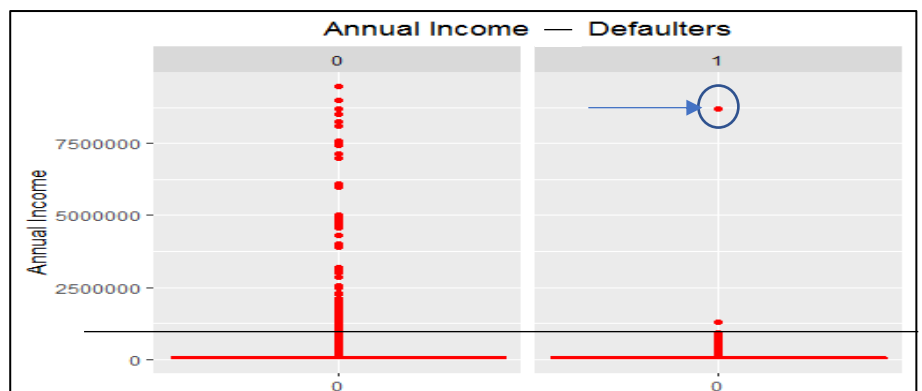


Figure 4.1 Box Plot of annual income(full data set)

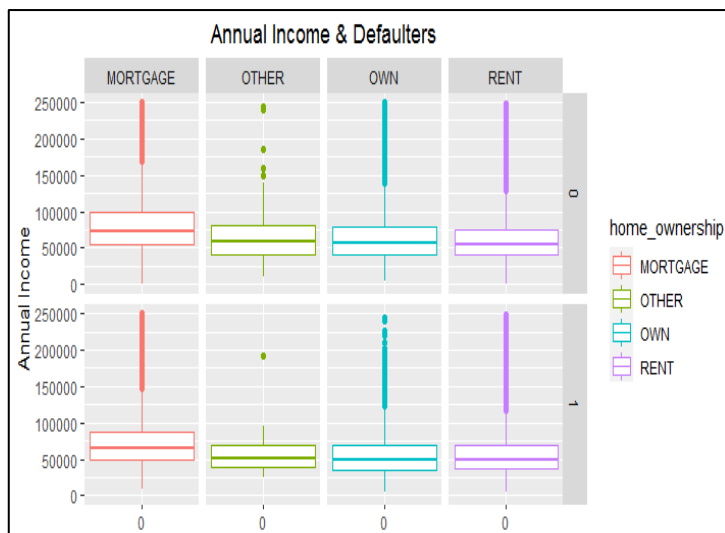


Fig 4.2(a) Box-Plot to understand annual income(<300k)

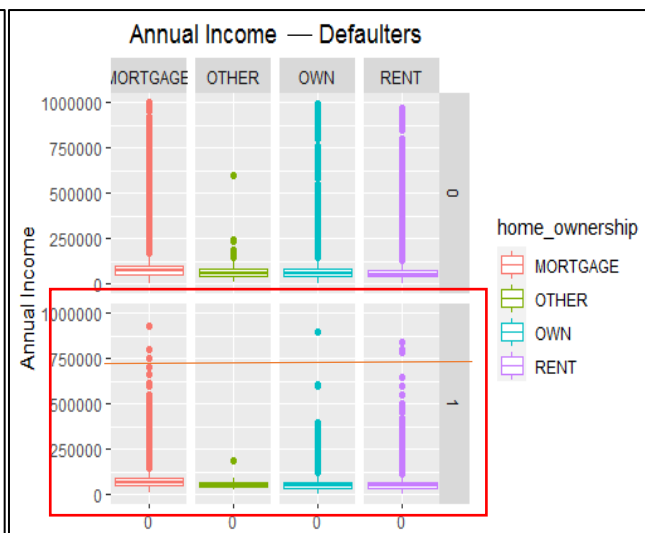


Figure 4.2(b) Box plot of annual income (<1000k)

One interesting factor noticed during this analysis is, one applicant with **annual salary more than 75000000** has been a defaulter. He/She might have started a business and unfortunately didn't turn out to be a profitable one. He/She is an **outlier** in this data.

annual_inc	application_type	loan_amnt	emp_length	home_ownership	sub_grade	verification_status	purpose	title	emp_title	addr_state	default_ind
8706582	INDIVIDUAL	8000	10+ years	MORTGAGE	C3	Source Verified	credit_card	Credit card refinancing	Correctional Sgt.	IL	1

Figure 4.3 Outlier and interesting record

4.2 APPLICATION TYPE

```
> table(finalData$application_type) / sum(table(finalData$application_type))

INDIVIDUAL    JOINT
0.9994836262 0.0005163738
> table(finalData$application_type)

INDIVIDUAL    JOINT
855527       442
>
```

Figure 4.4 Table function in R to show the count of application type

In the dataset, only **.05% of the applicants** have got the loan as type **"JOINT"** while the rest **99.95%** of the applicants have applied the loan **individually** which can be seen from the figure 4.4.

Now, the analysis has been done for **both the types separately**.

It is understood that applicants with type as **"joint"** have always paid the loan as shown in the figure 4.5 while there are chances that the applicant **might become a defaulter** if the type is **"individual"** as shown in figure 4.6. This makes sense as applicants with type **"joint"** both be the earning members and hence they will be able to repay the loan easily when compared to single person earning and closing the loan.

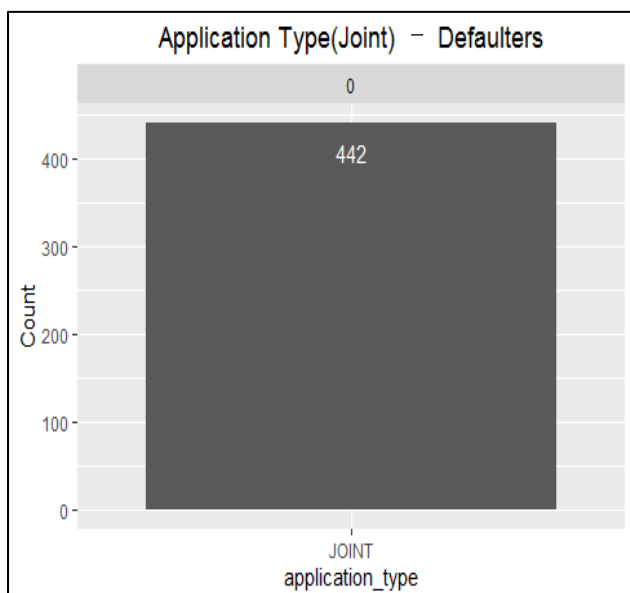


Figure 4.5 Application Type "Joint"

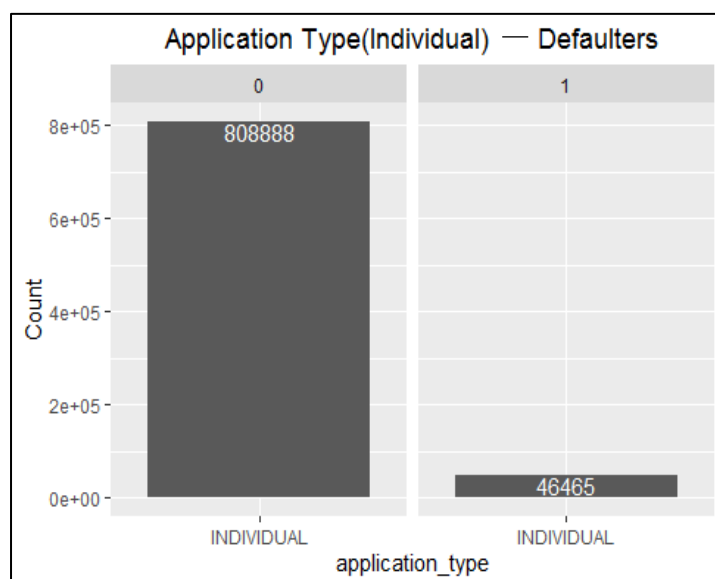


Figure 4.6 Application Type "Individual"

4.3 HOME OWNERSHIP:

From the **analysis outcome based on the annual income in box plot**, it is found that most of the people who are earning less **than 300k as their annual income are found to be defaulters**. Considering the outcome from that analysis, I did a **subset** of data for applicants earning less than **1000k annual income** to analyse in depth on how home ownership has influence on their loan repayment.

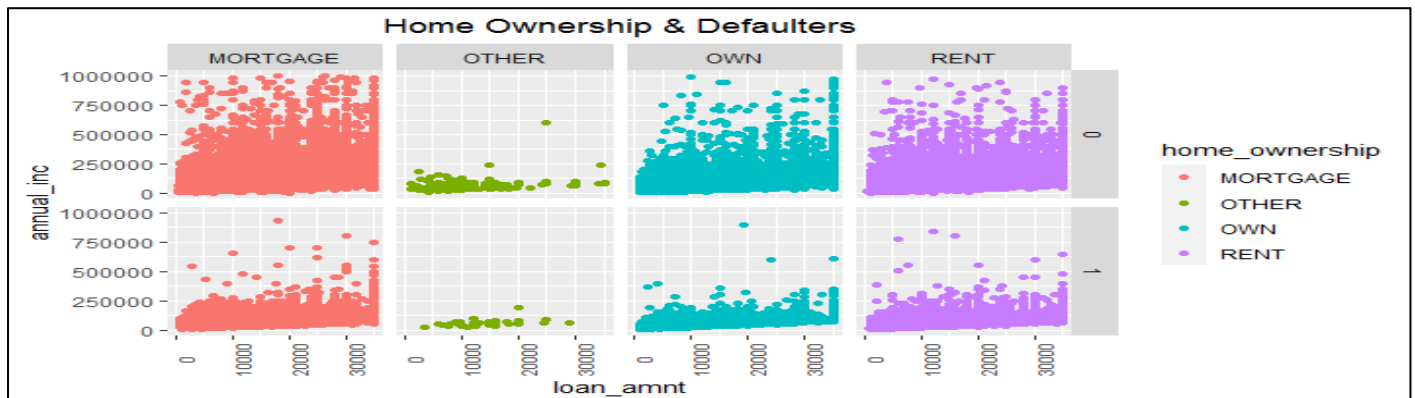


Figure 4.7 House Ownership Vs Defaulters wrt annual income

From the figure 4.7 it is understood that applicants whose **annual income** is less than **250k** has **their house ownership as “mortgaged” or “rented”**.

It can also be seen that applicants who also **owned the house are also subjected to be defaulters**, but **their annual income is less** when compared to applicants whose house is either “mortgaged” or “rented”.

4.4 EMPLOYMENT LENGTH:

From the figure 4.8, I concluded that **most of the loan applicants have 10+ years of experience** and **employment length has no significance to determine if they are going to be a defaulter or not**. The more the number of applicants with 10+ years of experience, the number of defaulters is going to be relatively high when compared to other categories. Thus **employment length has no significance** in determining if the applicant is going to be a defaulter or not.

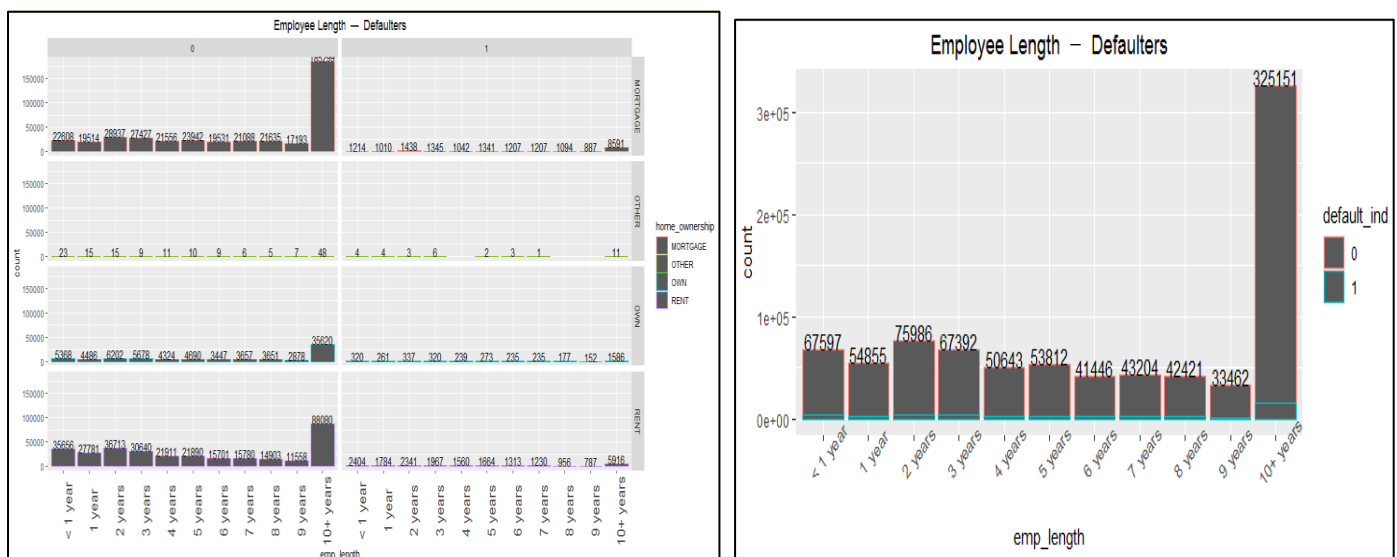


Figure 4.8 Employment length count plotted wrt house ownership and defaulters

4.5 OTHER INSIGHTS AND ANALYSIS PERFORMED WITH THE DATASET:

4.5.1 COUNTRY:

When plotted the map I saw that I **was analysing applicants from The United states**. On Analysing, I saw an **interesting fact** that applicants from the mentioned states **North Dakota and Maine have not defaulted the loan**. Few applicants from other states have been defaulters while others have not. The same was understood from the figure 4.9 below. Also, there are **almost 35% of applicants** from the states **New York, Texas, California and Florida** as the sum of applicants from these states alone were more than approximately 300,000.

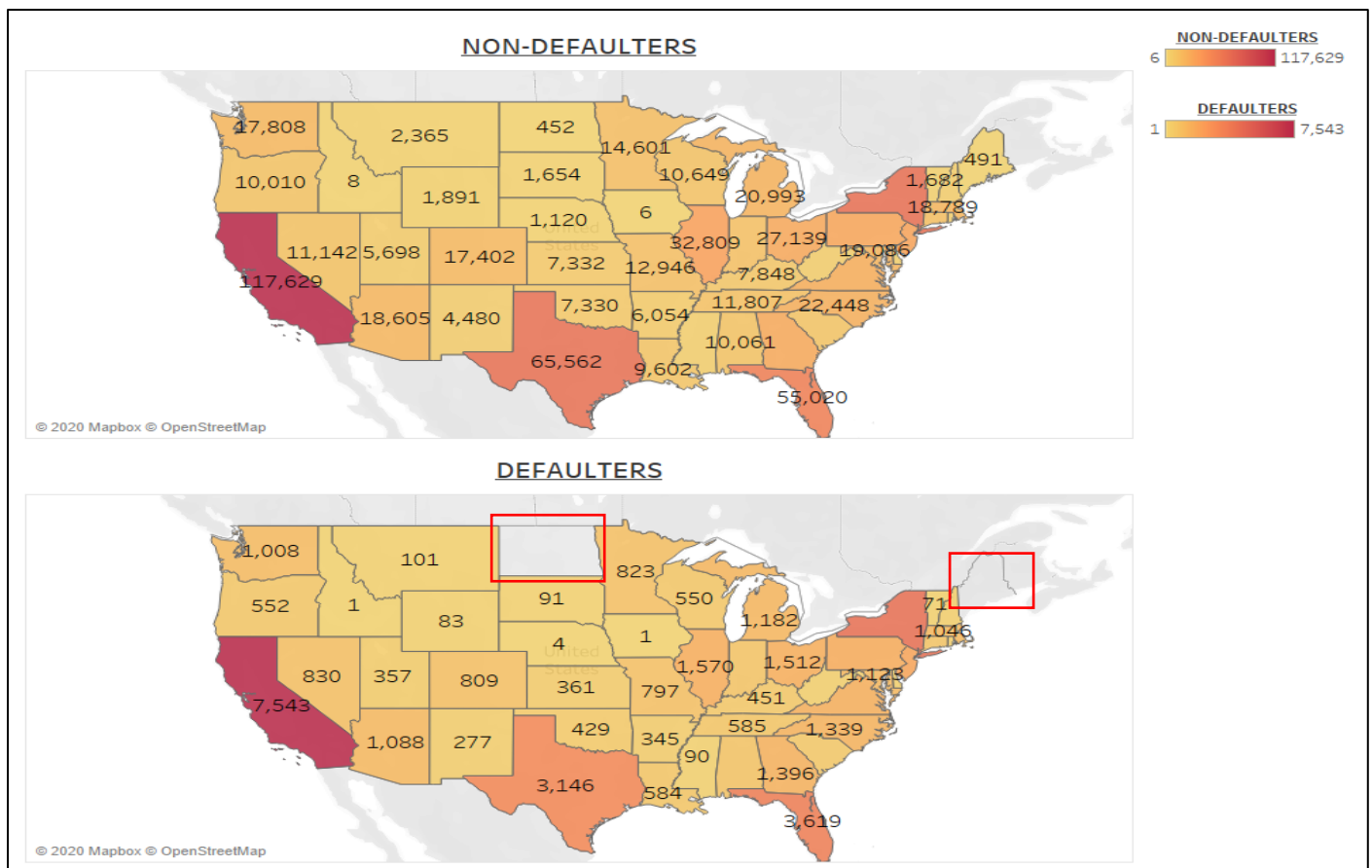


Figure4.9 showing defaulter status state wise

4.5.2 ATTRIBUTE “ANNUAL INCOME” VS “LOAN AMOUNT”

When comparing between the attributes annual income and loan amount to find a relation if any, I was able to conclude that **loan amount sanctioned to each of the applicants were almost same**. But an applicant tends to be a defaulter **based on his/her annual income**. From the below figures 4.10 and 4.11, it can be found that on an **average**, the annual income of a **non-defaulter is close to 1000k** while average remains less than **300k for a defaulter**.

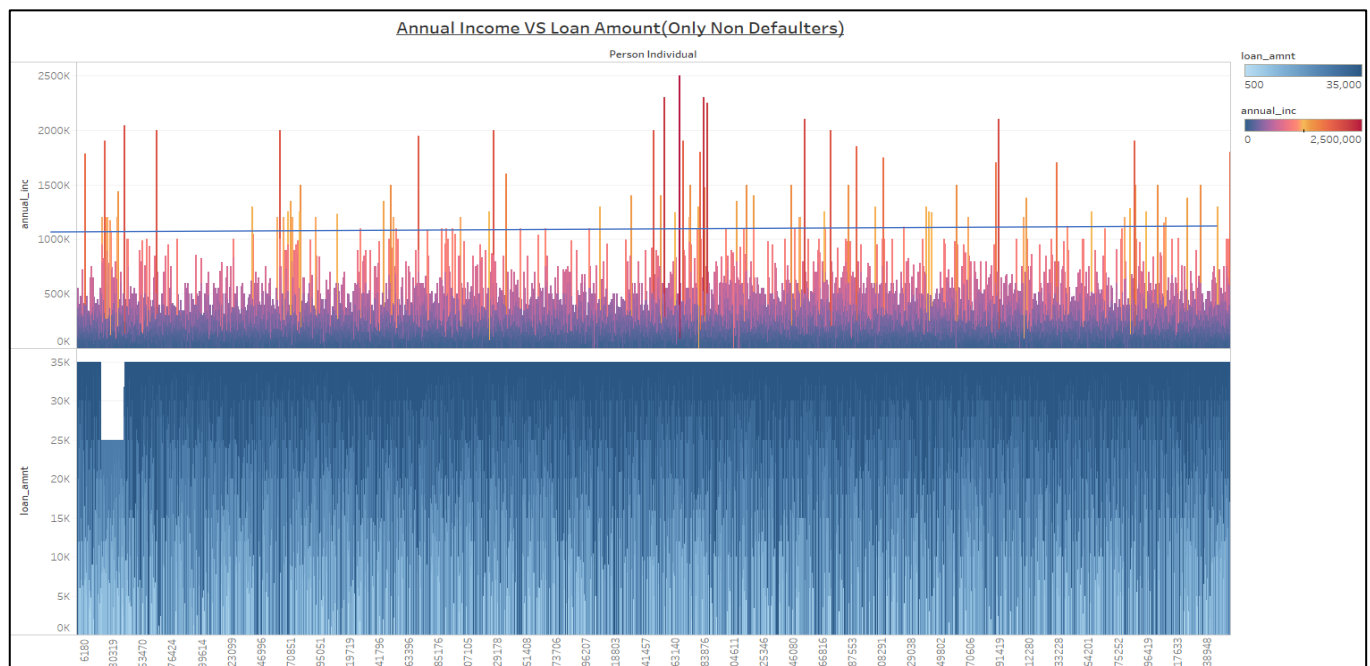


Figure 4.10 Non-Defaulter's annual income vs the loan amount

Figure 4.10 shows the annual income compared to loan amount for non-defaulters and the figure 4.11 shows the annual income compared to loan amount for defaulters.

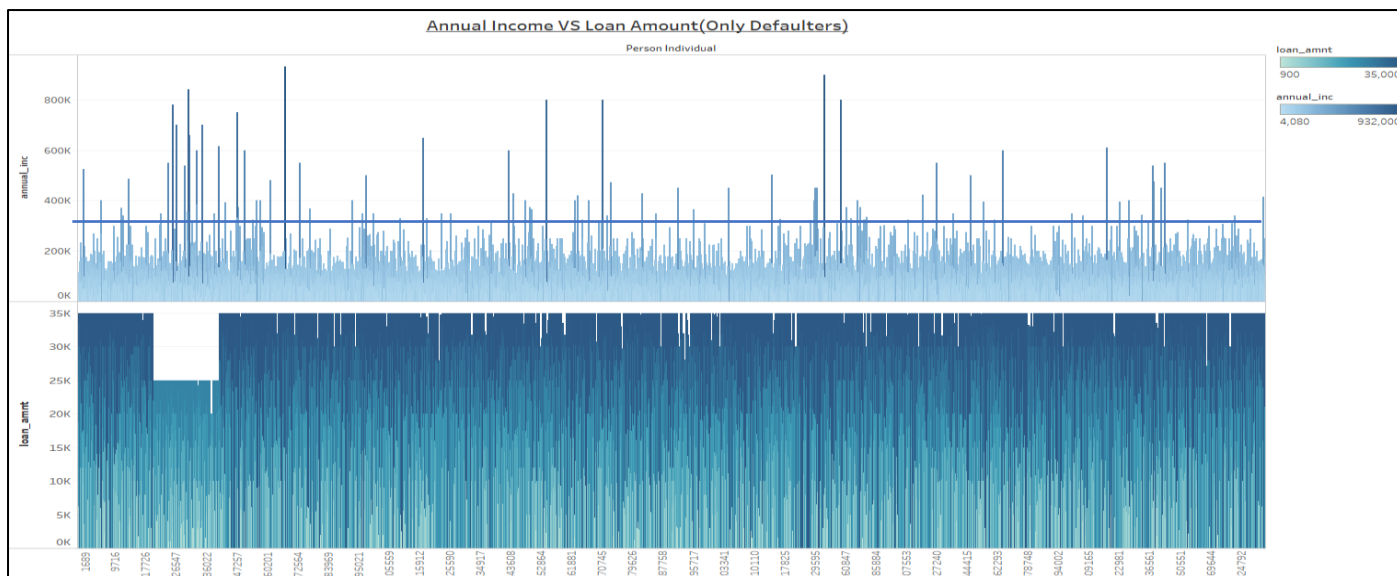


Figure 4.11 Defaulter's annual income vs the loan amount

4.5.3 PURPOSE:

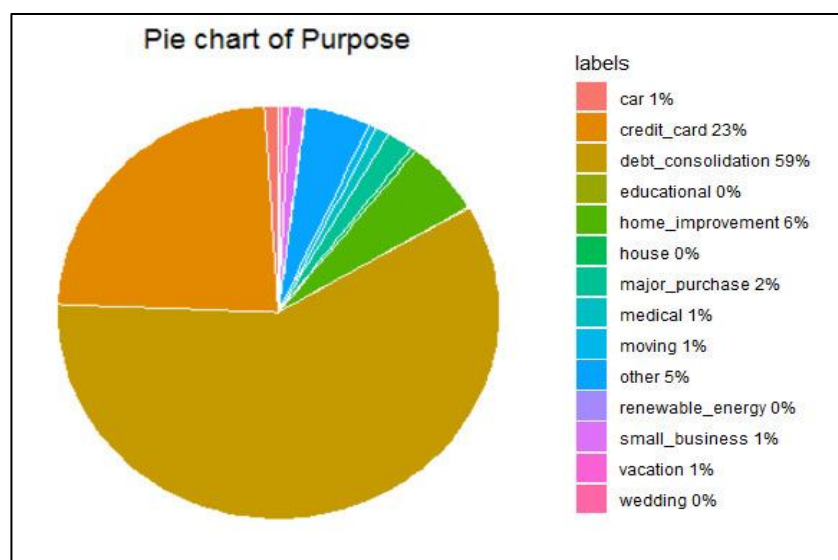


Figure 4.12 Pie showing the purpose of loan (full dataset)

To analyse for what purpose the applicants are applying for a loan, a pie chart got created in R with `geom_bar()` and `coord_polar()` and it is found that **almost 60% of the applicants** have applied loan for the **purpose of debt consolidations** (debt consolidation is process of getting a loan to close of other loans and manage financially).

Next reason for applying loan is to repay their **credit card debts**. Thus, from the figure 4.12, we can conclude that **almost 75%** of the applicants get loan only to **close their previous loans**.

To get some more insights on the purpose, I took the subset of the data excluding both the major purposes (debt consolidation and credit card) and analysed using treemap with `treemap()` package in R. It was found that the next **major purpose to get the loan was for house improvement and other** from the figure 4.13.

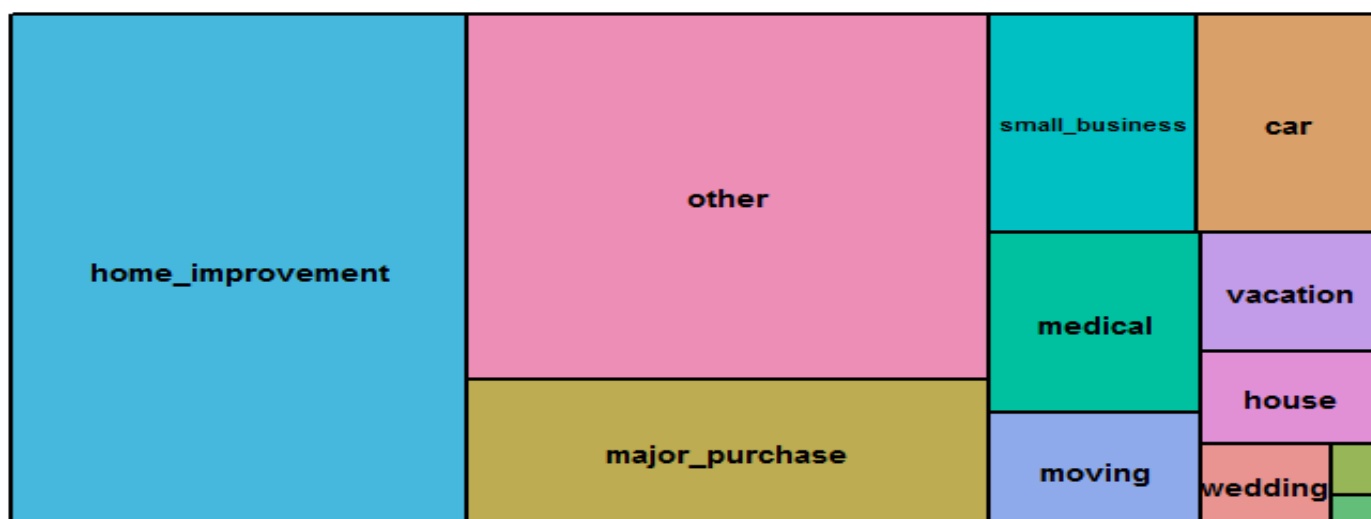


Figure 4.13 Treemap showing the purpose of loan (Subset of dataset excluding debt consolidations and credit card)

The attribute **emp_title** was a simple text that consisted of punctuations, numbers and some special characters using *tm* and *removePunctuation()* libraries and functions in R as shown in the figure 4.14. Later, a word cloud was created to analyse applicant from which profession has applied for the loan. It can be seen from the figure 4.15 that **most of the applicants** are one among **Teacher, Project Managerial roles, Owner, Sales person, Engineer, accountant and supervisors**.

Punctuations, special characters and numbers(for those who worked in more than 1 jobs) were there in the attribute employee title.

All those were removed as part of cleaning and word cloud was created in R to analyse the most employee title who are applying for loan.

Figure 4.14 DataFrame showing punctuations and special characters



Figure 4.15 Word Cloud for the attribute employee title

Statistical tests like chi-squared tests (for two categorical type attributes) and variance tests (for two continuous type attributes) were performed to understand the correlation between the variables as shown in the figure 4.15. It was found that, all the **variables have less correlation value** (as expected) and the correlation between *annual income*, *home ownership*, *application type* with *defaulte_ind* were more when compared to other attributes. Sample statistical output is shown in the figure 4.15.

```

> chisq.test(finalData$home_ownership,finalData$default_ind,correct = T)

Pearson's Chi-squared test

data:  finalData$home_ownership and finalData$default_ind
X-squared = 1115.8, df = 3, p-value < 2.2e-16

> var.test(finalData$annual_inc,finalData$loan_amnt)

F test to compare two variances

data:  finalData$annual_inc and finalData$loan_amnt
F = 58.179, num df = 855968, denom df = 855968, p-value < 2.2e-16
alternative hypothesis: true ratio of variances is not equal to 1
95 percent confidence interval:
 58.00519 58.35380
sample estimates:
ratio of variances
 58.1791

> chisq.test(finalData$application_type,finalData$emp_length)

Pearson's Chi-squared test

data:  finalData$application_type and finalData$emp_length
X-squared = 17.285, df = 10, p-value = 0.06828

```

Figure 4.16 Showing statistical test results done in R to find the correlation categorical variables

5.CONCLUSION:

On analysing various attributes which could affect the applicant to be a defaulter or not, below main attributes has their own significance and influence:

- **Annual Income:** Applicants with annual income less than 250k are most likely to be a defaulter. Also, we cannot access the applicant's application only with their annual income. Other attributes should also be taken into considerations.
- **Application Type:** When the application type is "Joint", then those applicants are not going to default as the earning is twice than that of applicants with application type "individual". When an application is of type individual, several other attributes are required to proceed and process the loan request.
- **Employment length:** The attribute employment length has no greater influence in deciding if the person is going to be a defaulter or not. But one quite interesting fact found is, applicants with more than 10 years of experience are applying for loan. This makes sense because the applicant may have planned their future based on their financial status to make progress in their life.
- **Home Ownership:** Applicant who are residing in a mortgaged house or in a rented house are likely to default their loan provided their annual salary is less than 250k. Hence only this particular attribute has some influence on determining if the applicant is going to default or not and has greater influence when other attributes are analysed alongside.
- **Other key attributes and findings:** Country's state was a new and interesting finding during the analysis. Other attributes were also analysed but per my analysis, they didn't seem to have much impact on determining if a person is going to default a loan or not.

6.REFLECTIONS:

Upon completion of this analysis, I gained knowledge on data wrangling, hypothesis testing, data cleaning, data modelling and data visualization using R and Tableau. Knowledge about various libraries like plyr, dplyr, wordcloud, RcolorBrewer, treemap, ggplot2, randomcoloR, corrplot, leaflet was also gained and also about its appropriate uses.

7. BIBLIOGRAPHY

1. <http://www.sthda.com/english/wiki/text-mining-and-word-cloud-fundamentals-in-r-5-simple-steps-you-should-know> (Word Cloud in R)
2. <https://rstudio.com/wp-content/uploads/2016/11/ggplot2-cheatsheet-2.1.pdf> (About GGLOT 2 graphs)
3. [https://en.wikipedia.org/wiki/Default_\(finance\)](https://en.wikipedia.org/wiki/Default_(finance)) (Domain knowledge)
4. <https://www.pexels.com/search/finance/> (Cover Photo)
5. <https://datascience.stackexchange.com/> (Hypothesis and statistical tests)