**upGrad**

*#LifeKoKaroLift*

# Lending Club Case Study

Learners:

Karishma Kumari

Vinay Karandi

upGrad

# Content

- Problem statement and the analysis approach
- Results of univariate, bivariate analysis
- Visualizations and Summary

**upGrad**

**Problem Statement**
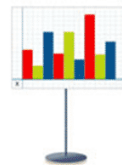
## What is Lending Club?

Lending Club is a marketplace for personal loans that matches borrowers who are seeking a loan with investors looking to lend money and make a return.



## How Lending Club Works

**Borrowers** apply for loans.
**Investors** open an account.

**Borrowers** get funded.
**Investors** build a portfolio.

**Borrowers** repay automatically.
**Investors** earn & reinvest.

When the company receives a loan application, the company has to make a decision for loan approval based on the applicant's profile. Two types of risks are associated with the bank's decision:

- If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company

- If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company
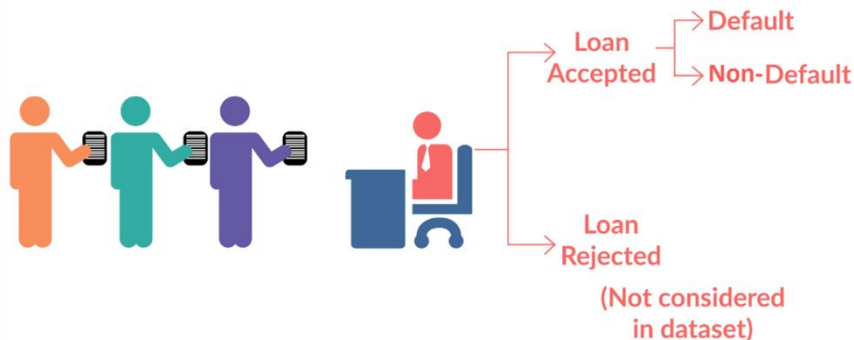
**Problem Statement**

If one is able to identify these risky loan applicants, then such loans can be reduced thereby cutting down the amount of credit loss. Identification of such applicants using EDA is the aim of this case study.

In other words, the company wants to understand the **driving factors (or driver variables)** behind loan default, i.e. the variables which are strong indicators of default. The company can utilize this knowledge for its portfolio and risk assessment.

**Problem Statement**



**LOAN DATASET**

**Fully paid:** Applicant has fully paid the loan (the principal and the interest rate)

**Current:** Applicant is in the process of paying the instalments, i.e. the tenure of the loan is not yet completed. These candidates are not labelled as 'defaulted'.

**Charged-off:** Applicant has not paid the instalments in due time for a long period of time, i.e. he/she has defaulted on the loan

**Analysis Approach**

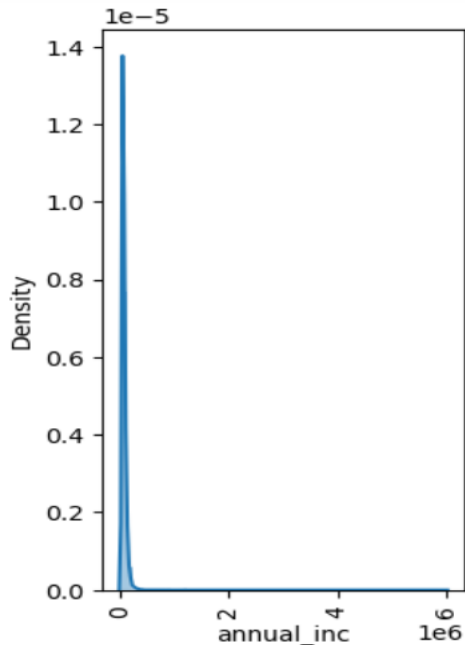There are four major steps while doing the case study:

1. Data understanding
2. Data cleaning (cleaning missing values, removing redundant columns etc.)
3. Data Analysis- Univariate and bivariate analysis
4. Summary and conclusion

**Data Cleaning**

1. Checked the duplicate values
2. Checked the percentage of missing values – the columns having more than 20% NAN values are removed
3. Dropped the rows having NAN values
4. Fixed the datatype of numerical and date columns
5. Removed the columns having just 1 unique value from the analysis

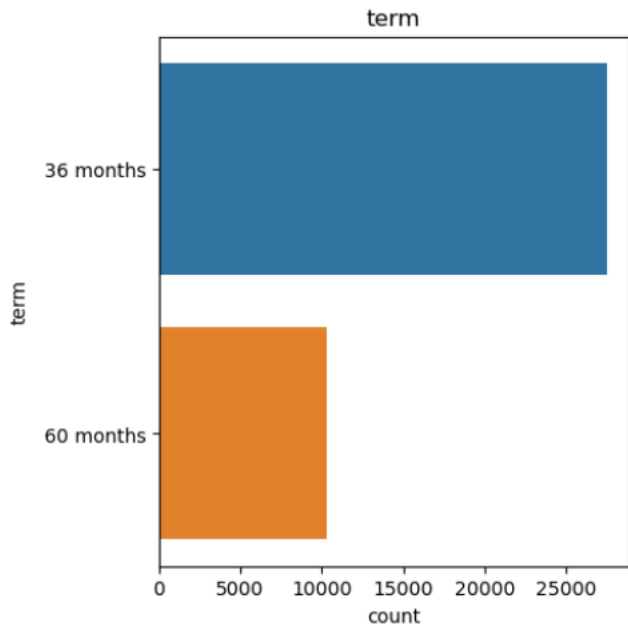**Data Analysis – Univariate Analysis – Continuous variables**



From Univariate analysis of annual_inc variable, it is observed that the data for this is skewed. Hence, creating buckets for this variable.

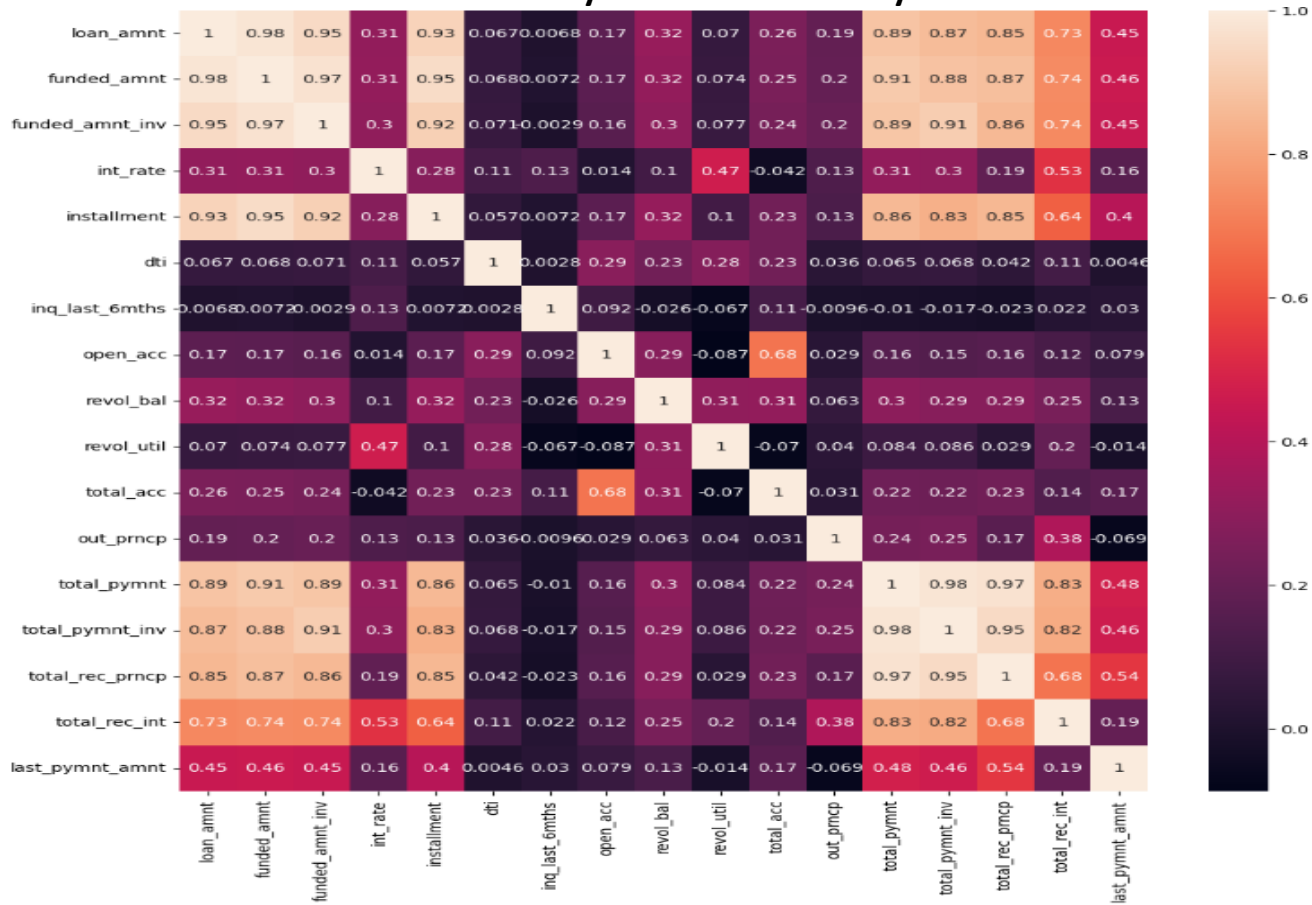Similarly, buckets have been created for the following variables:

1. Annual income
2. delinq_2yrs
3. pub_rec
4. out_prncp_inv
5. total_rec_late_fee
6. recoveries
7. collection_recovery_fee
8. pub_rec_bankruptcies

**Data Analysis – Univariate Analysis – Categorical variables**



To analyse the categorical variables, countplot has been used to see the frequency of each variable.
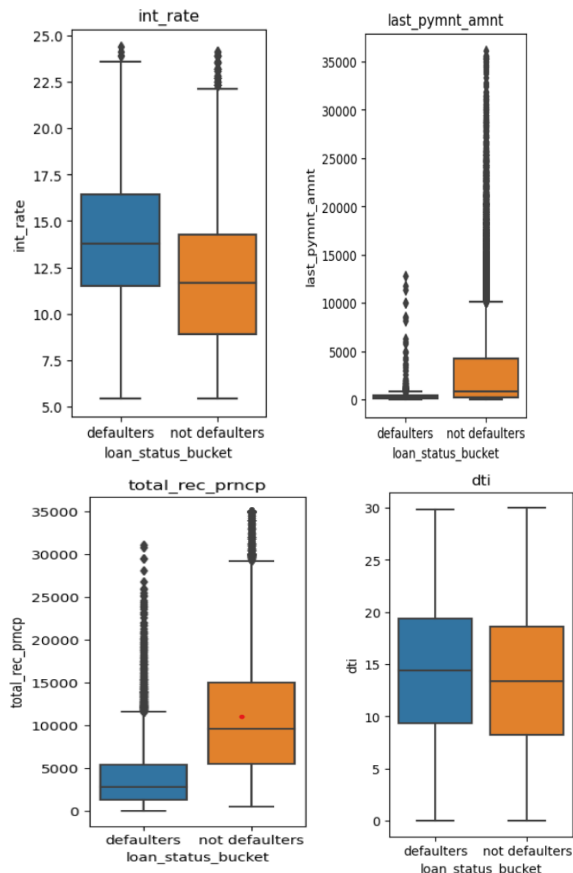
## Data Analysis – Bivariate Analysis – Continuous variables



To analyse which all Continuous variables are correlated, correlation matrix has been used. It can be observed that the following set of variables are high correlated:

- Loan Amnt, Funded Amnt and Funded Amnt_inv are highly correlated
- 'total_pymnt', 'total_pymnt_inv' and 'total_rec_prncp' are highly correlated
- funded_amnt and total_pymnt are highly correlated
- total_pymnt_inv and funded_amnt_inv are highly correlated

**Data Analysis – Bivariate Analysis – Continuous variables**



To analyse the Continuous variables' impact on loan_status, box plots for each level of loan_status_bucket (categorical variable) has been drawn.
It can be observed that the following Continuous variables are influencing loan_status_bucket:

- int_rate, dti, inq_lat_6mths, revol_util, total_pymnt, total_pymnt_inv, total_rec_prncp and last_pymnt_amnt
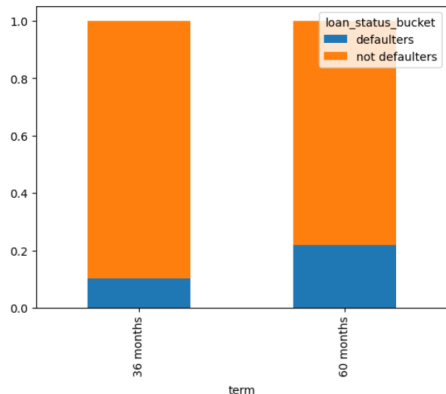
Since from the data dictionary definitions, we understood that the following variables are the customer behavior variables which are not available at the time of loan application, and thus they cannot be used as predictors for credit approval. Hence, they can not be taken into consideration to identify the chances of defaults and the risky loan applicants:

- total_pymnt, total_pymnt_inv, total_rec_prncp and last_pymnt_amnt

Hence, the final list of continuous variables influencing loan_status_bucket are:

- **int_rate, dti, inq_lat_6mths, revol_util**

## Data Analysis – Bivariate Analysis – Categorical variables



To analyse the Categorical variables' impact on loan_status, stacked bar charts for different categories have been created.

It can be observed that the following Categorical variables are influencing loan_status_bucket:

- Term, grade, sub grade, purpose, addr_state, zip code, ann_inc_bucket, 'out_prncp_cat', 'out_prncp_inv_cat', 'total_rec_late_fee_cat', 'recoveries_cat', 'collection_recovery_fee_cat'
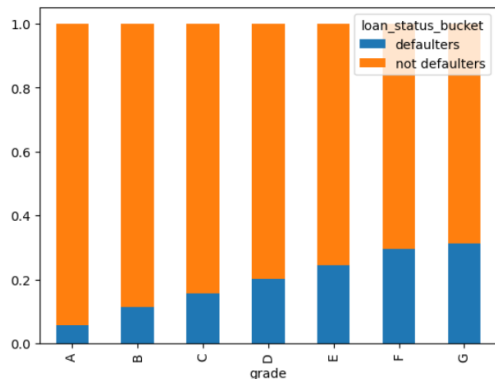
Since from the data dictionary definitions, we understood that the following variables are the customer behavior variables which are not available at the time of loan application, and thus they cannot be used as predictors for credit approval. Hence, they can not be taken into consideration to identify the chances of defaults and the risky loan applicants:

- 'out_prncp_cat', out_prncp_inv_cat, total_rec_late_fee_cat, recoveries_cat, collection_recovery_fee_cat
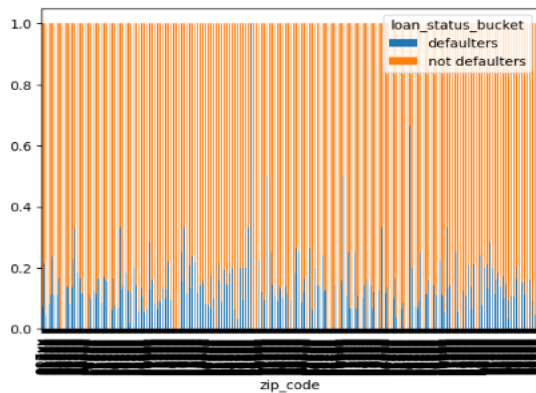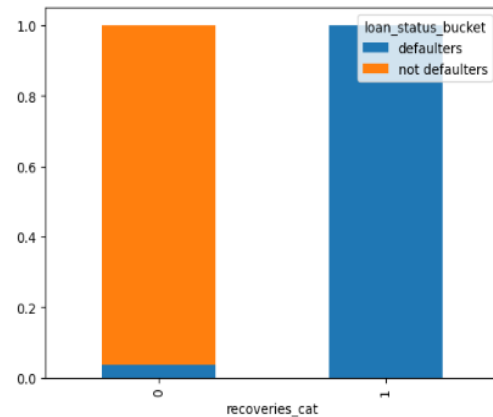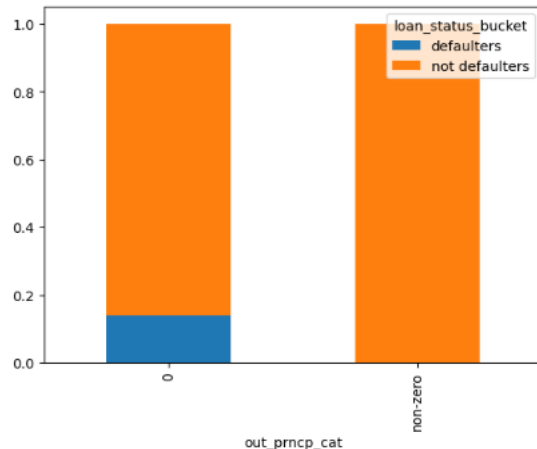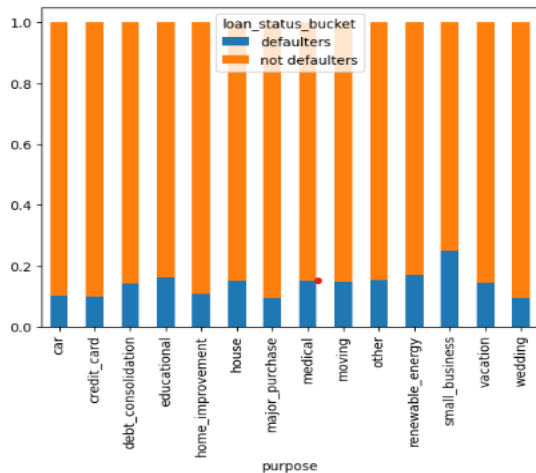
Also, since grade and sub-grade are related, taking the most granular info into consideration i.e. Sub-Grade. Similarly, since addr_state and zip code are geographic related variables, taking the most granular info into consideration i.e. zip code.

Hence, the final list of categorical variables influencing loan_status_bucket are:

- **Term, sub grade, purpose, zip code, ann_inc_bucket**



13

## Data Analysis – Bivariate Analysis – Categorical variables

## Conclusion

Final list of categorical variables influencing loan_status_bucket are:

- term
- sub grade
- purpose
- zip code

Final list of continuous variables influencing loan_status_bucket are:

- int_rate
- dti
- inq_lat_6mths
- revol_util
- ann_inc