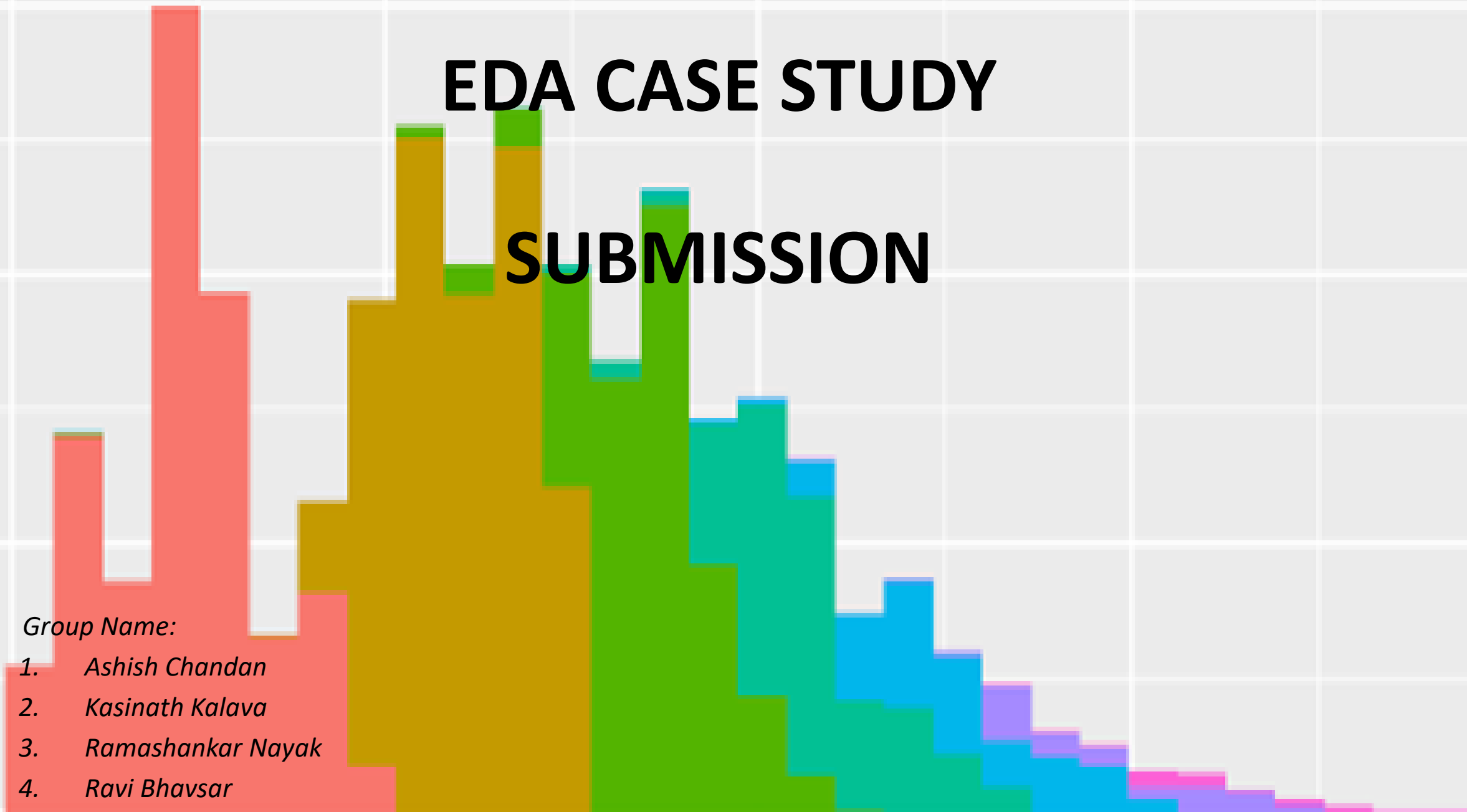


# EDA CASE STUDY

## SUBMISSION

*Group Name:*

1. *Ashish Chandan*
2. *Kasinath Kalava*
3. *Ramashankar Nayak*
4. *Ravi Bhavsar*



**Business Objective:** The objective is to understand the driving factors behind loan default, i.e. variables which are strong indicators of default.

## Methodology:

- Understanding the provided loan dataset with the objective to identify various variables.
- Perform EDA on the loan dataset.
- Determine the driving factors and attributes that are strong indicators for loan defaults.
- Finding trends on the basis of Count and Ratio of Charged Off loans using CrossTable in the below format and corresponding plot:

Each cell contains Count of combination and Ratio of Count with Row Total to determine %age of Charged off/Fully Paid

Count
Ratio

**Count** refers Count of combination

**Ratio** refers Ratio of Count with row total

## Data Exploration:

1. Data set contains 111 columns, however, we considered following for our analysis :

id	loan_amnt	term	int_rate	installment	grade
sub_grade	emp_length	home_ownership	annual_inc	verification_status	issue_d
loan_status	purpose	addr_state	dti	delinq_2yrs	inq_last_6mths
open_acc	pub_rec	revol_bal	revol_util	revol_util	

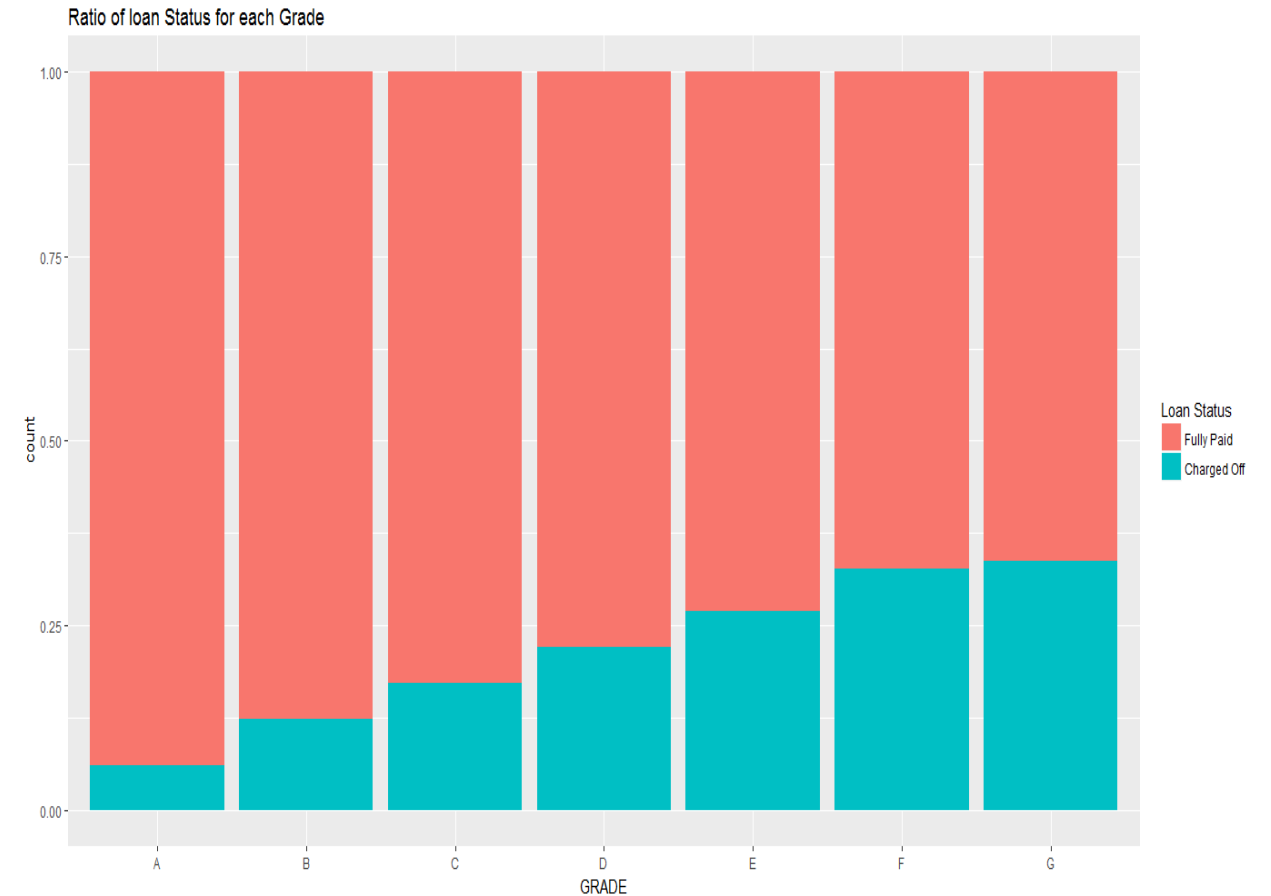
- We have eliminated any fields that would not have been known at the time of issuance as we'll be trying to make decisions on loan investments using available pre-issuance data.
- We have also eliminate a few indicative data fields that are repetitive or too granular to be analyzed.

2. We have to find driving variables to determine defaulters, so just keeping "Charged Off" and "Fully Paid" loan\_status and have removed "Current" as we cannot categorize the record into defaulter/non-defaulter.

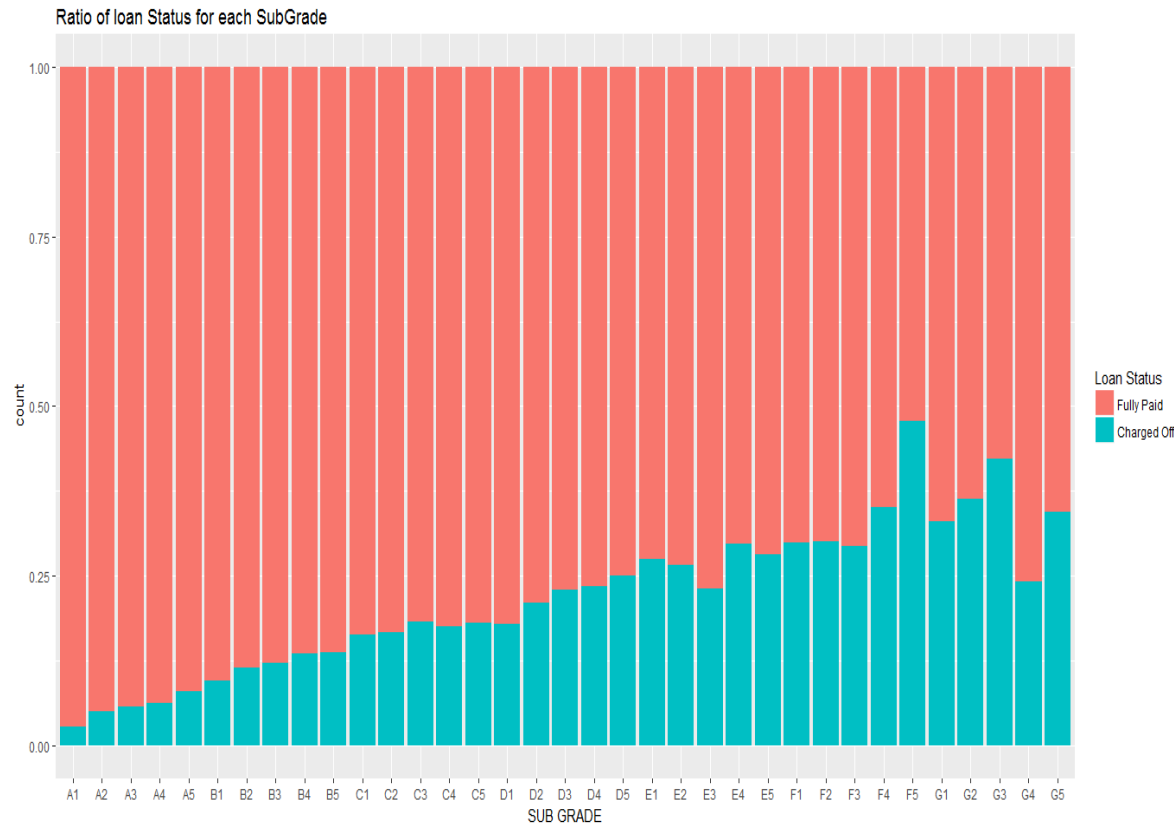
**Data Cleaning and Data Manipulation:** Possible data inconsistencies are as follows:

- Checking for any duplicate ids
- Check for missing values
- Removing extra space
- Extracting year and month from issue date
- Convert a required variables to factors
- Convert int\_rate and revol\_util into numeric

loan\$grade	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
A	602 0.060	9443 0.940	10045 0.260
B	1425 0.122	10250 0.878	11675 0.303
C	1347 0.172	6487 0.828	7834 0.203
D	1118 0.220	3967 0.780	5085 0.132
E	715 0.268	1948 0.732	2663 0.069
F	319 0.327	657 0.673	976 0.025
G	101 0.338	198 0.662	299 0.008
Column Total	5627	32950	38577



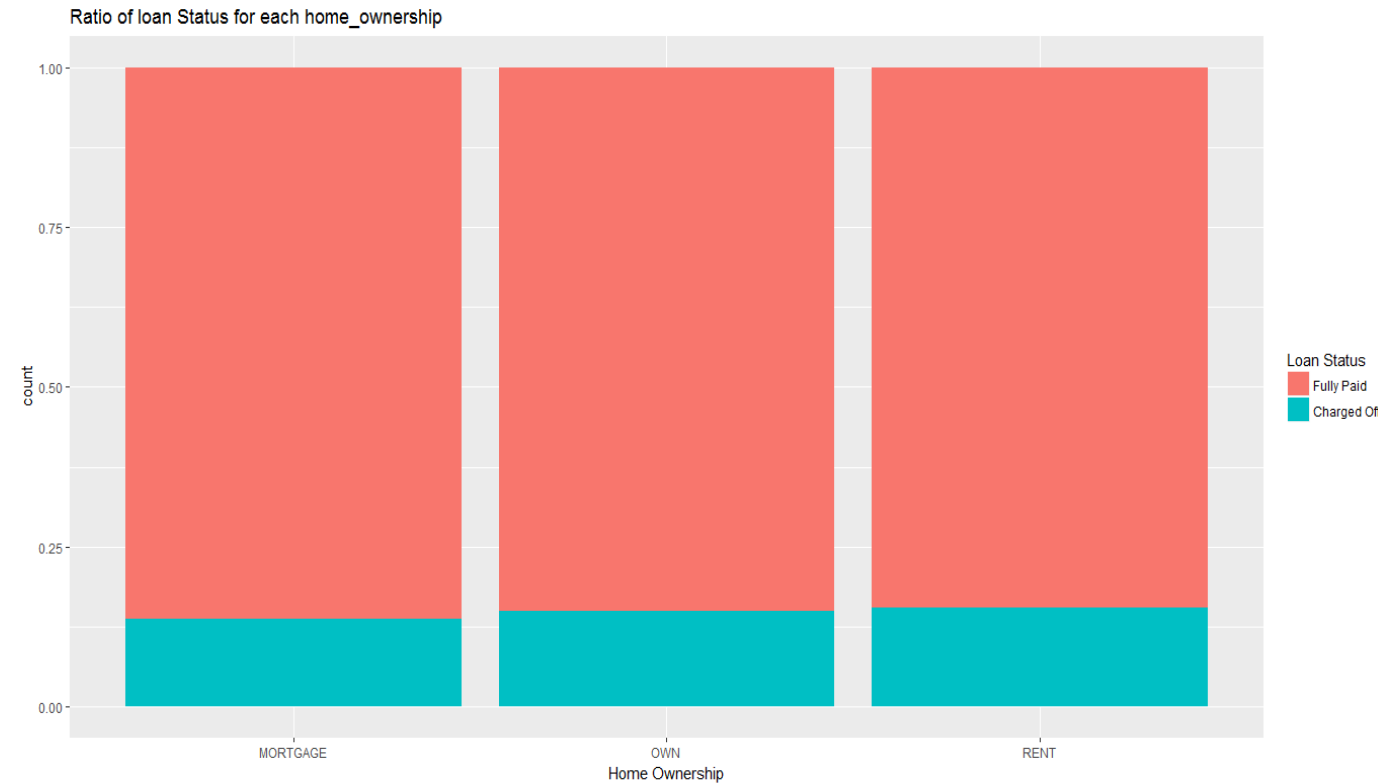
**Default rate steadily increases as the loan grades worsen from A to G.**



Plot shows similar pattern to Grades, although the trend begins to weaken across the G1-G5. In contrast to data points for the A1 to F5, there are fewer data points for G1-G5 and differences are not large enough to be significant. Hence, G1-G5 sub-grades can be ignored.

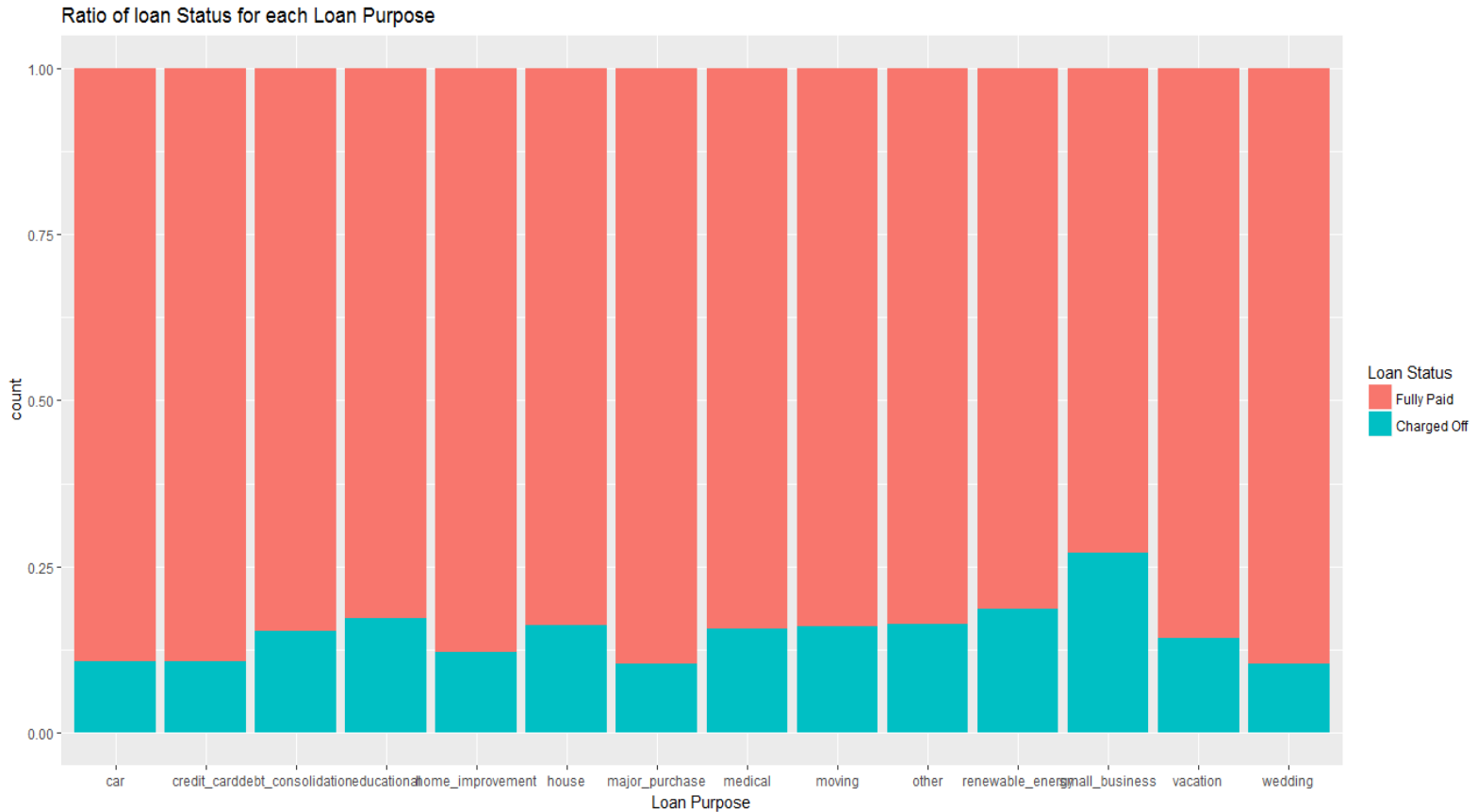
Frequency				Ratio			
	Charged	Off	Fully Paid		Charged	Off	Fully Paid
A1	30		1109	A1	0.0053		0.0337
A2	74		1434	A2	0.0132		0.0435
A3	103		1707	A3	0.0183		0.0518
A4	178		2695	A4	0.0316		0.0818
A5	217		2498	A5	0.0386		0.0758
B1	171		1626	B1	0.0304		0.0493
B2	228		1773	B2	0.0405		0.0538
B3	341		2484	B3	0.0606		0.0754
B4	329		2108	B4	0.0585		0.0640
B5	356		2259	B5	0.0633		0.0686
C1	336		1719	C1	0.0597		0.0522
C2	321		1610	C2	0.0570		0.0489
C3	270		1218	C3	0.0480		0.0370
C4	212		994	C4	0.0377		0.0302
C5	208		946	C5	0.0370		0.0287
D1	167		764	D1	0.0297		0.0232
D2	271		1015	D2	0.0482		0.0308
D3	256		860	D3	0.0455		0.0261
D4	215		703	D4	0.0382		0.0213
D5	209		625	D5	0.0371		0.0190
E1	198		524	E1	0.0352		0.0159
E2	163		451	E2	0.0290		0.0137
E3	119		397	E3	0.0211		0.0120
E4	126		298	E4	0.0224		0.0090
E5	109		278	E5	0.0194		0.0084
F1	91		214	F1	0.0162		0.0065
F2	70		163	F2	0.0124		0.0049
F3	51		123	F3	0.0091		0.0037
F4	53		98	F4	0.0094		0.0030
F5	54		59	F5	0.0096		0.0018
G1	31		63	G1	0.0055		0.0019
G2	28		49	G2	0.0050		0.0015
G3	19		26	G3	0.0034		0.0008
G4	13		41	G4	0.0023		0.0012
G5	10		19	G5	0.0018		0.0006

loan1\$home_ownership	loan1\$loan_status		Row Total
	Charged off	Fully Paid	
MORTGAGE	2327 0.137	14694 0.863	17021 0.442
OWN	443 0.149	2532 0.851	2975 0.077
RENT	2839 0.154	15641 0.846	18480 0.480
Column Total	5609	32867	38476



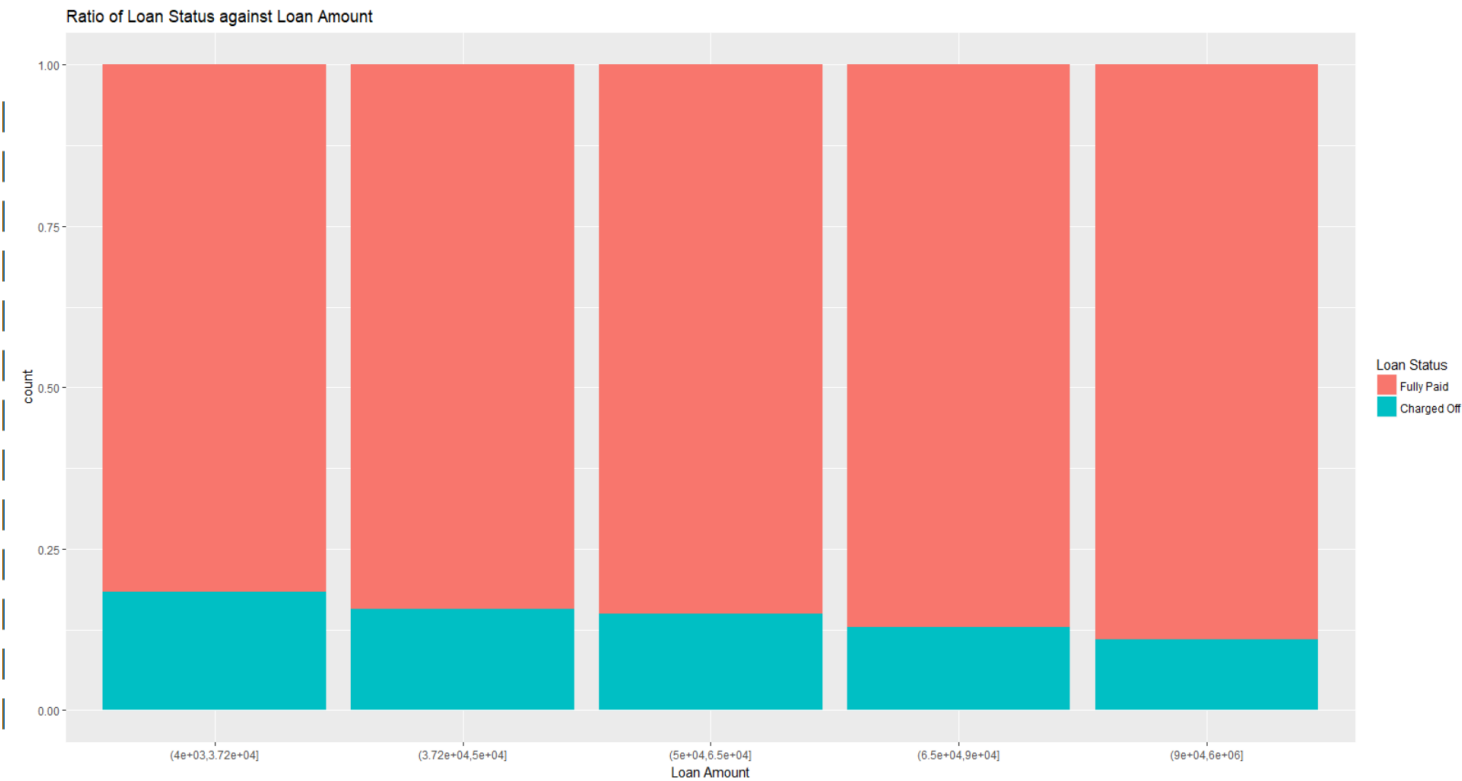
**Those with Mortgages default (Charged off) the least, followed by those who own their own homes.**

loan\$purpose	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
car	160 0.107	1339 0.893	1499 0.039
credit_card	542 0.108	4485 0.892	5027 0.130
debt_consolidation	2767 0.153	15288 0.847	18055 0.468
educational	56 0.172	269 0.828	325 0.008
home_improvement	347 0.121	2528 0.879	2875 0.075
house	59 0.161	308 0.839	367 0.010
major_purchase	222 0.103	1928 0.897	2150 0.056
medical	106 0.156	575 0.844	681 0.018
moving	92 0.160	484 0.840	576 0.015
other	633 0.164	3232 0.836	3865 0.100
renewable_energy	19 0.186	83 0.814	102 0.003
small_business	475 0.271	1279 0.729	1754 0.045
vacation	53 0.141	322 0.859	375 0.010
wedding	96 0.104	830 0.896	926 0.024
Column Total	5627	32950	38577



**Loan purpose refers to the borrower's stated reason for taking out the loan.**  
**Here, loan status for Small business is very poor.**

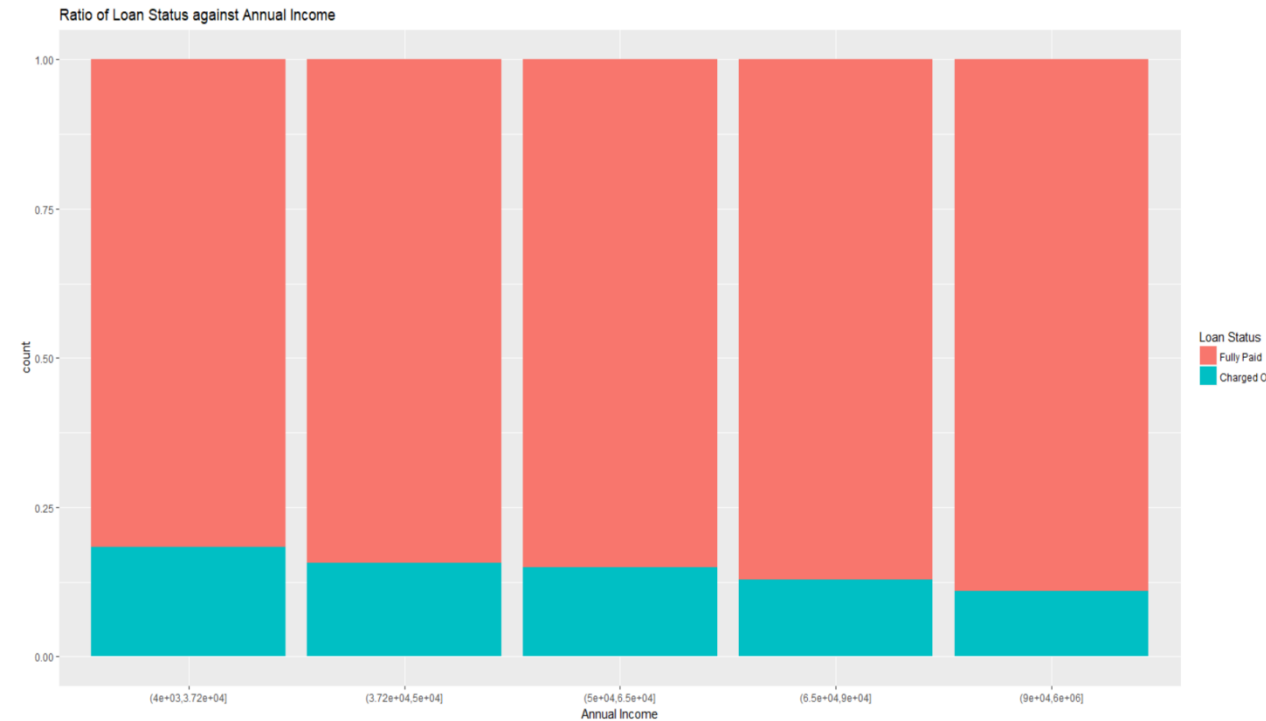
loan\$loan_status			
loan\$loan_amnt_bkt	Charged Off	Fully Paid	Row Total
(0,15000]	4011	25946	29957
	0.134	0.866	0.777
(15000,30000]	1436	6449	7885
	0.182	0.818	0.204
(30000,35000]	180	555	735
	0.245	0.755	0.019
Column Total	5627	32950	38577



As the amount borrowed increases, we see increasing rates of defaulting loans

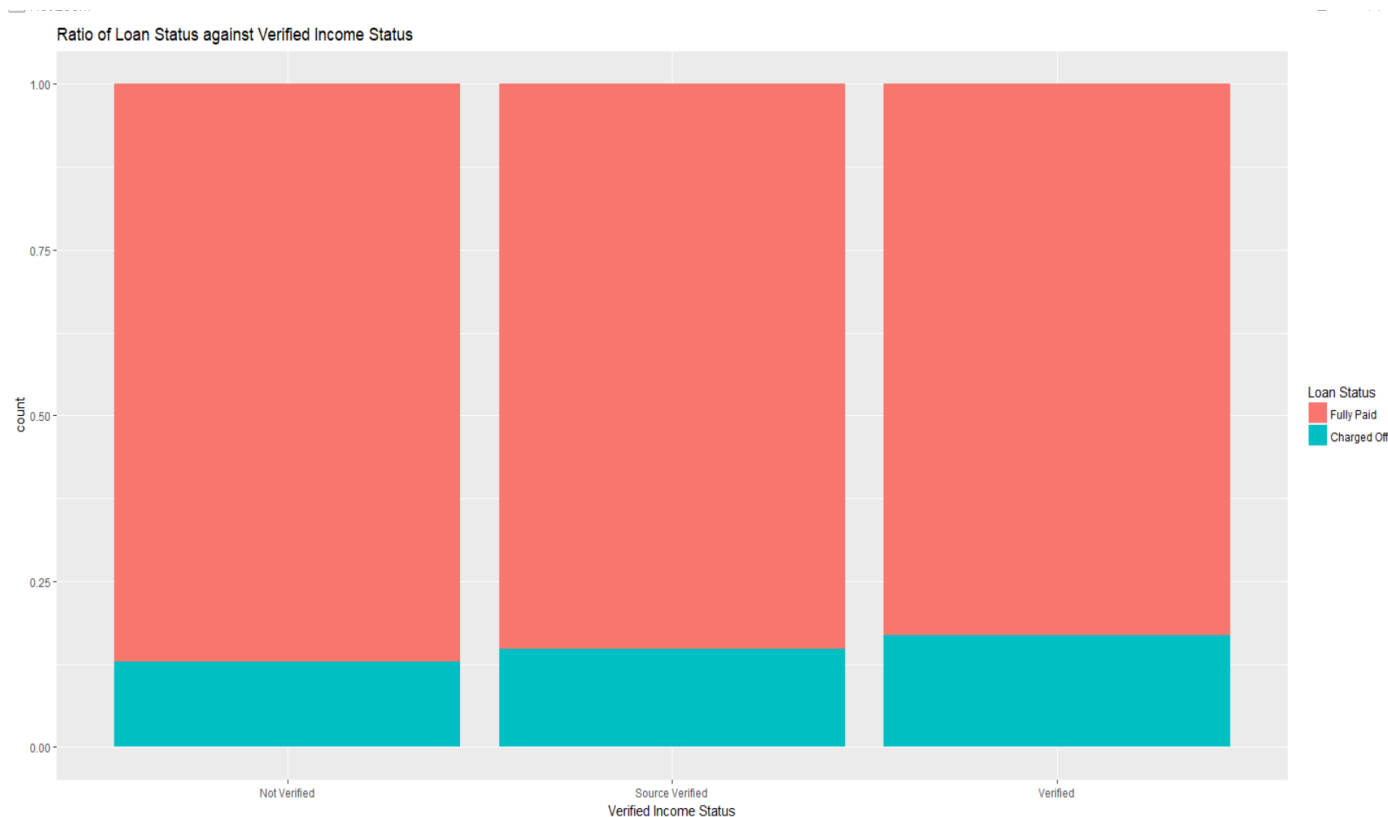


loan\$annual_inc_bkt	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
(4e+03,3.72e+04]	1413 0.183	6302 0.817	7715 0.200
(3.72e+04,5e+04]	1219 0.157	6549 0.843	7768 0.201
(5e+04,6.5e+04]	1149 0.150	6530 0.850	7679 0.199
(6.5e+04,9e+04]	1028 0.130	6901 0.870	7929 0.206
(9e+04,6e+06]	818 0.109	6667 0.891	7485 0.194
column Total	5627	32949	38576



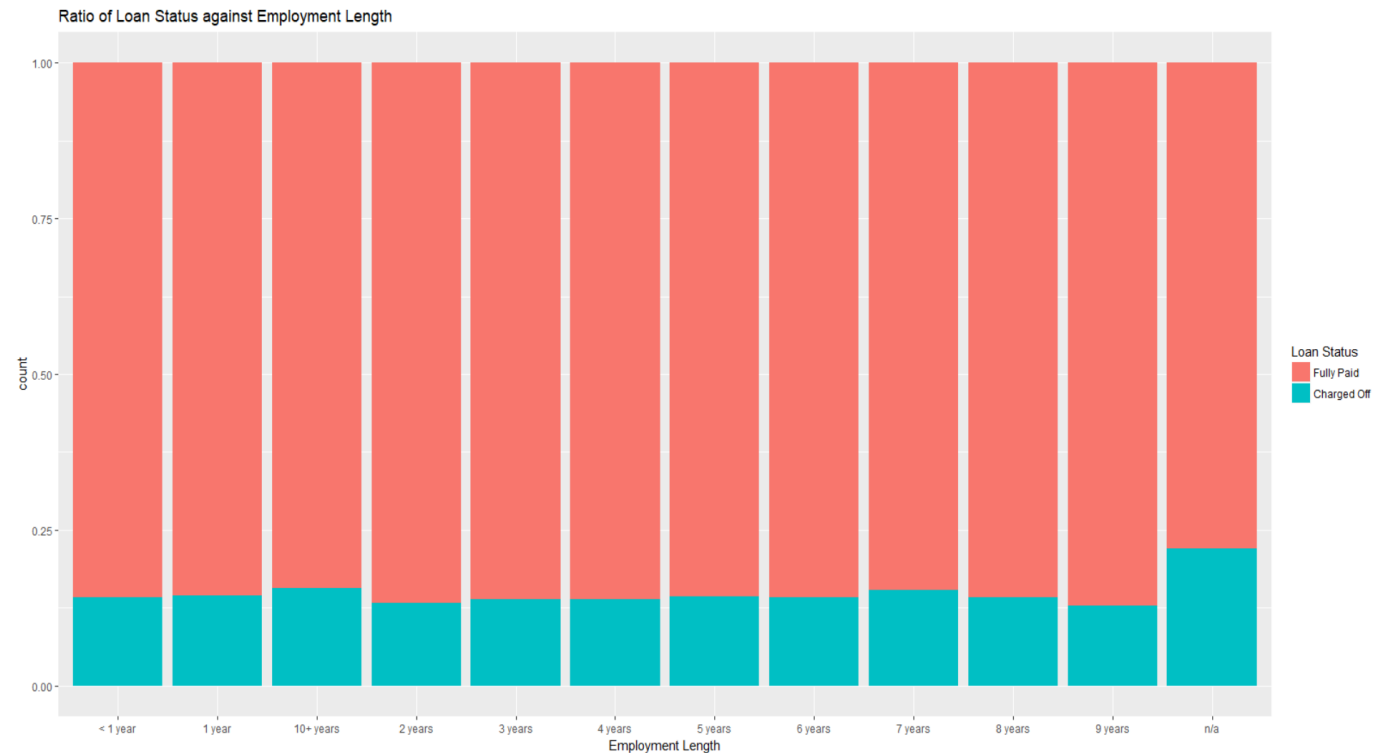
**The higher a borrower's annual income the less likely they are to default to repay their loans.**

loan\$verification_status	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
Not Verified	2142 0.128	14552 0.872	16694 0.433
Source Verified	1434 0.148	8243 0.852	9677 0.251
Verified	2051 0.168	10155 0.832	12206 0.316
Column Total	5627	32950	38577



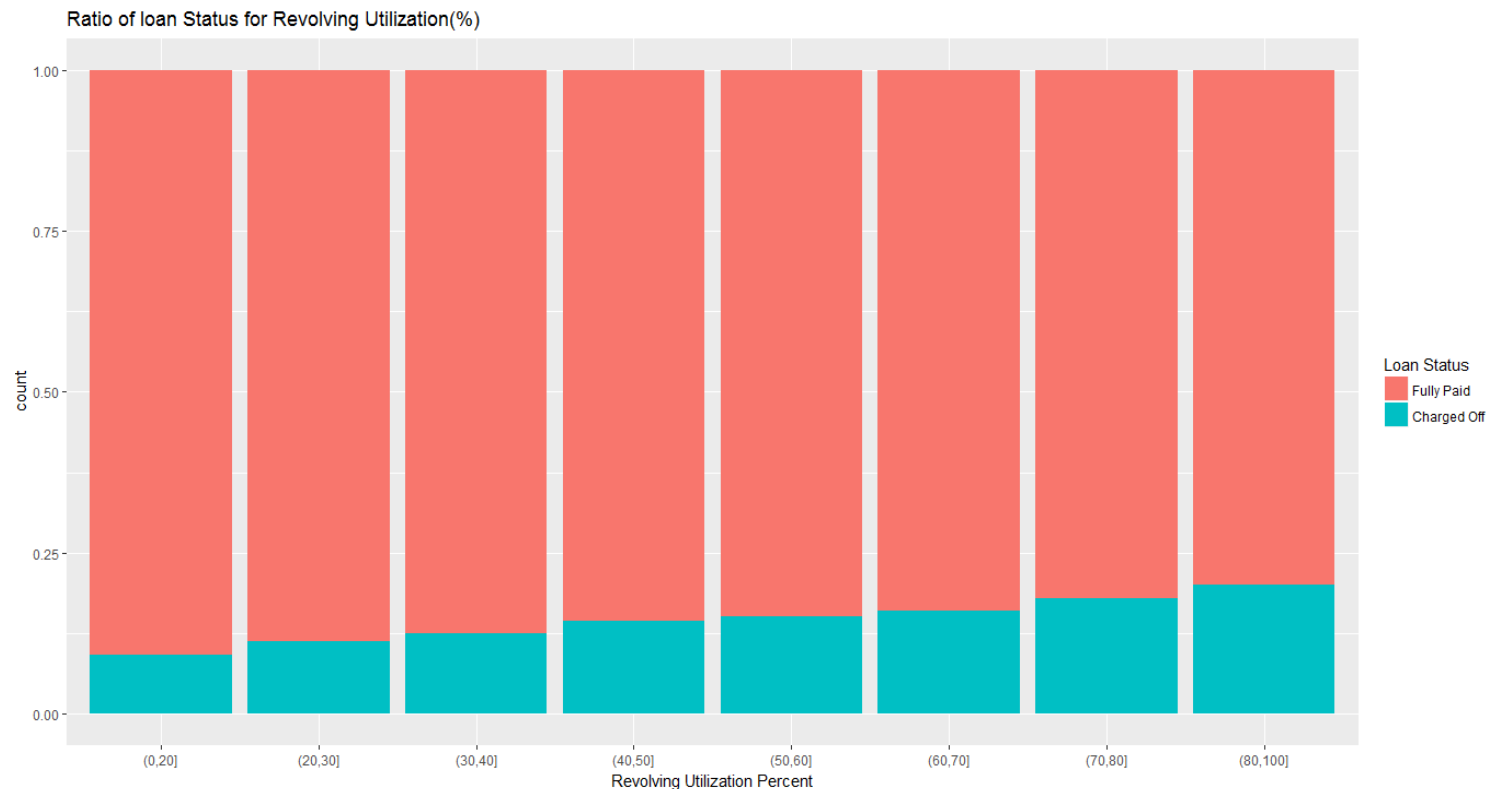
Interestingly, we see that as income verification increases, the loan performance actually worsens.

loan\$emp_length	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
< 1 year	639 0.142	3869 0.858	4508 0.117
1 year	456 0.144	2713 0.856	3169 0.082
10+ years	1331 0.157	7157 0.843	8488 0.220
2 years	567 0.132	3724 0.868	4291 0.111
3 years	555 0.138	3457 0.862	4012 0.104
4 years	462 0.138	2880 0.862	3342 0.087
5 years	458 0.143	2736 0.857	3194 0.083
6 years	307 0.142	1861 0.858	2168 0.056
7 years	263 0.154	1448 0.846	1711 0.044
8 years	203 0.141	1232 0.859	1435 0.037
9 years	158 0.129	1068 0.871	1226 0.032
n/a	228 0.221	805 0.779	1033 0.027
Column Total	5627	32950	38577



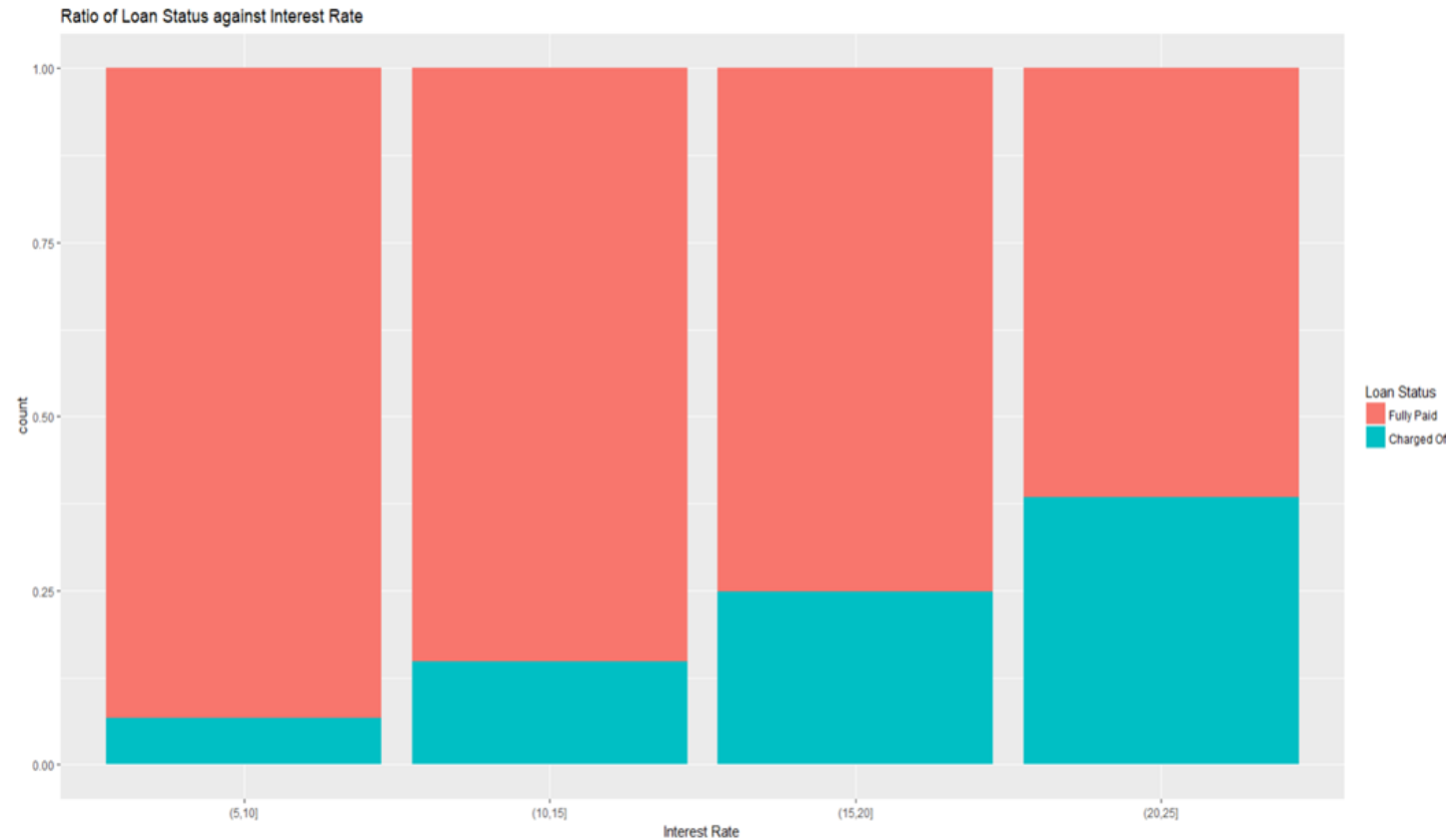
Assuming NA as Unemployed/Self-employed/Not mentioned, it clearly depicts that the unemployed are most likely to default.

loan\$revol_util_bkt	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
(0,20]	626 0.092	6180 0.908	6806 0.181
(20,30]	414 0.112	3290 0.888	3704 0.099
(30,40]	501 0.125	3500 0.875	4001 0.106
(40,50]	606 0.143	3618 0.857	4224 0.112
(50,60]	630 0.151	3554 0.849	4184 0.111
(60,70]	659 0.160	3458 0.840	4117 0.110
(70,80]	697 0.179	3200 0.821	3897 0.104
(80,100]	1334 0.201	5306 0.799	6640 0.177
Column Total	5467	32106	37573



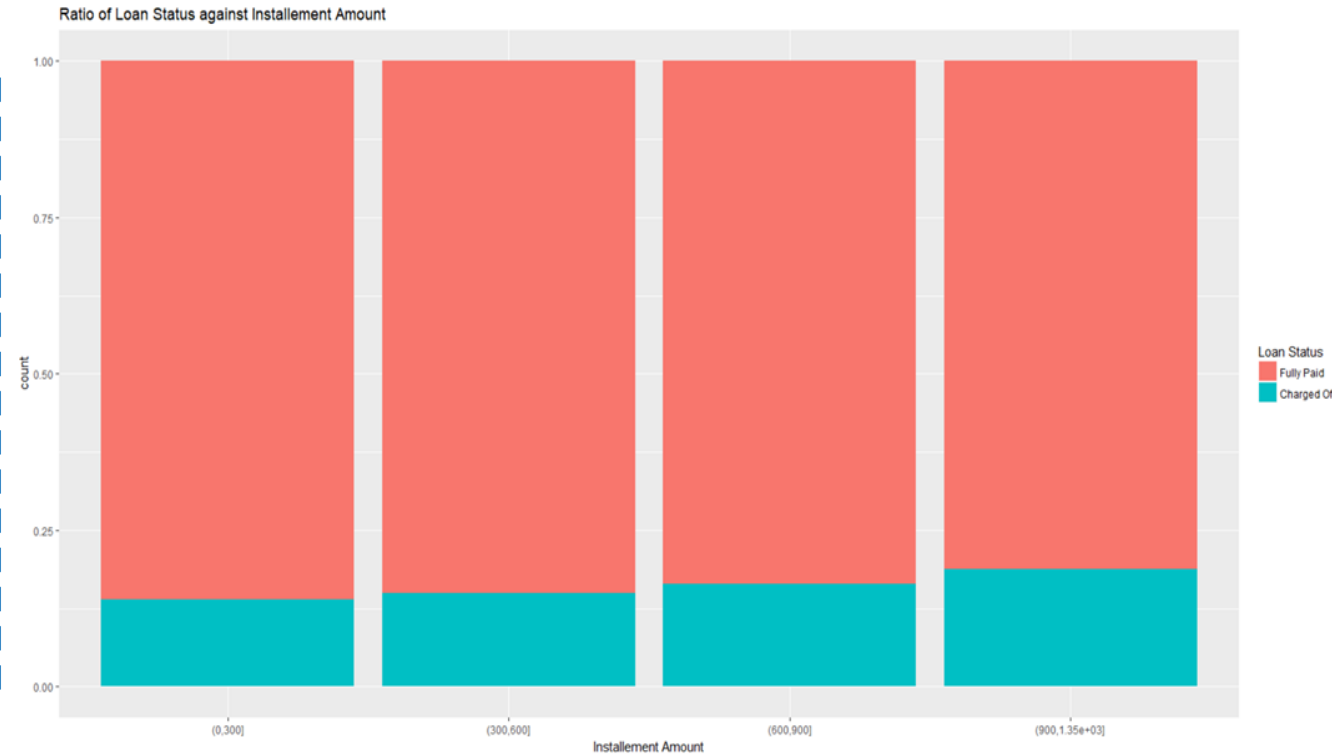
The percentage of defaulting loans steadily increases with utilization rate. Borrowers with high utilization rates are more likely to have high fixed credit card payments which might affect their ability to repay their loans.

loan\$int_rate_bkt	loan\$loan_status		Row Total
	Charged off	Fully Paid	
(5,10]	830 0.067	11486 0.933	12316 0.319
(10,15]	2707 0.148	15558 0.852	18265 0.473
(15,20]	1794 0.248	5432 0.752	7226 0.187
(20,25]	296 0.384	474 0.616	770 0.020
column Total	5627	32950	38577



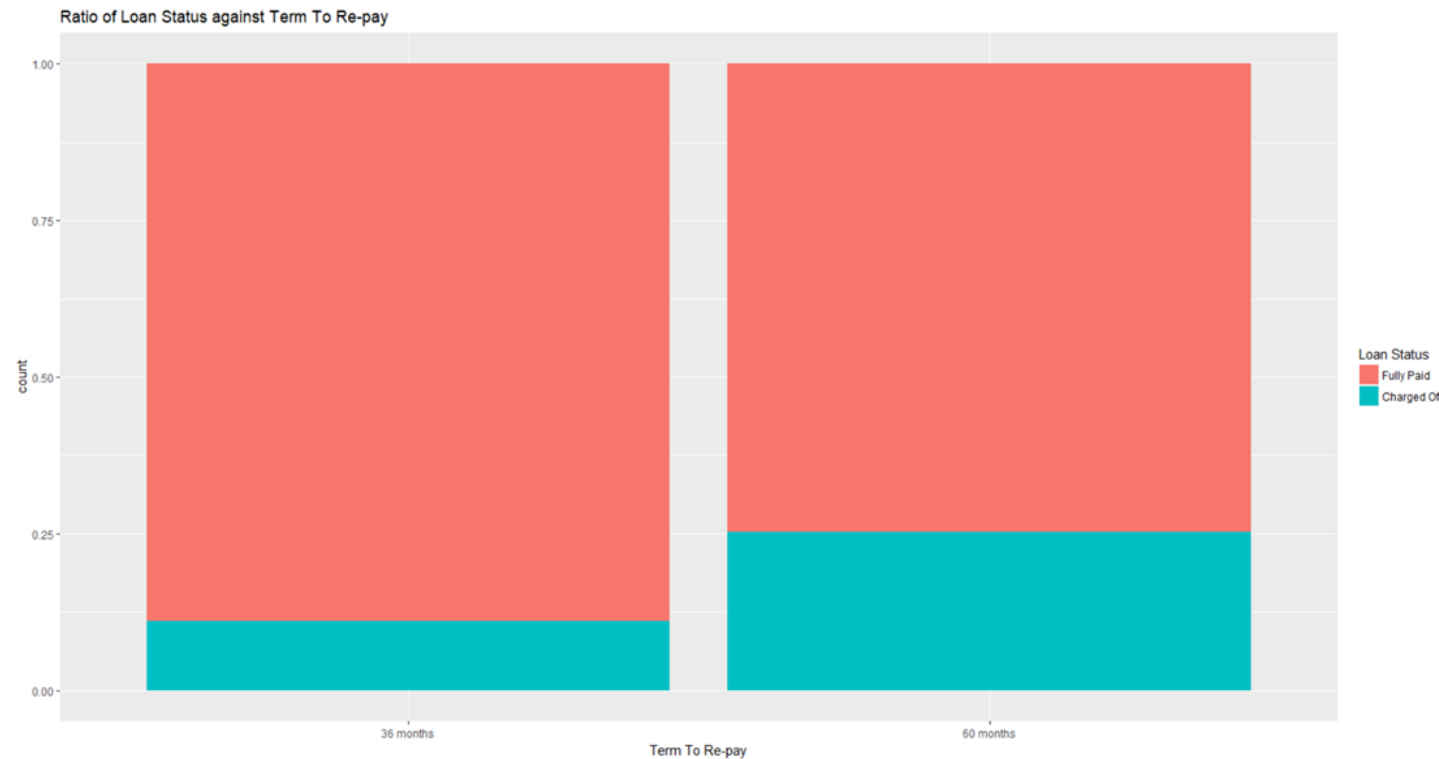
**Number of Defaulters increase with increase in interest rates.**

loan\$loan_status			
loan\$installment_bkt	Charged Off	Fully Paid	Row Total
(0,300]	2881	17801	20682
	0.139	0.861	0.536
(300,600]	2042	11645	13687
	0.149	0.851	0.355
(600,900]	606	3082	3688
	0.164	0.836	0.096
(900,1350]	98	422	520
	0.188	0.812	0.013
Column Total	5627	32950	38577



**As the instalment amount increases the defaulter's percentage increases as well.**

loan\$term	loan\$loan_status		Row Total
	Charged Off	Fully Paid	
36 months	3227 0.111	25869 0.889	29096 0.754
60 months	2400 0.253	7081 0.747	9481 0.246
Column Total	5627	32950	38577



The plot clearly shows that 36 month term has more probability of getting Charged Off.

**Grade** and **Sub-grade** variables provide the most predictive power for determining expected loan performance.

A large number of the other variables also provide strong indications of expected performance, most descriptive among them are:

- Grades/ Sub-grades /Interest Rates
- Home Ownership Status
- Loan Purpose
- Loan Amount
- Annual Income
- Verified Income Status
- Employment Length
- Revolving Utilization Percent
- Instalment Amount
- Term to Re-pay

Verified income status and show results opposite from what we would expect. This is likely due to increased standards on borrowers with poorer credit history, so all else equal we see outperformance in these loans.