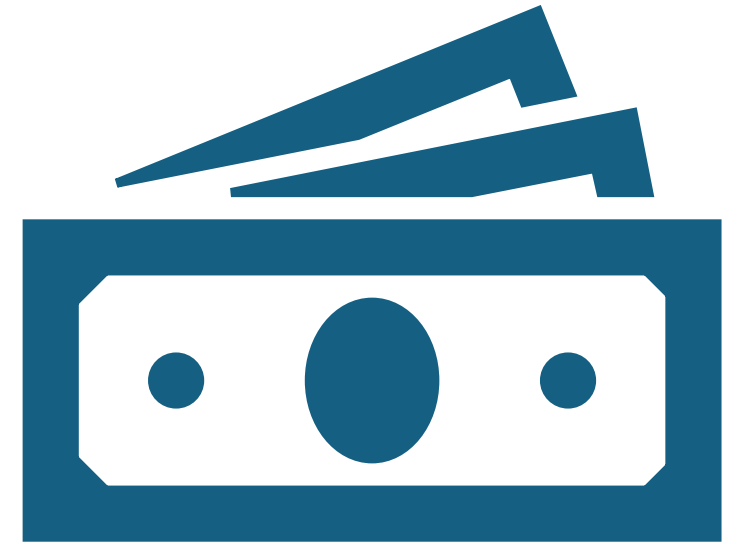


Lending Club case study

- Sathiyathan Subramanian
- Ushasis saha



Problem Description



Our consumer finance company faces the challenge of credit loss, where loan approvals to unlikely repayers may result in credit losses.



The dataset given contains information on past loan applicants and their default status.



Our objective is to understand how consumer and loan attributes influence loan default tendencies. Our goal is to mitigate credit loss by identifying the driver variables behind loan default

Approach

- Ignored all null columns – There are 55 columns that are completely null
- Ensured target variable integrity: No incorrect/null values identified
- Removed columns with only one unique value across all records.
- Columns with Categorical variables and high values for all records were removed
Ex: 'member_id' this would have unique value for all rows unique
- Removed variables related to post-loan recovery (e.g., 'recoveries', 'collection_recovery_fee') as they don't influence loan approval. Confirmed unique member IDs, ensuring each user has only one loan with the company.

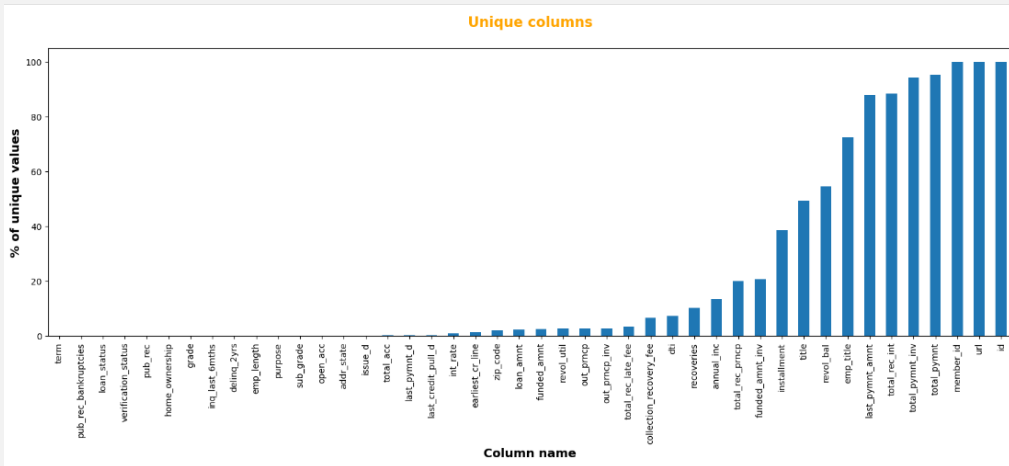
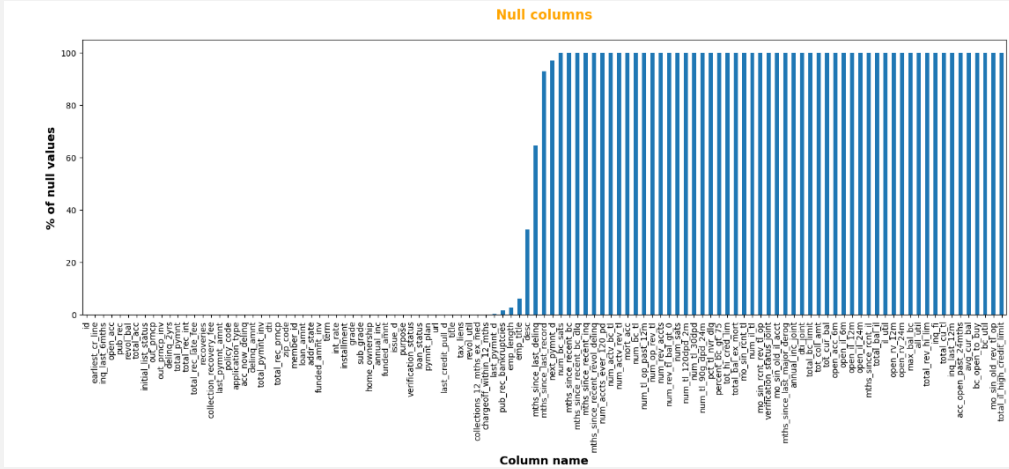
Ignored all null columns – There are 55 columns that are completely null

Ensured target variable integrity: No incorrect/null values identified

Removed columns with only one unique value across all records.

Columns with Categorical variables and high values for all records were removed
Ex: 'member_id' this would have unique value for all rows unique

Removed variables related to post-loan recovery (e.g., 'recoveries', 'collection_recovery_fee') as they don't influence loan approval. Confirmed unique member IDs, ensuring each user has only one loan with the company.



Approach



Manipulation of Datatypes:

Handling Date and String Variables
Appropriately



Data Quality Checks:

Identification and Treatment of Missing
Values

Removal of Rows with 'current' Loan
Status

Treat outliers



Data Segmentation:

Splitting the Data into Two Parts:

- Defaulters
- Fully Paid Cases

Approach

Classified Variables as Numeric and Categorical (Ordered, Unordered)

Univariate Analysis:

- Numeric Variables: Distribution Plot comparing Paid vs. Defaulters
- Ordered Categorical Variables: Segmented by Bucketing, Analyzed with Bar and Pie Charts
- Unordered Categorical Variables: Analyzed with Grouped Line Chart, Bar Chart, Pie Chart, and Grouped Columns

Identified 10 Driver Variables Based on Univariate Analysis

Approach



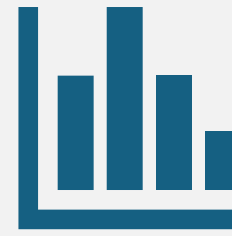
Bivariate Analysis:

Numeric Variables:

- Pair Plot: Visualizes correlations between numeric variables.

Categorical Variables:

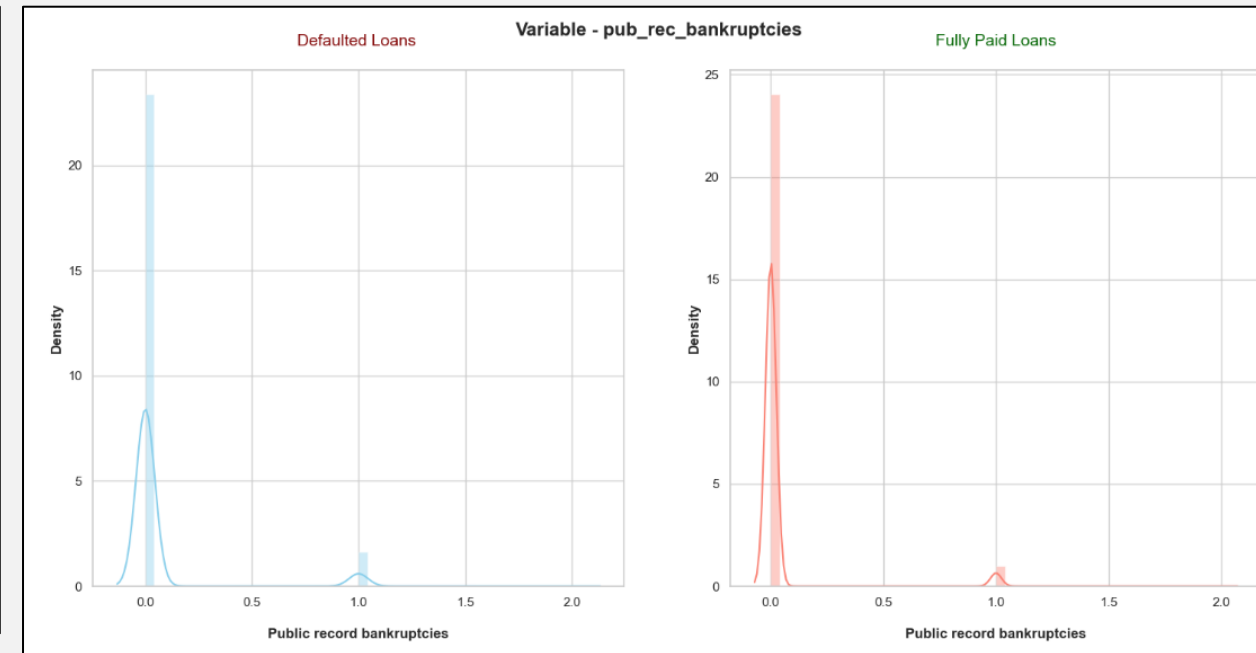
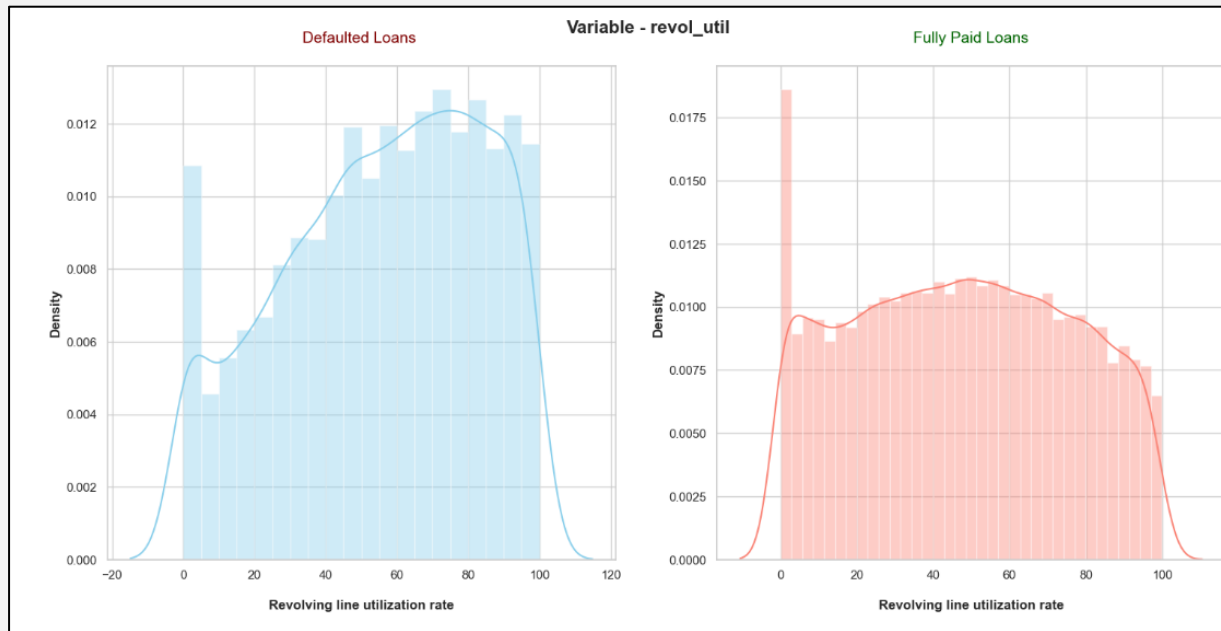
- Pivot Table + Heatmap: Illustrates correlations between categorical variables.



Identified Relationships:

Discovered 2 related categorical variables, 2 related numerical variables; reducing total driver variables to 8.

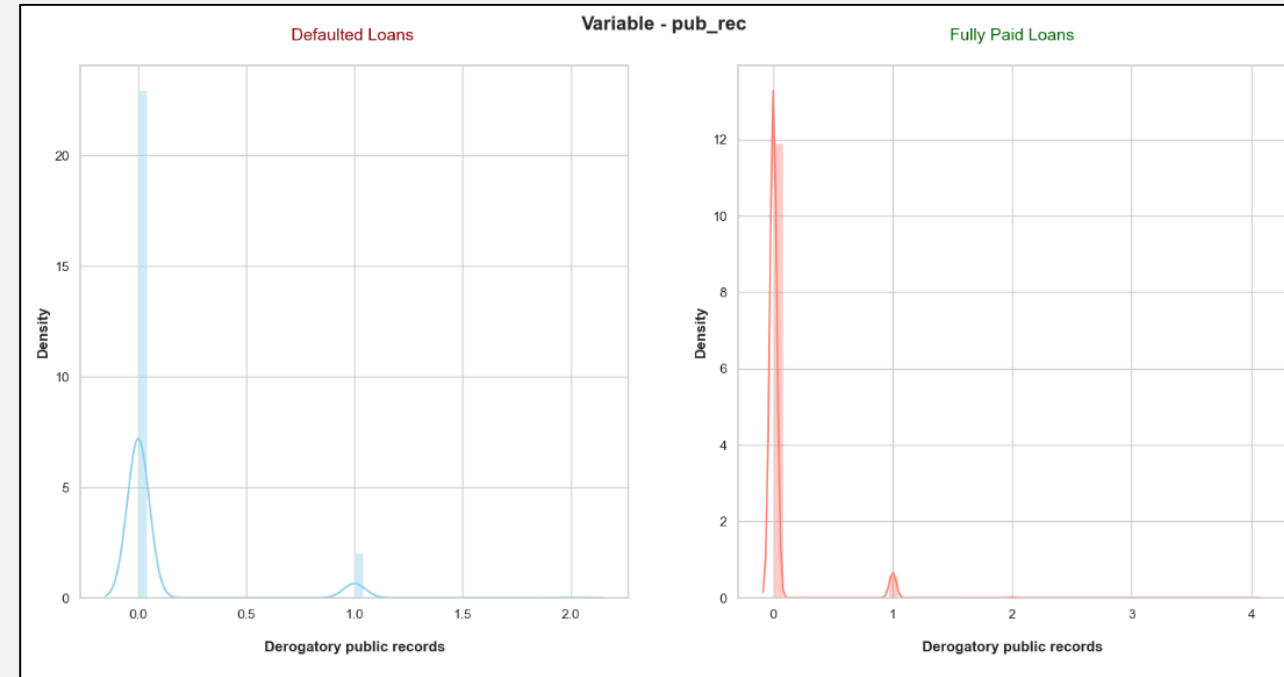
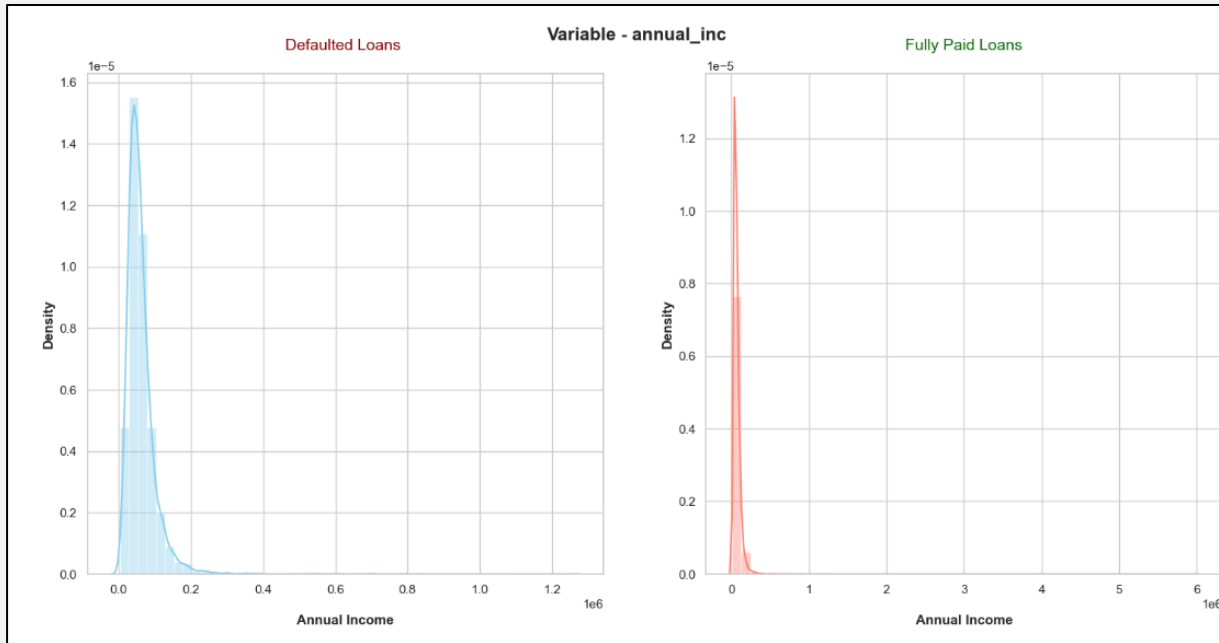
Numeric Variables



revol_util: We could see that PDC is different. Density peaks especially around 60-80%. Therefore this column will be considered

pub_rec_bankruptcies: We could see that PDC is different. Therefore, this column will be considered

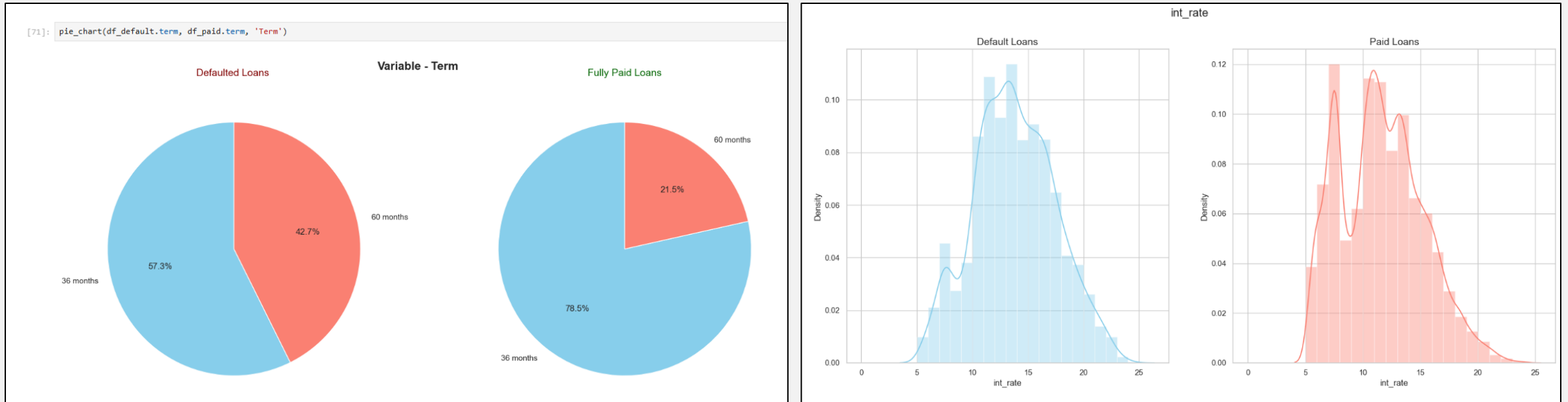
Numeric Variables



annual_inc: Unable to comment on annual_inc at this moment, the spike seems to be similar; but the range of annual_inc is different. Could be due to the outliers, we need to revisit to check the outliers. We will consider this column

pub_rec: We could see that PDC is different. Therefore this column will be considered

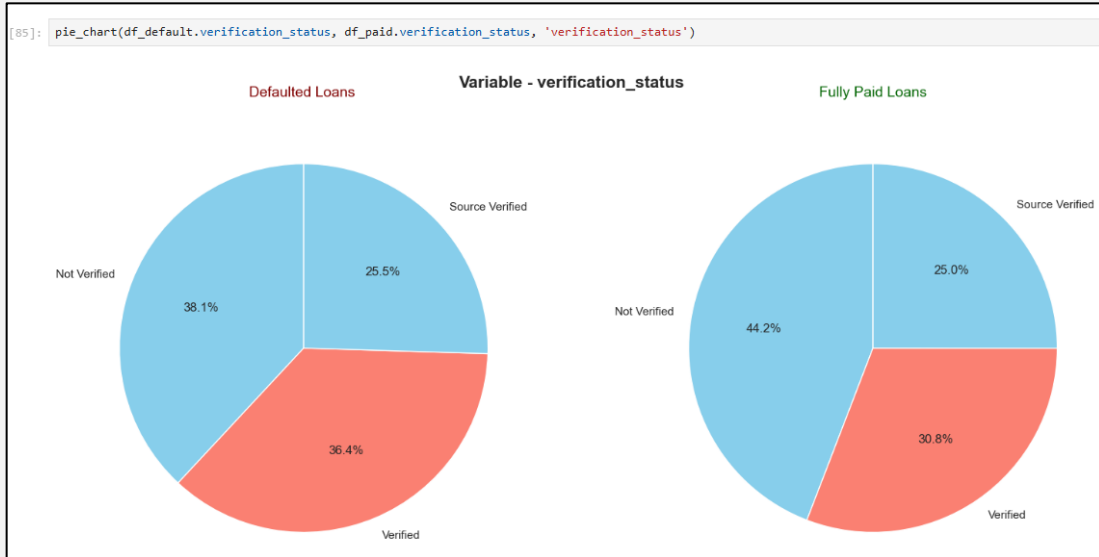
Ordered categorical Driver variables



term: From a side-by-side bar plot, we could see that there is a difference in distribution for the 60-month segment, hence drilled further into this by plotting a pie chart, we could see that in the case of Defaulted Loans, the 60-month term contributes a lot compared to Fully Paid.

int_rate: We could see that there is a difference in PDC. In fully paid loans, there is a dip at 8, 9% but in default, it is a steady spike. Let's consider this variable as a driver variable.

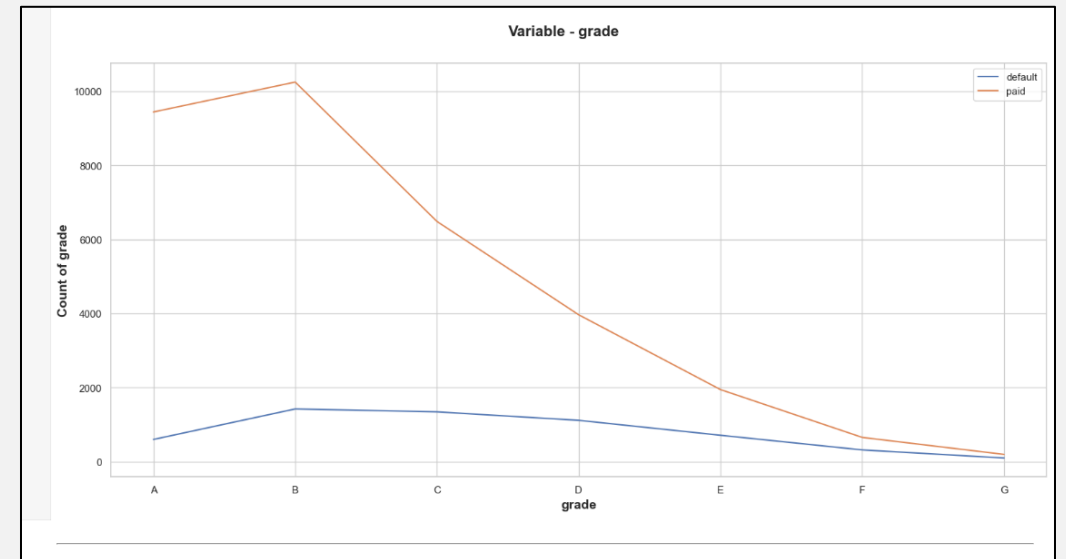
Unordered categorical Driver variables



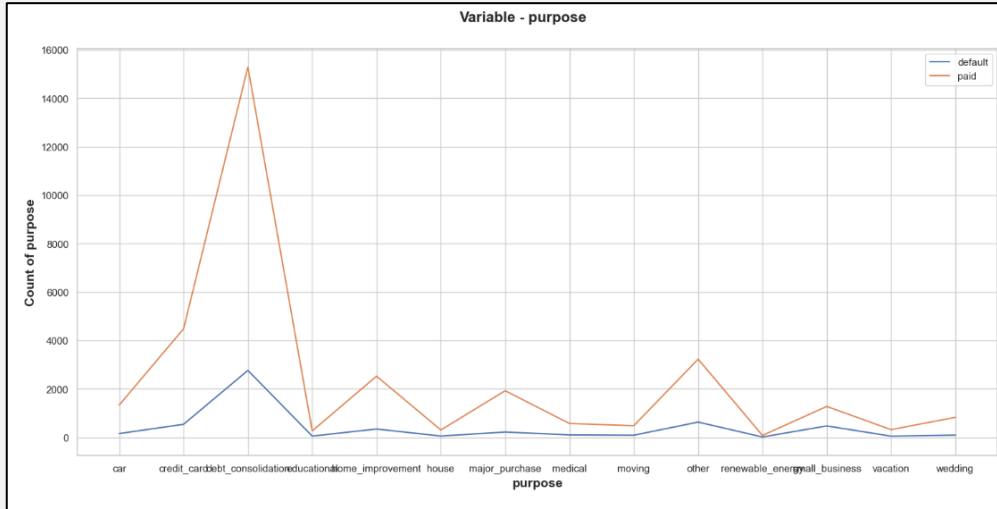
verification_status: From Bar plot, We could see a difference in pattern especially for verified, lets dig further with Pipe chart. The distribution might look similar, but there is a 6% share difference b/w verified and non-verified. Therefore, lets consider this variable for driver variable.

grade: From a side-by side bar plot, We can see a difference in the pattern especially for Grade A

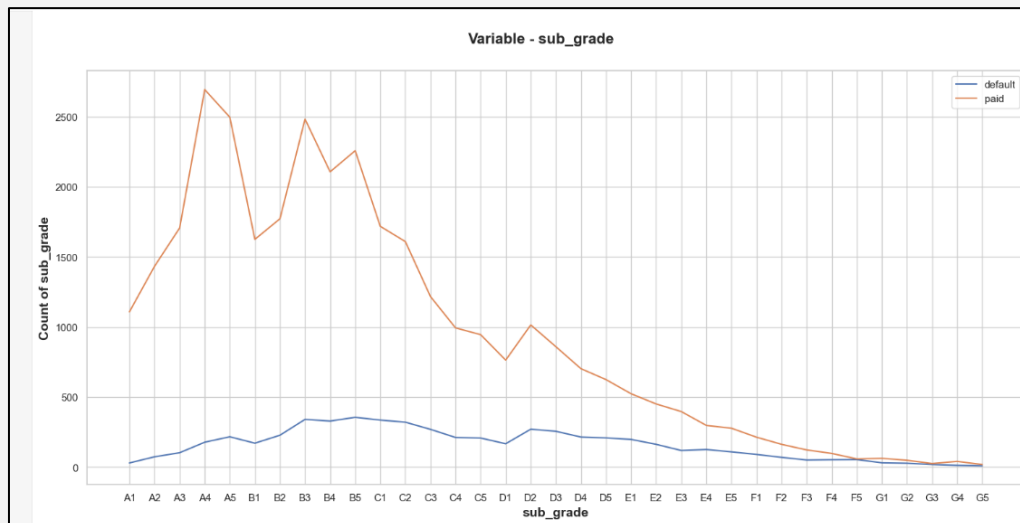
To drill further, plotted this in group line chart, we could see that the difference in the line plot. The difference b/w paid, and default keeps reducing towards right of the plot. Therefore, let's consider this column for driver variable



Unordered categorical Driver variables



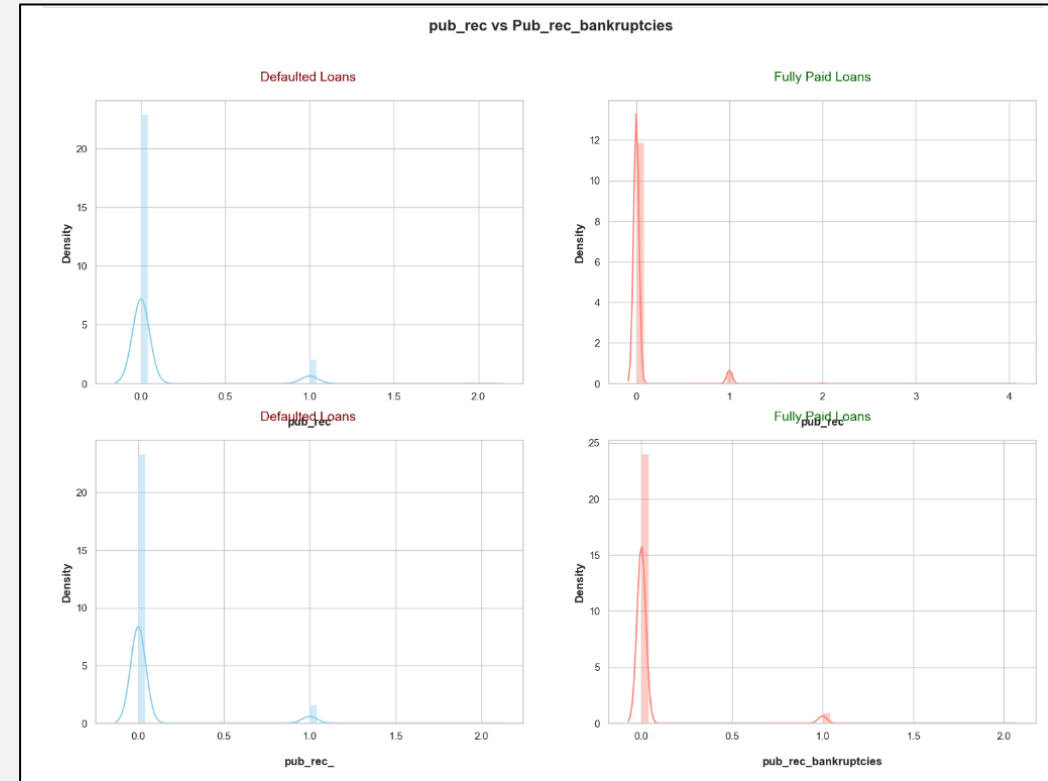
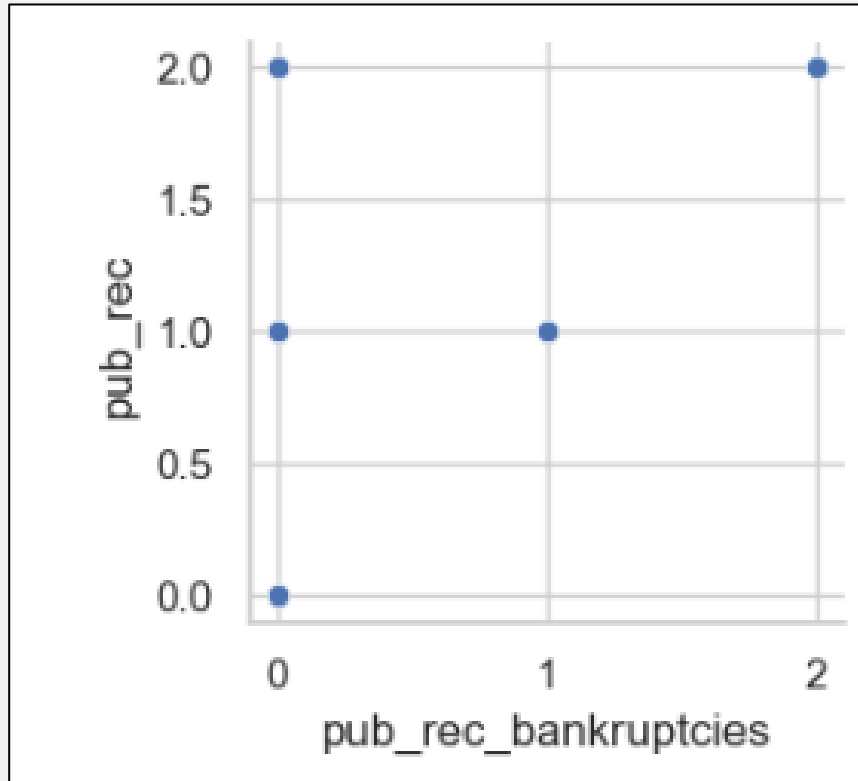
purpose: The pattern from a bar chart looked similar, but looking closely, we could see small_business proposition is high in Defaulted loans than fully paid. Plotting it in group column chart, we could see the count of small_business defaulters is 50% equivalent to that of paid ones. Debt consolidation seems to be safer one. Therefore, let's consider this variable as well to investigate further.



sub_grade: similar difference like observed in grade. The difference b/w paid, and default keeps reducing towards right of the plot. Therefore, this column will be considered as a driver variable

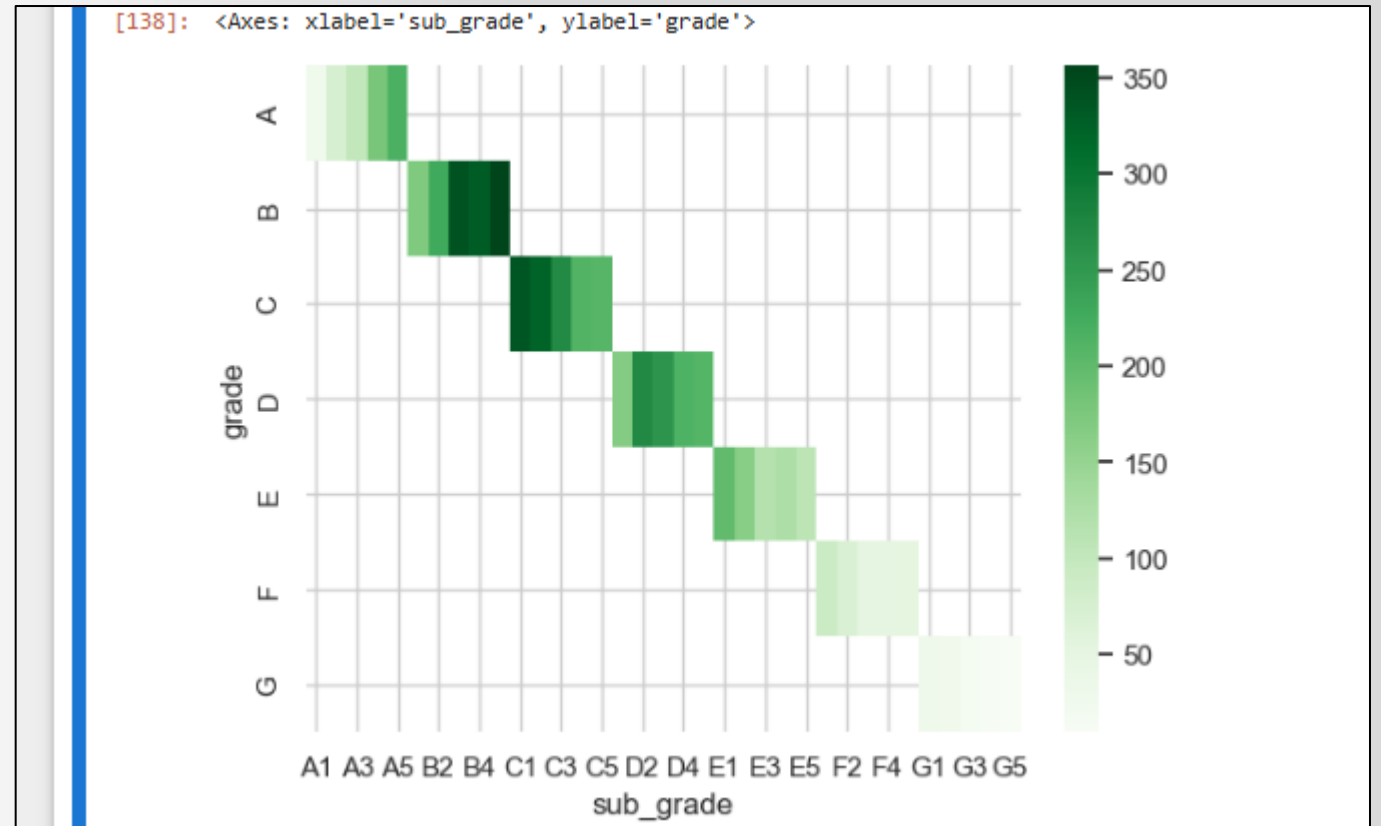
Bivariate – Numeric correlation

Pub_rec vs Pub_rec_bankruptcies, has linear progression relationship with each other, we shall ignore Pub_rec_bankruptcies from driver variables



Bivariate - categorical

- Grade, sub_grade – grade can be dropped
- Output is logical, because sub_grade is component of grade



Inferences/Relationships from other Bivariate analysis

Relationships:

- Pub_rec and pub_rec_bankruptcies are related to each other, one can be dropped
- Similarly, Grade – sub_grade is a component of Grade; grade can be dropped

Few Inferences:

- int_rate vs term: Loans of 36 months with interest 10-14% seems to be risky
- grade vs term: Loans of 36 months with sub_Grade B3,B5 seems to be high risky
- verification_status vs term: Loans of 36 months not verified is risky
- purpose vs term: Debt Consolidation seems to be risky - highest for 36 months and next for 60 months

Recommendations & Conclusions

- We have identified 8 variables that drive the target variable (loan_status)
- Also, inferential relationships can be identified, which can be used in machine learning model later

Variable	Variable type
term	Ordered Categorical
int_rate	Ordered Categorical
sub_grade	Unordered Categorical
annual_inc	Numerical
verification_status	Unordered Categorical
purpose	Unordered Categorical
pub_rec	Numerical
revol_util	Numerical



Thank you