# Gramener Case Study

By

Raghunandan Vellanki

Vyom Bhatt

Vasanthkumar Subramanian

# Abstract

A consumer finance company, the largest online loan market place, has many risk applicants who default the loans and thereby causing credit loss to the company and its investors.

The primary objective of the study is to identify potential applicants who may cause credit loss to the company by defaulting the loan.

The analysis that we have done is

- Using EDA, understanding how the consumer and loan attributes influence the tendency of loan default
- Identifying the driving factors which cause the applicants to default the loans

Stages involved in the Analysis

- Data Cleansing and Manipulation
- Exploratory Data Analysis
- Summarizing the output and answering the business problem

UpGrad

# Data Cleaning

For the analysis, we have taken into consideration the following **assumptions**:

- Nothing outside the scope of the analysis is to be considered for the desired output

- All the monetary figures in the dataset are in the same currency

In order to get the data prepared for our EDA, we went ahead and cleansed the data in the following manner:

- Excluded all the columns that were filled with over 90% null values

- On further observation, we got rid of columns that didn't seem fit for analysis based on the data description

- Dropping all columns that had only 1 unique value – as these won't be of any impact to our final output

- Corrected the format for the date columns to a date format for processing in Python

- Cleaning the columns such as:

    - Term column to exclude the ' months' string and change it to numeric

    - Employee length column to change it to numeric

    - Interest rate column to change it to decimal format

- Bucketing the interest rates, income levels, DTI for understanding trends in these variables later

UpGrad

# Deriving New Metrics

Based on the data provided, we went ahead and derived new variables to assist our analysis further. These include:

- Extracting the year from issue date column

- Extracted the month from the issue date column

We went ahead and created a new metric that would give us a holistic view of how every variable is impacting the rate at which loans are getting charged off i.e **Charged-Off Index.**

> This index takes into account the distribution of users that are charged off and the distribution of the total population i.e. pct_charged_off / pct_total_population. The reason for doing this is to ensure that we don't ignore the entire population set and just consider the charged_off/fully_paid users, as this will create a bias in case the population is skewed towards one entity

In addition to the above, we also created a business driven metric to understand the performace of each variable with regard to the nature of the business. The metric **payback percentage** which is essentially (total_rec_prncp / loan_amnt)*100
This metric will tell us what percentage of the loans are getting paid back – the lower it is, the riskier the loan gets as we haven't received the amount.
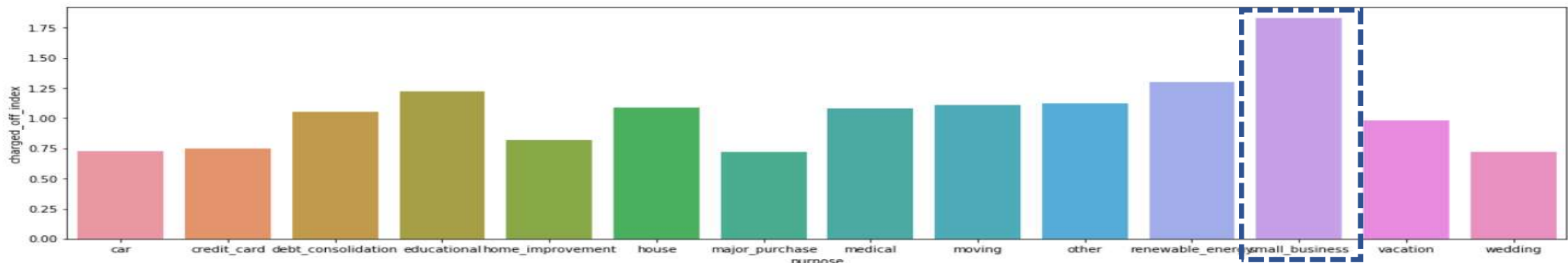
# EDA

Running through the code submitted, we have gone through an extensive piece of data cleaning that involves:

- Data Cleaning and Preparation (as discussed previously)

- Univariate analysis

    - Outlier Analysis and treatment

    - Analyzing distributions for independent variables

- Segmented Univariate analysis

    - Analyzing independent variables by segmenting them with other variables to observe trends and impact on the dependent variable

- Bivariate Analysis

    - Analyzing multiple variables to identify how they impact the dependent variable

    - Prepared a correlation matrix for all the variables that we considered

- Using the Dervied Metrics to perform a deep-dive analysis on important variables identified
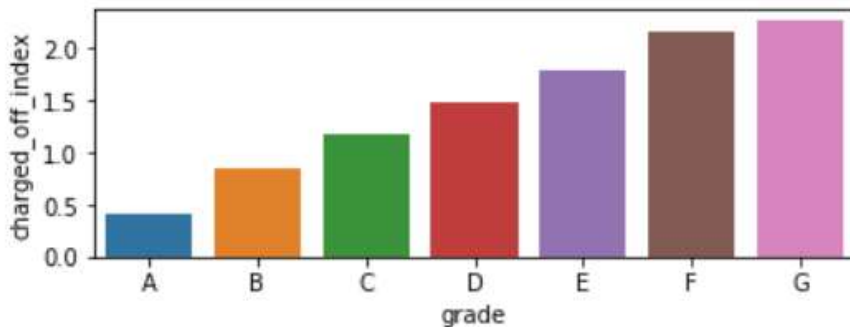
In the upcoming sections of this presentation, we will present the major findings of our EDA, and present the shortlisted variables that we identified as the most important in terms of impacting whether a loan would get defaulted or not
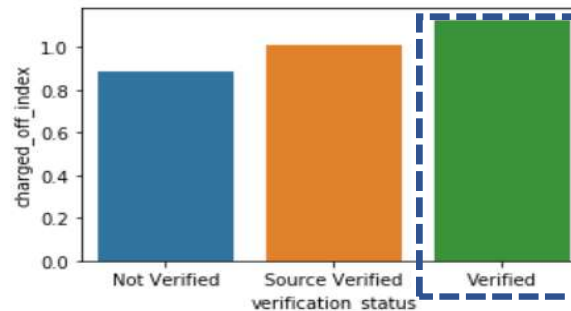
UpGrad

# EDA Findings - 1

1. Among all the purposes, the 'Small Businesses' loan purpose is the riskiest with the highest Charged-off index, however the most loans are taken for the Debt Consolidation purpose



2. We observe a trend in the 'Grade' variable, as the grade reduces, the risk of the loan getting defaulted increases
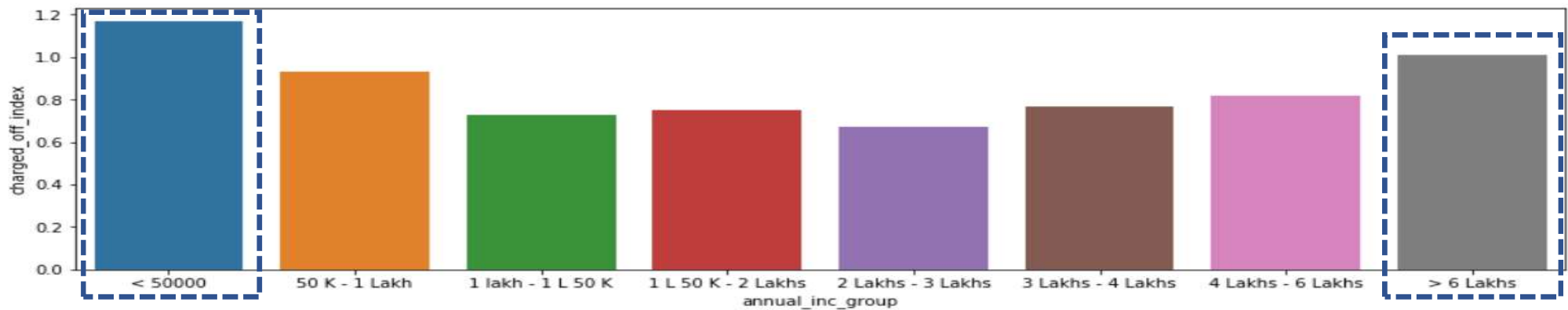


3. the data for verification is one that stands out as it goes against what one would expect: the loans with a verified status have a higher risk compared to not-verified ones.
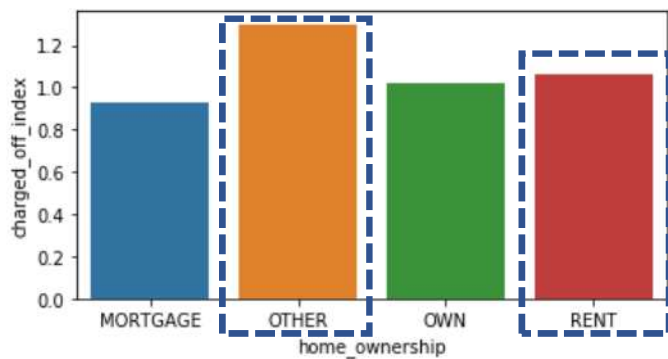


Perhaps the firm should look to make the verification process more stringent to improve this so that the verified applicants are actually ones that don't end up defaulting

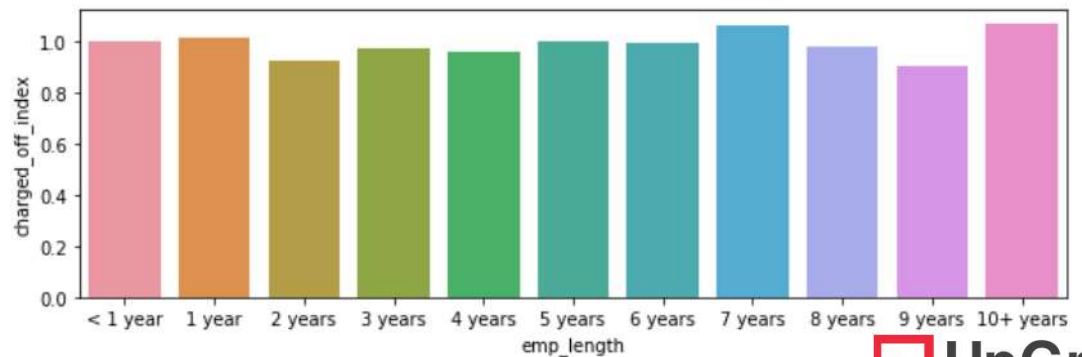UpGrad

# EDA Findings - 2

4. The income levels is **not** the most important variable in determining the chances of the loan getting defaulted as it doesn't follow any trend



5. For the home ownership variable, the **other** and **rent** are the riskiest home ownership types
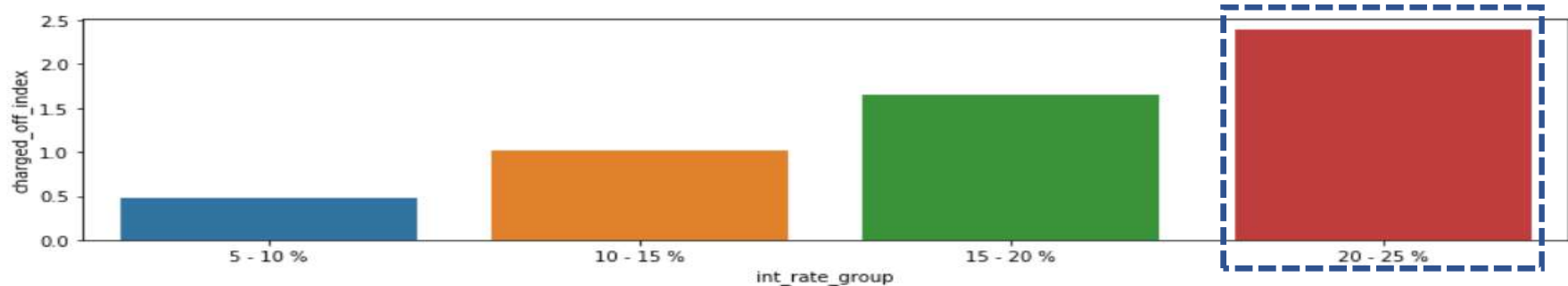


6. The employment length in years variable does not show to have any major impact on the chances of a loan getting charged-off
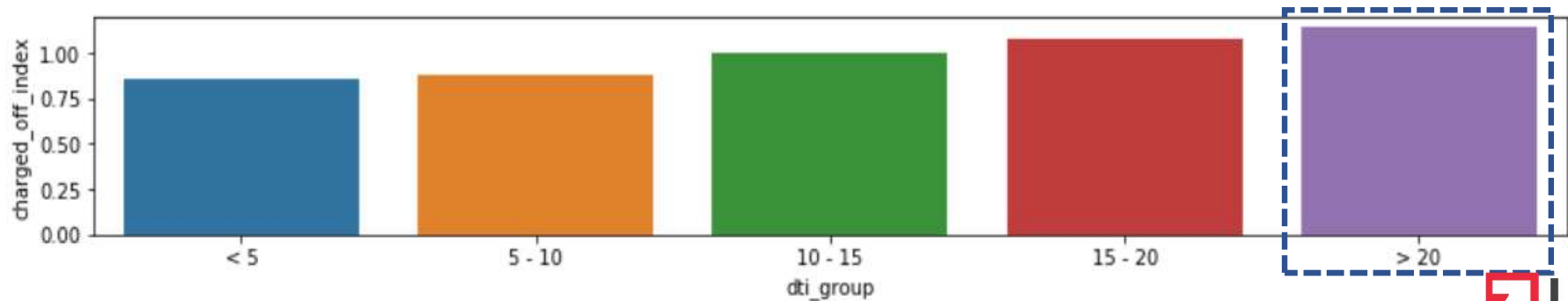


UpGrad

# EDA Findings - 3

7. Analysis on the public bankruptcy level tells us that there is a high rate of loans getting defaulted for applicants who have a high number of bankruptcy. However, **94%** of the data provided was for users that had 0 numbers of public bankruptcy, hence actioning based on this wouldn't be advised due to population skewness

8. As the interest rate increases, the rate at which the users default also increases. This, however, could be as a result of the action that the bank takes for charging higher interest rates to users that are assumed to be of high risk
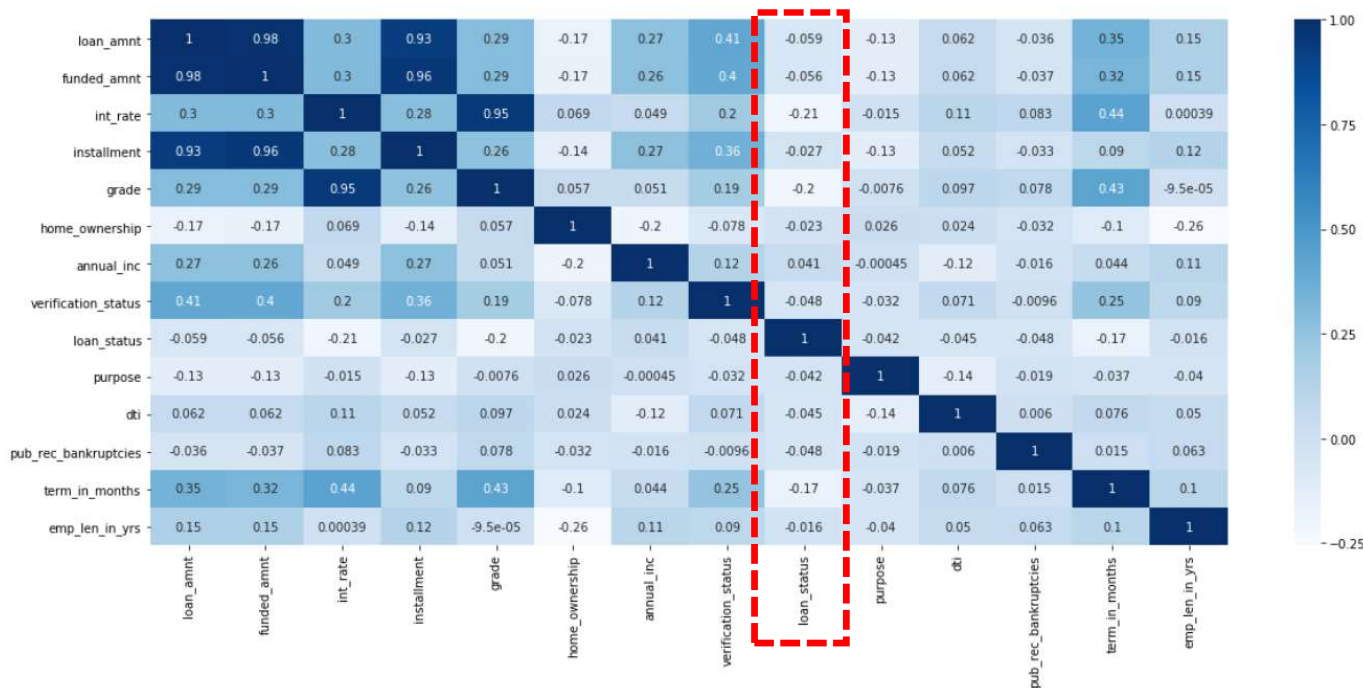


9. The DTI groups also have a direct impact on the rate at which loans get charged off. As the DTI rate increases, the loan gets more riskier



UpGrad

# EDA Findings - 4

On performing Bi-variate analysis for every variable against each other, we observed correlations between the variables to see what variables are representing string relationships against each other.

The dataset with the dependent variable as loan status was filtered out only for loans that had status as 'charged off' or 'fully paid' to observe the relationship every other variable was having with respect to this dependent variable. The categorical variables here were fist converted to factors before the correlations were found.



Looking at the correlation matrix above, we observe that the variables correlating the most with our dependent variable "loan status" are: interest rate, grade, and month.
However, all of these are at most mildly correlated with the dependent variable

# The Driving Factors

Based on the exhaustive exploratory data analysis done above, the factors that we identified as the most crucial ones that impact the chances of loan defaults are:

- ❑ Grade
- ❑ Purpose
- ❑ Verification Status (This process needs to be looked at and ,made more stringent by the firm)
- ❑ Home Ownership
- ❑ Interest Rate
- ❑ DTI

UpGrad