# Lending Club Case Study

Suresh Kumar Yatirajula

Sundar Shankar
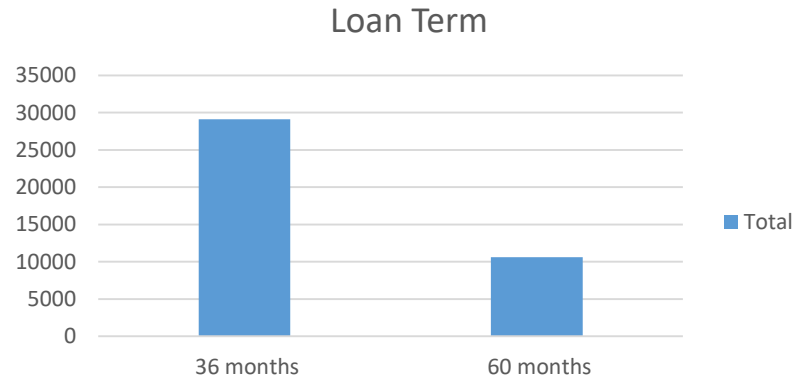
Srinivasa Marappa

Prateek Pande
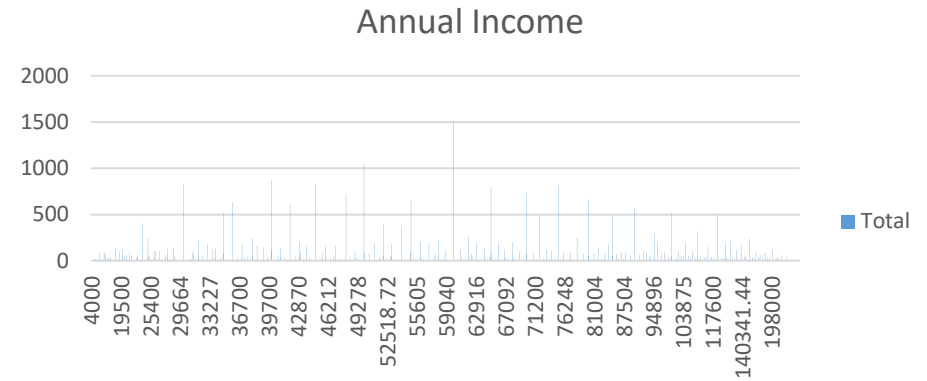
# Data Understanding

# Some Interesting Observations from Data
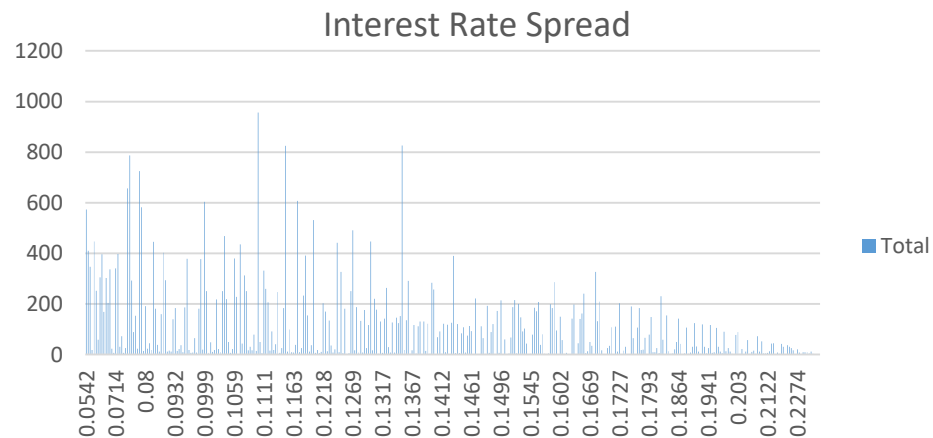
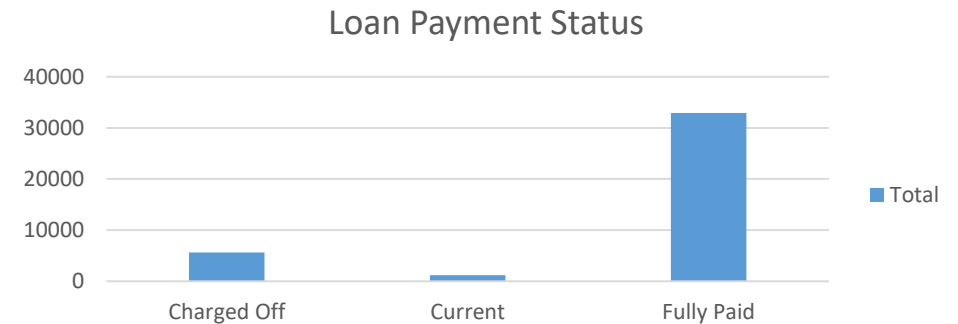- "Loan Term" is either 36 months or 60 months

### Loan Term


- Income of borrowers range from 4K to 6 million, with a mean salary around 68K

### Annual Income


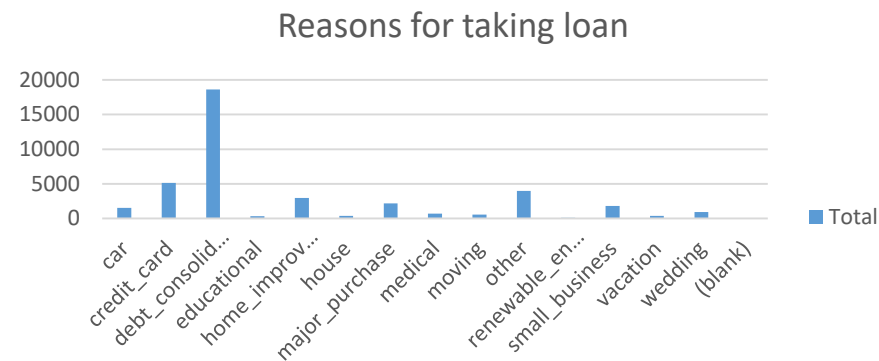- Minimum interest is 5.42% and max interest rate is 24.59%

### Interest Rate Spread


- Loan Delinquent Borrowers Percentage - ~ 14%

### Loan Payment Status

# Some Interesting Observations from data

- The most common reason for taking loan is: "Debt Consolidation"

### Reasons for taking loan



- The loans are categorized from A to Z and then further sub-graded from 1 to 5

### Loan Categories Sanctioned



- Charged Off Loans in Dollars is around ~15%

### Loan Status by Amount

# Data Understanding

- The input loan data file has 39,718 records and 111 columns

- Columns with all NULL values:
  - There are around 54 columns where all values are NULL/NaN

    'mths_since_last_major_derog', 'annual_inc_joint', 'dti_joint', 'verification_status_joint', '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', 'il_util', 'open_rv_12m',
    'open_rv_24m', 'max_bal_bc', 'all_util', 'total_rev_hi_lim', 'inq_fi', 'total_cu_tl', 'inq_last_12m', 'acc_open_past_24mths',
    'avg_cur_bal', 'bc_open_to_buy', 'bc_util', 'mo_sin_old_il_acct', 'mo_sin_old_rev_tl_op','mo_sin_rcnt_rev_tl_op', 'mo_sin_rcnt_tl',
    'mort_acc', 'mths_since_recent_bc', 'mths_since_recent_bc_dlq', 'mths_since_recent_inq', 'mths_since_recent_revol_delinq',
    'num_accts_ever_120_pd', 'num_actv_bc_tl', 'num_actv_rev_tl', 'num_bc_sats', 'num_bc_tl', 'num_il_tl', 'num_op_rev_tl',
    'num_rev_accts', 'num_rev_tl_bal_gt_0', 'num_sats', 'num_tl_120dpd_2m', 'num_tl_30dpd', 'num_tl_90g_dpd_24m', 'num_tl_op_past_12m',
    'pct_tl_nvr_dlq', 'percent_bc_gt_75', 'tot_hi_cred_lim', 'total_bal_ex_mort', 'total_bc_limit', 'total_il_high_credit_limit'

- Columns with Majority NULL values:
  - Out of remaining 57 columns, there are 3 columns, where more than 90% values are NULL

    'mths_since_last_delinq','mths_since_last_record','next_pymnt_d'

- Columns with only 1 Distinct Value:
  - There are 6 columns where there is only a one distinct value

    ['pymnt_plan', 'initial_list_status', 'policy_code', 'application_type', 'acc_now_delinq', 'delinq_amnt']

# Data Understanding

- ## Columns with NULL values in more than 90% values
  - Few columns have more than 90% of data as NULL. These columns are not very useful in the analysis

  ```
  Index(['pymnt_plan', 'initial_list_status', 'policy_code', 'application_type', 'acc_now_delinq',
  'delinq_amnt'])
  ```

- ## Columns with either 0 or NaN values:
  - Some columns have only either 0 or NaN as values. These columns are not very useful in the analysis.

  ```
  'chargeoff_within_12_mths', 'pub_rec_bankruptcies', 'tax_liens'
  ```

- ## Columns with Missing Values:
  - There are few columns with missing values. These columns needs to be filled for further analysis

- ## Categorical data which can be converted to Numerical Data:
  - There are some categorical columns like "Interest Rate percentage", "loan term", "employment length" etc… which can be converted to numerical fields for better analysis.

# Data Cleaning and Manipulation

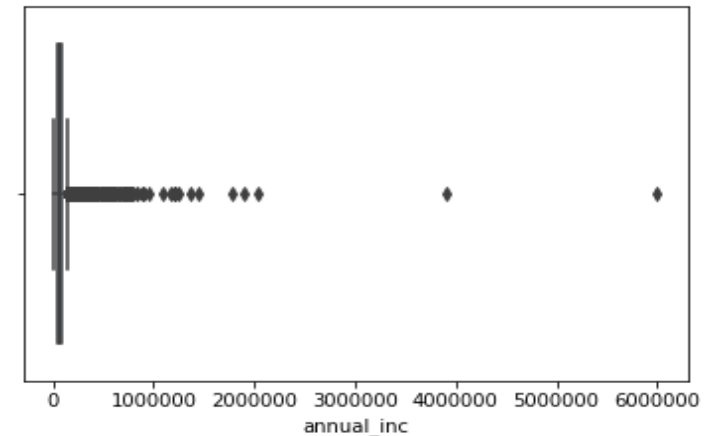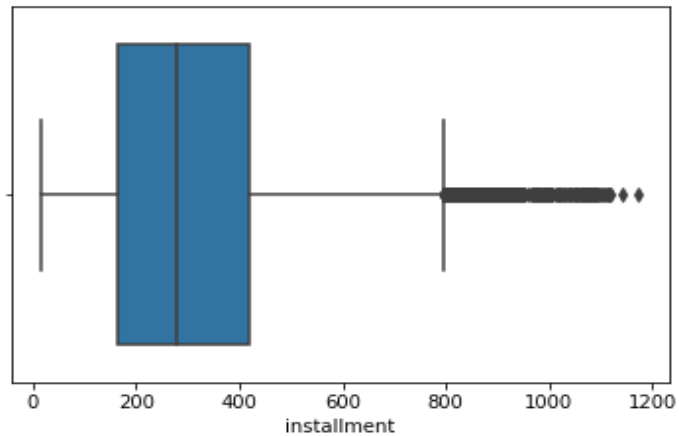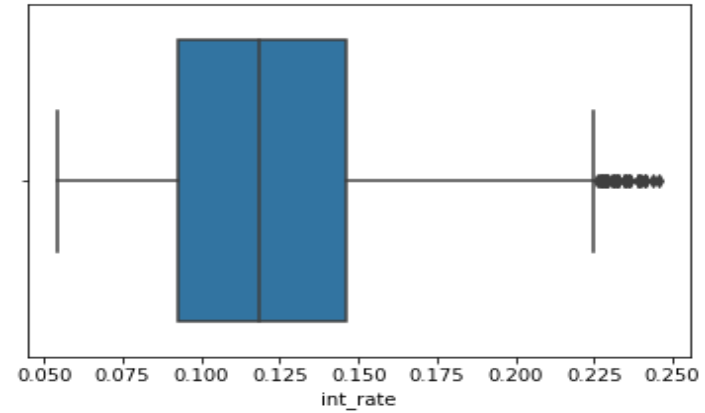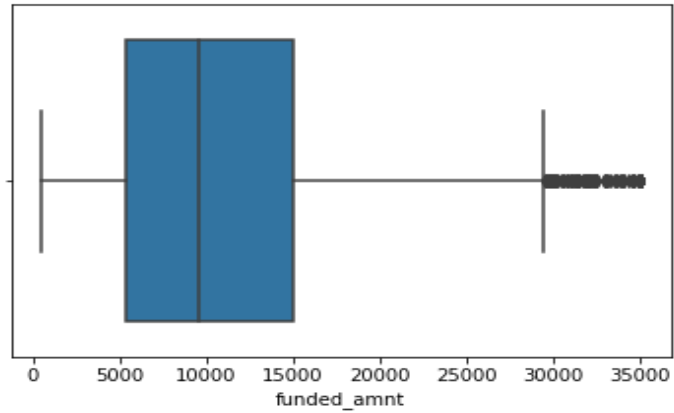# Data Cleaning and Manipulation

- The following cleaning/Manipulation steps are done on the input data:

  - The Columns following properties are <u>dropped</u> :
    - All values as Null or Empty
    - only 1 Distinct Value
    - Values either NaN or 0
    - Less than 10% of valid values

  - The following columns that have redundant or unwanted data are <u>dropped</u>:
    - <u>URL:</u> The information available in url are already part of other columns and hence redundant
    - <u>loan_amt</u>: Similar information is available in column "funded_amt" and doesn't seem to add any value
    - <u>funded_amt_inv, out_prncp_inv , total_pymnt_inv :</u> These columns are strongly correlated to "funded_amt" and doesn't seem to add any value

  - The following <u>categorical</u> columns are converted to <u>Numerical</u> columns for better analysis
    - <u>Term:</u> remove the string "months" from "36 months" & "60 months". So, the values of the column now become numerical 36 & 60.
    - <u>int_rate:</u> Remove the symbol "%" from the values.
    - <u>revol_util:</u> Remove the symbol "%" from the values.
    - <u>emp_length:</u> Remove the string "Years" and special characters like "<" & "+".

## Data Cleaning and Manipulation

- Format and Convert Data Type of Date Fields:
  - The columns issue_d, last_pymnt_d, last_credit_pull_d & earliest_cr_line are formatted and converted to Date data type

- Fill missing values:
  - last_credit_pull_d: whenever values of this date field are missing, then the values from issue_d are taken
  - pub_rec_bankruptcies: Missing values of this field are filled with 0
  - emp_title: This is a string data type column and hence missing values are filled with "NA"

- Remove Outliers:
  - Outlier are identified and removed for the following numerical columns:
    - funded_amnt
    - int_rate
    - Installment
    - int_rate
    - annual_inc
    - delinq_2yrs
    - inq_last_6mths
    - open_acc
    - pub_rec
    - revol_bal

# Data Cleaning and Manipulation

- Example Outliers Box Plots for "funded_amt","int_rate","installment","annual_inc":

# Data Analysis

# Data Analysis

After cleaning and manipulating data, the number of columns have comedown from 111 to 37. The data analysis is done on this cleaned data and some important steps and results are presented below.

- Finding Numerical and Categorical Features:
  - There are 23 numerical and 14 categorical features in the data:

```
Numerical Fields:
Index(['funded_amnt', 'term', 'int_rate', 'installment', 'emp_length', 'annual_inc', 'dti', 'delinq_2yrs',
'inq_last_6mths', 'open_acc', 'pub_rec', 'revol_bal', 'revol_util', 'total_acc', 'out_prncp', 'total_pymnt',
'total_rec_prncp', 'total_rec_int', 'total_rec_late_fee', 'recoveries', 'collection_recovery_fee', 'last_pymnt_amnt',
'pub_rec_bankruptcies'], dtype='object')

Categorical Fields:
Index(['id', 'grade', 'sub_grade', 'emp_title', 'home_ownership', 'verification_status', 'issue_d', 'loan_status',
'purpose', 'title', 'addr_state', 'earliest_cr_line', 'last_pymnt_d', 'last_credit_pull_d'], dtype='object')
```
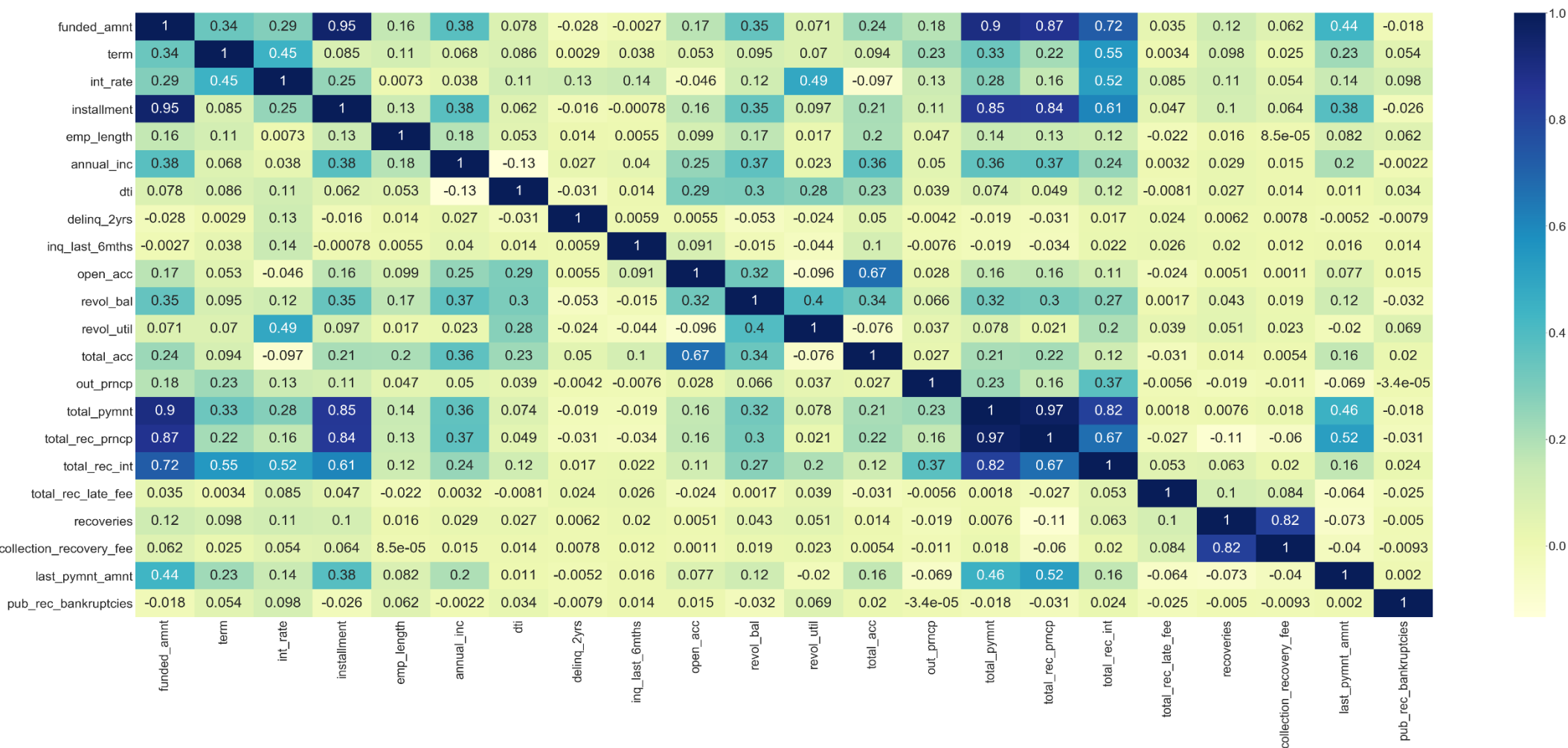
This data helps in understanding the type of variables and possible analysis like uni-variant/bi-variant analysis on them. This data also helps in determining the columns for correlation analysis.

- Analysis: Further analysis is done on various column combinations to find and arrive at the correlation between different variables and results are presented in following slides:
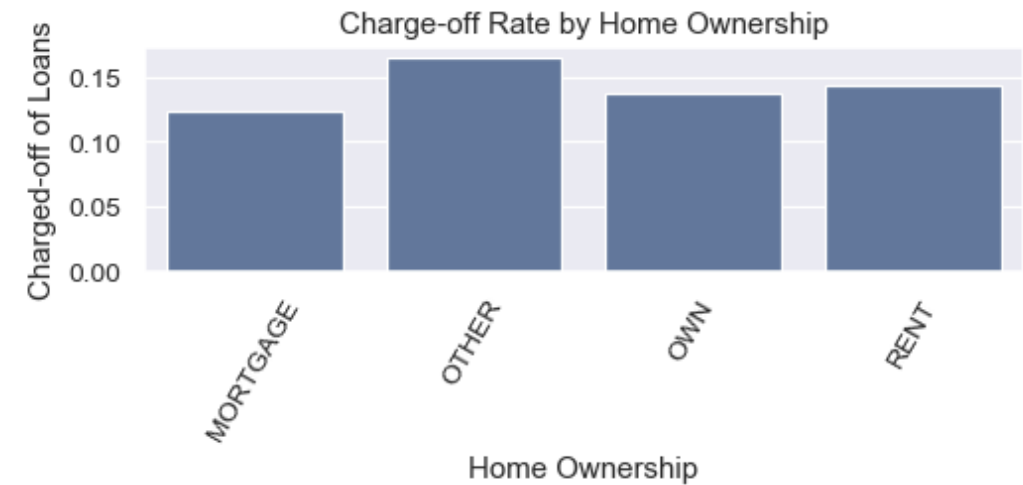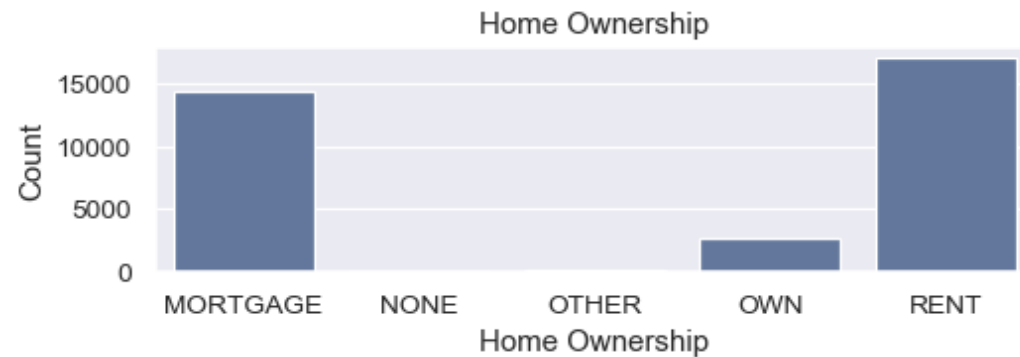
# Data Analysis
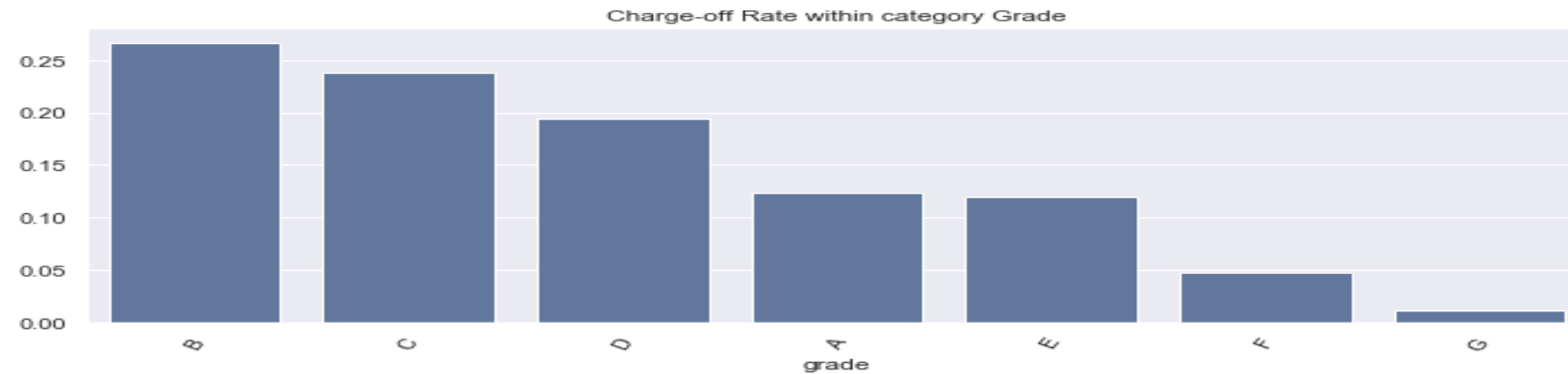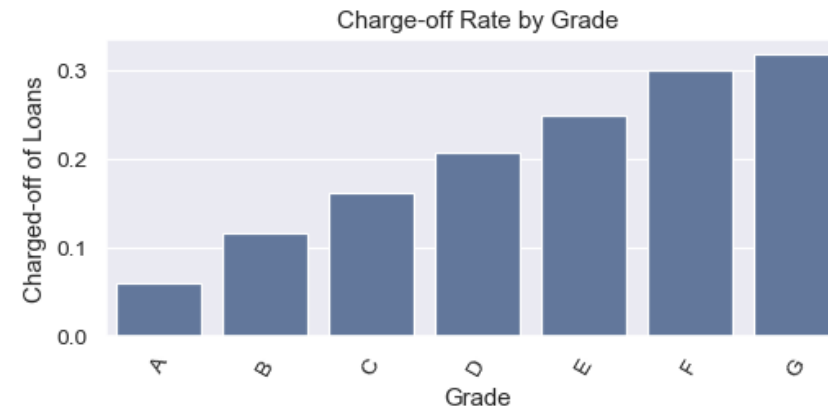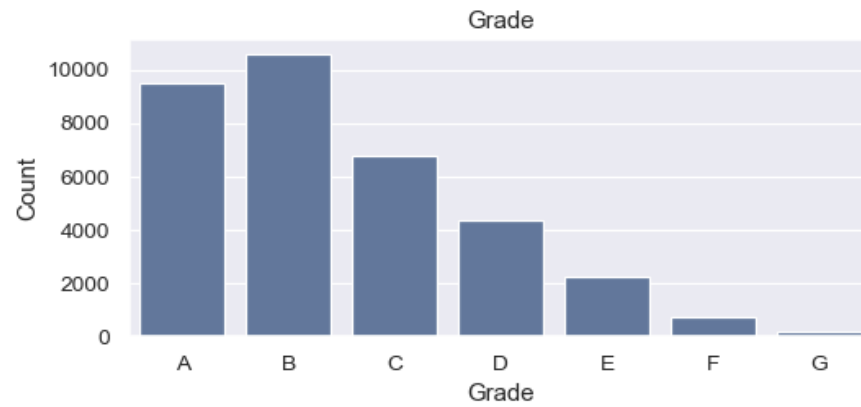
## Correlation plot for numerical data

# Data Analysis

- Uni-Variant and Bi-Variant analysis is done on different columns to analyze the data and their relationship
- Some important plots of various columns in relationship with loan default/charge-off
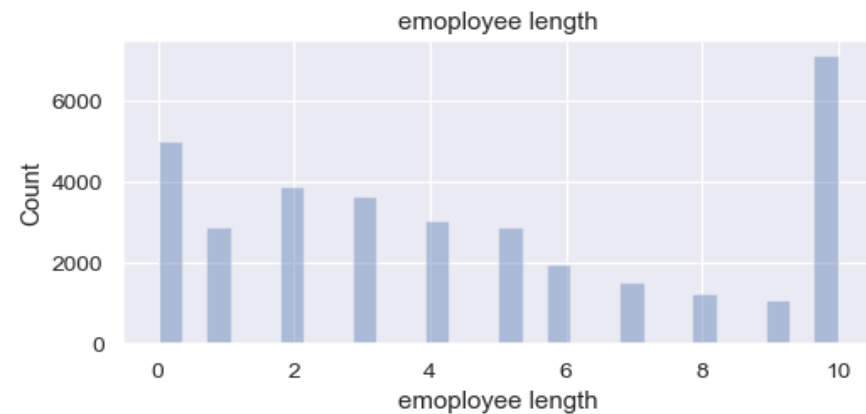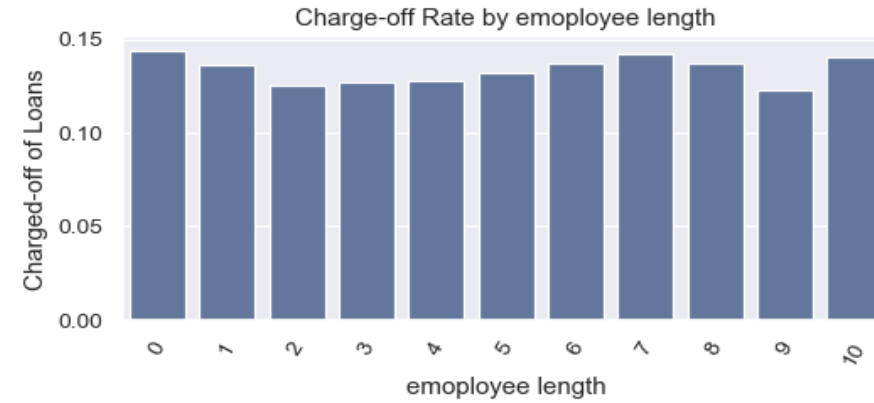
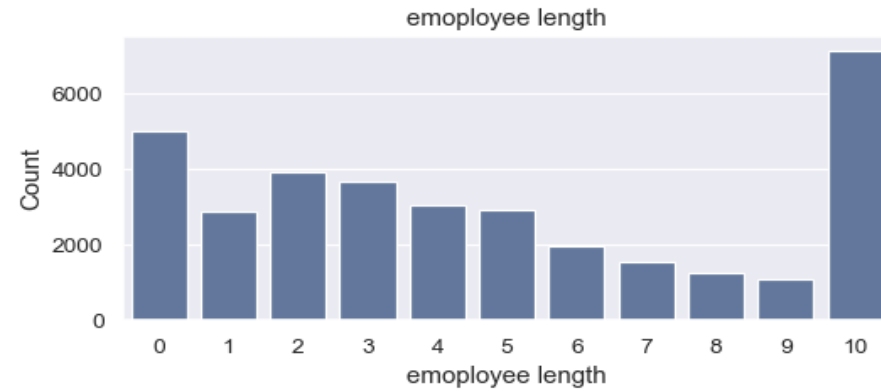- Plot: Charge-off rate by Home Ownership

# Data Analysis

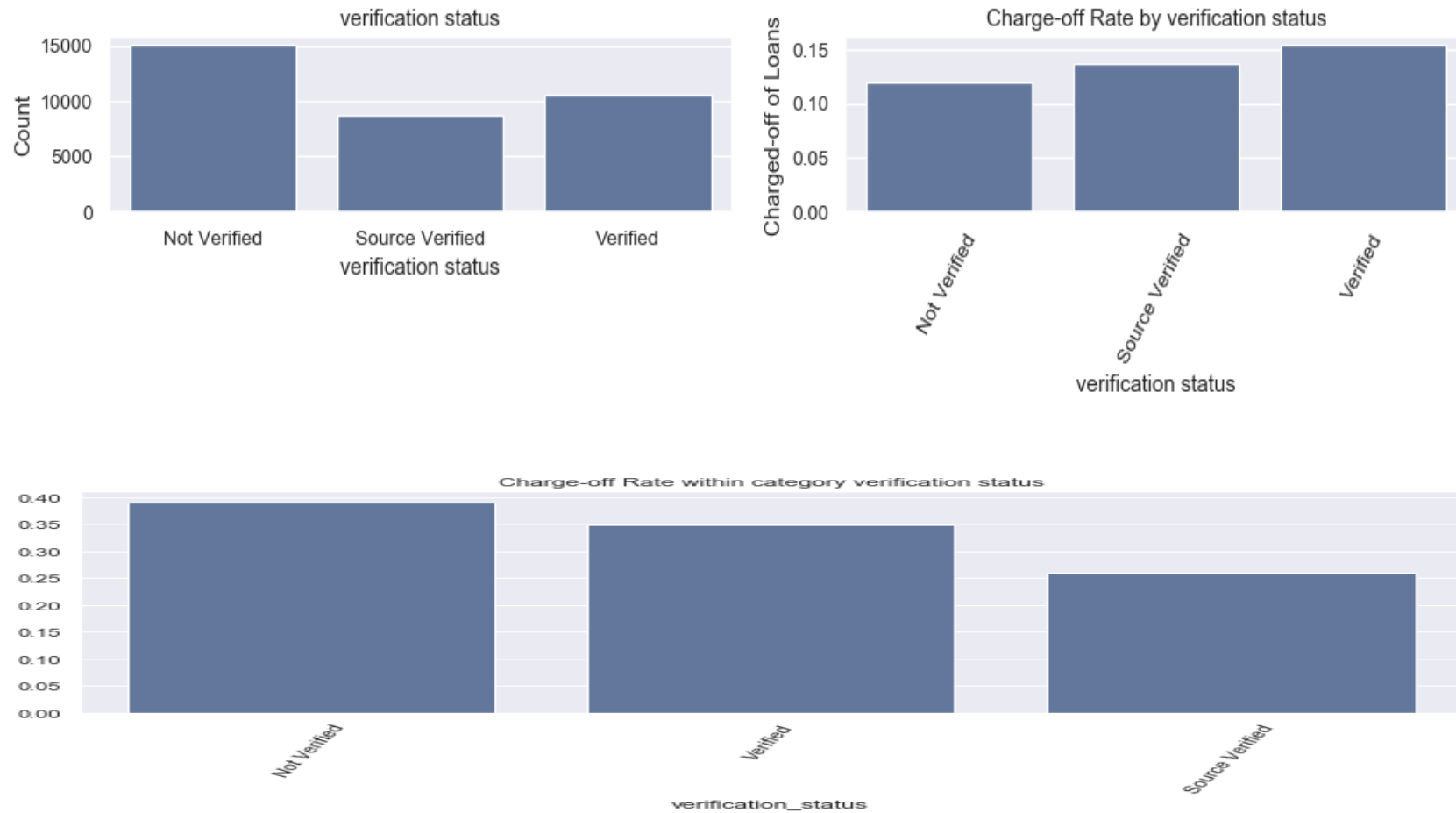- Plot: Charge-off rate by "Loan Grade"

# Data Analysis

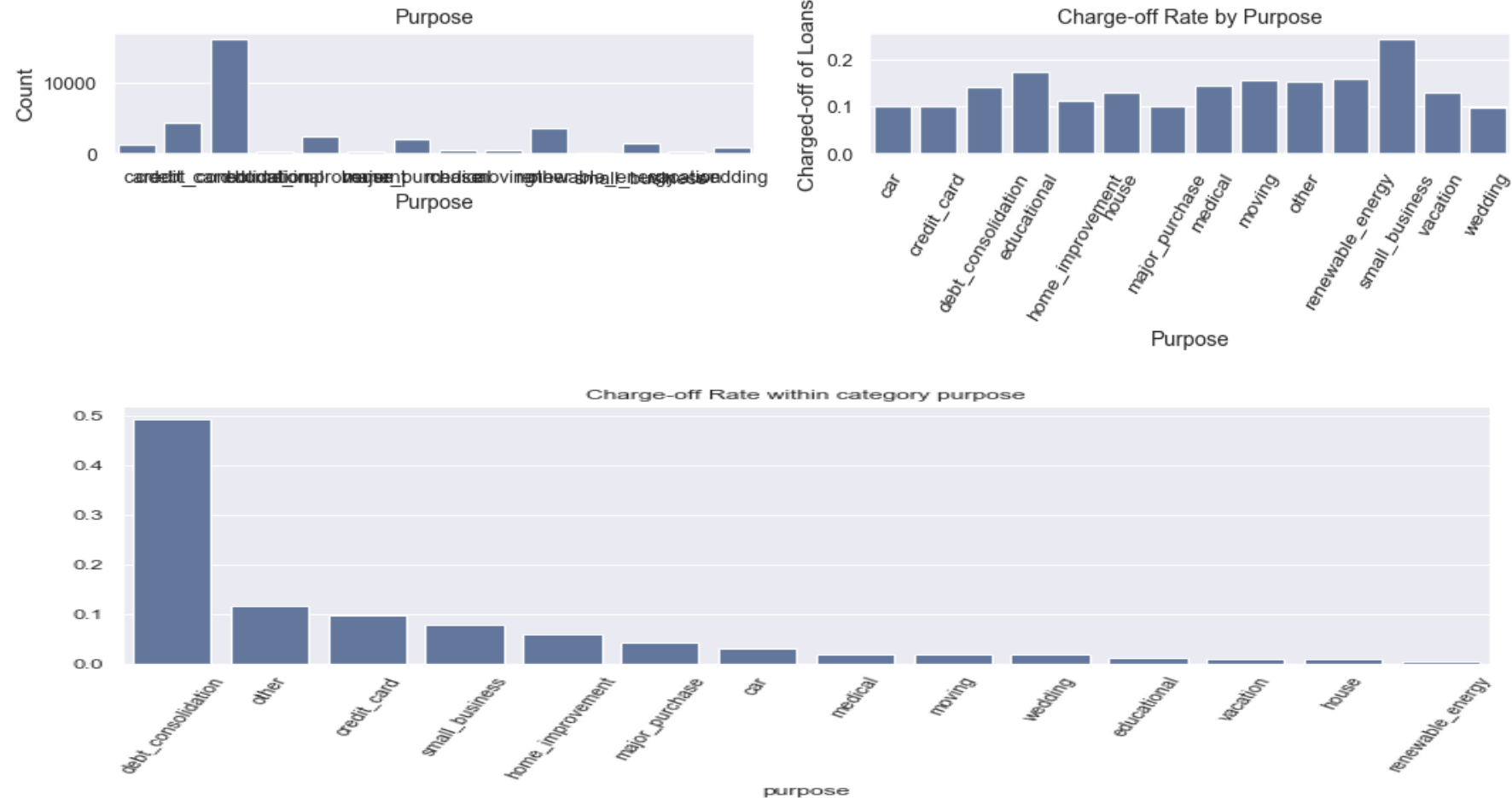- <u>Plot: Charge-off rate by "Employment Length"</u>

# Data Analysis

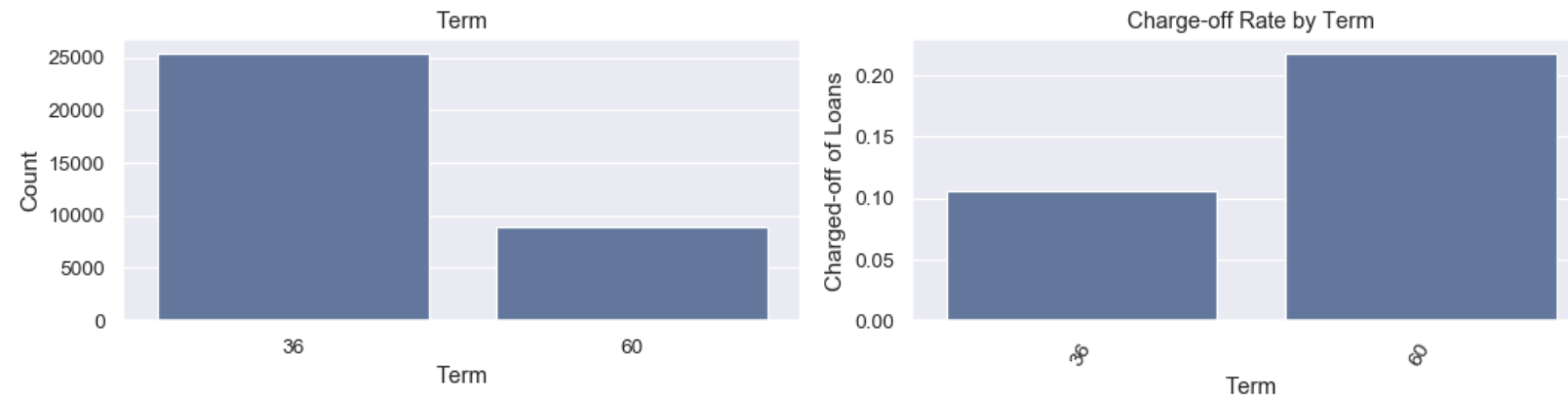- <u>Plot: Charge-off rate by "Verification Status"</u>

# Data Analysis

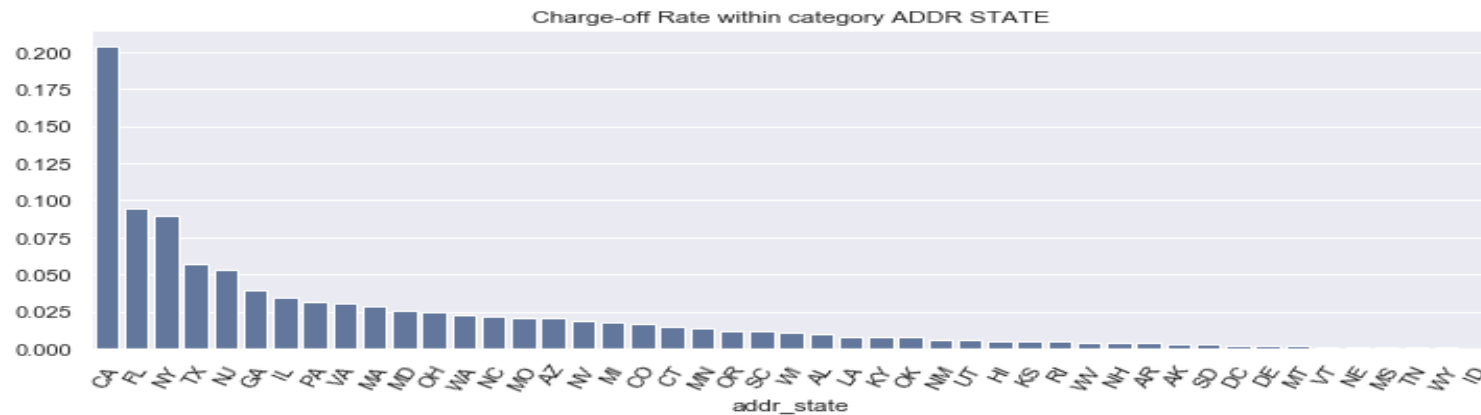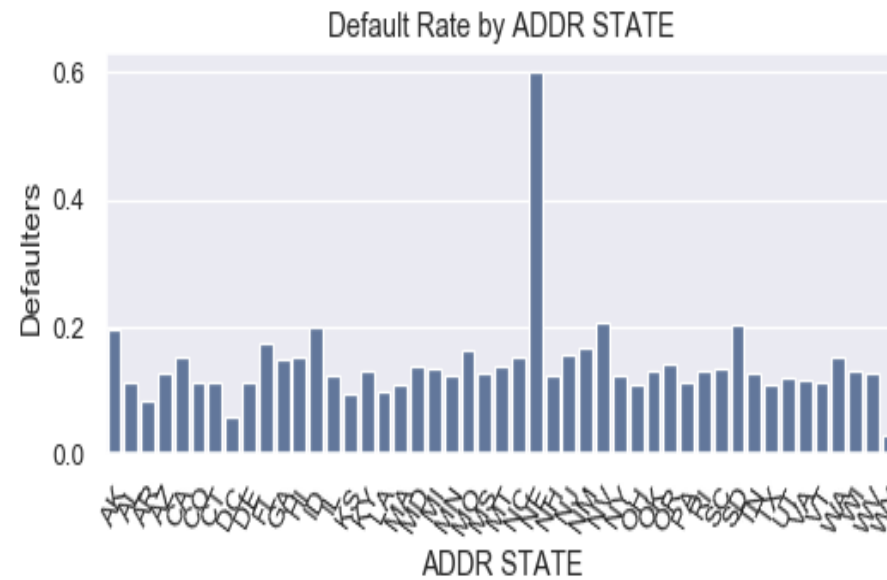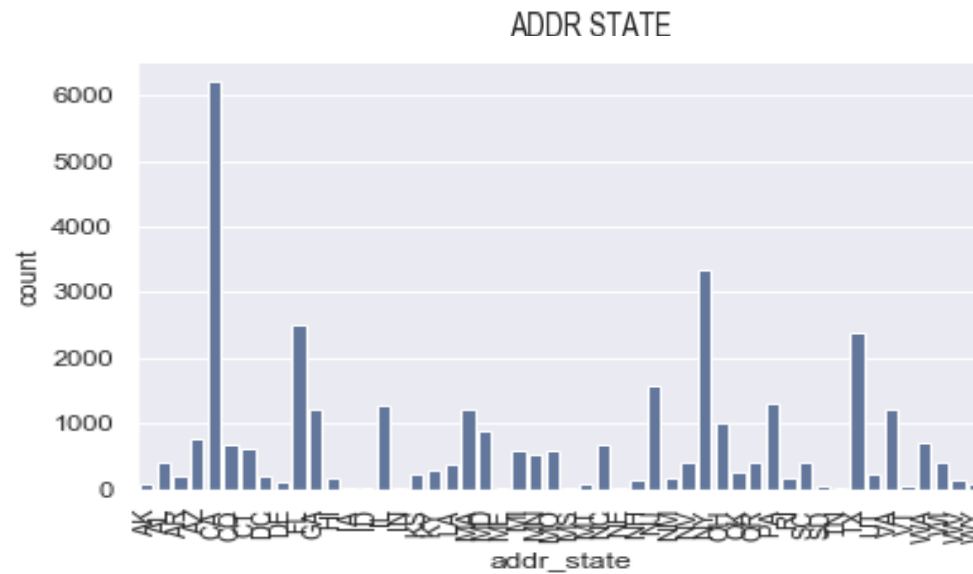- <u>Plot: Charge-off rate by "Purpose of Loan"</u>

# Data Analysis

- <u>Plot: Charge-off rate by "Loan Term"</u>

# Data Analysis

- Plot: Charge-off rate by "State"

# Data Analysis

- <u>Annova Analysis on Loan</u>

| | feature | f | p |
|---|---|---|---|
| 1 | grade | 507.23 | 0.00 |
| 2 | sub_grade | 105.42 | 0.00 |
| 5 | verification_status | 3401.19 | 0.00 |
| 8 | purpose | 257.75 | 0.00 |
| 14 | funded_amnt_cat | 118824.05 | 0.00 |
| 15 | annual_inc_cat | 286.03 | 0.00 |
| 11 | earliest_cr_line | 4.79 | 0.00 |
| 12 | last_pymnt_d | 14.74 | 0.00 |
| 4 | home_ownership | 209.47 | 0.00 |
| 6 | issue_d | 17.62 | 0.00 |
| 7 | loan_status | 349.89 | 0.00 |
| 13 | last_credit_pull_d | 3.86 | 0.00 |
| 9 | title | 1.15 | 0.00 |
| 10 | addr_state | 2.13 | 0.00 |
| 3 | emp_title | 1.03 | 0.05 |
| 0 | id | 0.00 | nan |

# Results & Recommendations

# Inferences, Results & Recommendations

Based on the above analysis:

Top inferences and driving factors behind loan defaults are marked in blue color:

Loan Provider Factors:
- Median of funded amount is higher for current loan than fully paid loan or charged off loan.
- Loans are not given for DTI more than 30.
- Also very less loans are provided for DTI higher than 25.

Loan Consumer Factors:
- People living in rent or mortgage tend to take more loans.
- **Percentage of people who are living in rent have highest defaults. After that people with mortgage and least are people with own house or others.**
- **Number of default is more in small business.**
- Maximum number of loans is taken from California.
- **Also highest default rate is from CA.**
- **As funding amount increases above 20K risk of default increases.**

Loan Consumer Behavior Factors -
- **Frequency of people getting loans is higher for higher grade. Also people with better grades have lower defaults.**
- People with 0 or 10+ years of experience take more number of loans
- **Within defaulters people with 0 or 10+ years experience tends to default more.**
- **Default rate is highest in loan for the purpose of debt-consolidation.**
- Maximum loans are availed by people with income below 100,000.
- As the income increase default rates are coming down.
- More number of people have taken loan with less term of 36 month.
- **Higher defaulters are found in long term loan of 60 months.**