# Stage C Report: Data Preparation and Model Selection

This report presents the progress made in the Data Preparation and Modeling step of our project, building on insights from the previous Data Ingestion and EDA report. We outline the data processing steps, considerations for selecting the appropriate model, and approach for evaluating multiple models. Findings are detailed in the report, with technical details available in the attached code file (stepC.py).

## Documenting Completed Work Assignments-

Our approach challenges the conventional data division into two types, acknowledging its limitations in capturing real-world complexity. Instead, we adopt a nuanced classification of seven types: useless, nominal, binary, serial, count, time, and interval. This empowers us to understand feature diversity and apply specialized techniques accordingly, capturing vital nuances and avoiding information loss. In this step, our goal is to optimize future feature selection. We exclude useless columns and focus on rigorously tested features. By preparing the remaining features for direct model integration, we minimize future data preparation efforts. This enables us to prioritize model refinement and insightful analysis, freeing us from laborious data preprocessing tasks.

1. **Progress in the data cleaning process as indicated in the previous report:**
   a. **Selecting Relevant Loans:** The raw data contains 434,407 instances that decrease in the previous step to 328,956 instances, with 99,404 instances removed due to loan status other than "Charged Off" and "Fully Paid," and 6,047 joint loan applications removed to mitigate bias (Refer to stage-B report for more information). At this stage, we removed 42 instances from our analysis due to unusual changes observed in the int_rate and installment columns between 2018 and 2019 snapshots. Considering the unlikely possibility of actual changes in these columns, we attributed these discrepancies to potential human error, which could potentially compromise the reliability and accuracy of our analysis. Therefore, we ultimately analyzed 328,914 instances.
   b. **Feature selection:** In this step, we conducted further investigation on the initial selection of features from the raw data containing 151 features, based on the following considerations (See attached Excel file for column feature updates):
      i. Features exhibiting leakage-We removed 43 features that identified as potential leakage features (refer to stage-B report for further details).
      ii. Missing values (NA's) per feature- We removed 23 features with more than 45% missing values (refer to stage-B report for further details).
      iii. Features business relevant- We removed 7 features that didn't align with our business understanding (see stage-B report for further details).
      iv. The distribution of features values- In addition to our previous assessment of entropy and low frequency values, we performed tests for standard deviation and high frequency values to assess attribute value distribution in this stage due this we removed 8 features. The reason we performed these additional tests was to identify highly scattered/not scattered values and values with a high frequency in one value, respectively. We made decisions on whether to remove or modify features

based on their relevance to the problem and the distribution by target variable (refer to Appendix A for further details for each column and the report of Step B).

2. **Handling missing data:** We addressed missing data by employing methods such as linear regression and clustering, and used statistical measures to impute some of the missing values. To avoid compromising the true distribution of each feature, we used two completion methods for each feature. Details and corresponding graphs for each column can be found in Appendix B. We plan to evaluate the effectiveness of these methods by using F1, r^2, and other metrics on our model in the following weeks.

3. <u>**Selection of models:**</u> In the previous step, we computed two target variables: loan yield (continuous) and yield above 2% (binary). To understand p2p loan returns comprehensively, we selected a regression model for quantitative yield estimation and a classification model to identify loans meeting the target return threshold (2%). These models will evaluate the profitability and risk of integrating p2p loans into investment portfolios. Over the next two weeks, we will test various models: K-Nearest Neighbors, Decision Tree, Random Forest, Support Vector Machines, and Logistic Regression for classification, and Linear Regression, Polynomial Regression, Decision Tree Regression, and Random Forest Regression for regression. In 4 weeks, following the performance evaluation outlined in the last paragraph of this report, we will select the optimal model for implementation.

4. <u>**Handling Outliers:**</u> We examined the distribution of variables to identify and handle outliers. Histogram and QQ plots were used to understand the distribution of each feature. For normally distributed variables, we employed the IQR method and considered values within three standard deviations from the mean as non-outliers. For non-normal distributions, we explored techniques such as log transformation and features and transformation. (refer to Appendix E for further details for each column)

5. <u>**Correlation between features:**</u> We considered data type, distribution, and variable type to choose the appropriate correlation method. For continuous features, the assumption of normal distribution and linearity associated with Pearson correlation coefficient was evaluated. However, as our numeric continuous features did not follow a normal distribution, we opted for Spearman's rank correlation coefficient, a non-parametric test. A correlation heatmap was created to visualize correlations, a threshold of 0.9 correlation score which indicates a strong linear relationship between the variables was chosen to remove correlated variables and scatter plots helped examine the highly correlated variables. To manage multicollinearity, two variables were removed and a new feature was created from correlated features. If any more features need to be removed, a lower threshold will be chosen. For categorical variables, we used the chi-square test to assess associations and uncover meaningful patterns and dependencies. (refer to Appendix F for further details for each column)

6. <u>**Feature Engineering:**</u> We aimed to improve the predictive power of the model by generating new features through various transformations applied to existing features. Our goal was to identify informative features while preventing overfitting. The need for this process arose due to low frequency of values, outliers, business understanding, high correlations between variables, and handling categorical variables etc.. The Phase B report details the entire process, with a list of new features created in Appendix C. In total, we generated 12 new features.

7. **Handling Categorical Features:** Continuing from the previous step, we addressed the issue of handling categorical variables in our model. To simplify the model, we adopted a general approach that involved reducing the dimensions of the categorical variables. We used several techniques such as category transformation based on their relevance to the target variable, one-hot encoding, and hashing. For a detailed explanation of each variable's treatment and corresponding graphs, please refer to Appendix D and the report of Step B.

8. **Standardization:** We applied min-max normalization to standardize the parameters and account for handled outliers. This method scales parameter values between 0 and 1, maintaining the relative relationships between data points. It ensures equal treatment and consistent scaling of all parameters, which is crucial for machine learning algorithms. We will continue testing and may make adjustments as needed due to sensitivity to outliers for this method.

**Potential Pitfalls:**
- Incorrect data provided by the borrower can lead to biased results.
- Inaccurate imputation methods that may introduce bias or distort the true distribution of features.
- Insufficient evaluation of missing value completion methods.
- Limited evaluation of models using relevant metrics and benchmarks.
- Inappropriate assumptions about variable distributions leading to ineffective outlier handling.
- Limited exploration of alternative outlier detection and treatment techniques.
- Insufficient assessment of multicollinearity's impact on model stability and interpretability.
- Overcomplicating the model with excessive new features without proper justification.
- Failure to address potential overfitting issues caused by new feature introduction.
- Incorrect type of transformation leading to incorrect or meaningless results.
- Incorrect encoding of categorical variables can lead to incorrect assumptions about their relationship with the target variable.
- Arbitrary dimensionality reduction methods that may result in information loss.
- Neglecting the impact of outliers on the chosen standardization method's effectiveness.

Our next steps involve training, optimizing, and evaluating the model's performance using metrics like accuracy, precision, recall, and F1-score. We will fine-tune the model through techniques like hyperparameter tuning and cross-validation to optimize performance and assess robustness. Through iterative improvement, we will refine the model by exploring different approaches, algorithms, or parameters. We will compare the performance of the regression and classification models, considering business and performance perspectives. Additionally, we will compare our model with the existing loan grade model, using methods like forecasting the probability of returns above 2% and divide them into 7 categories (A-G), comparing the percentage of loans predicted correctly in each grade. We will use Python to perform these operations. Anticipated timeframe for clear results is 5 weeks.

Best regards,
The DataDriven Portfolio Solutions Management Team

# Appendices

## Appendix A:

| Column Name | Treatment | |
|---|---|---|
| **delinq_amnt** | This variable contains 99.45% of the data in one category, where both the failure and return rates are close to the overall average. Therefore, we will remove this variable and thus prevent potential biases or inaccuracies in our analysis that could arise from including a variable with minimal variation and limited impact on the target variable. |  |
| **collections_12_mths_ex_med** | This variable contains 98% of the data in one category, where both the failure and return rates are close to the overall average. The remaining 2% of the data is dispersed and does not provide conclusive insights. Therefore, we will remove this variable and thus prevent potential biases or inaccuracies in our analysis that could arise from including a variable with minimal variation and limited impact on the target variable. |  |

The table for delinq_amnt shows:

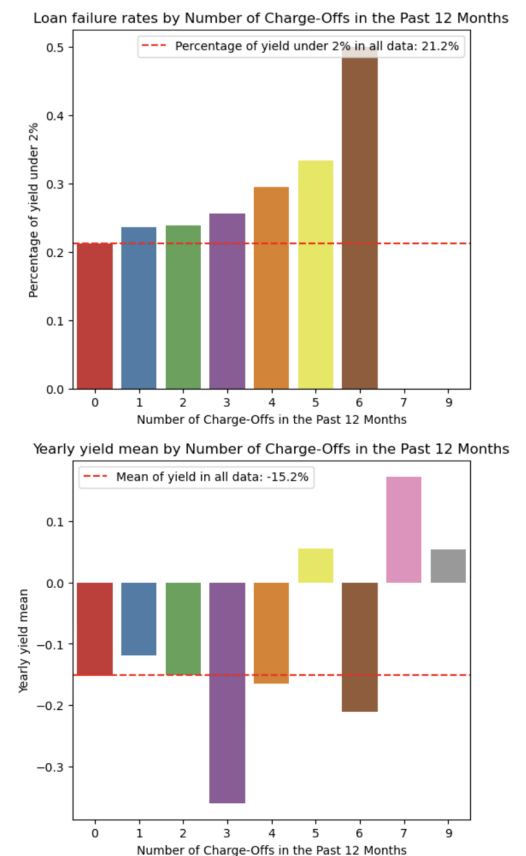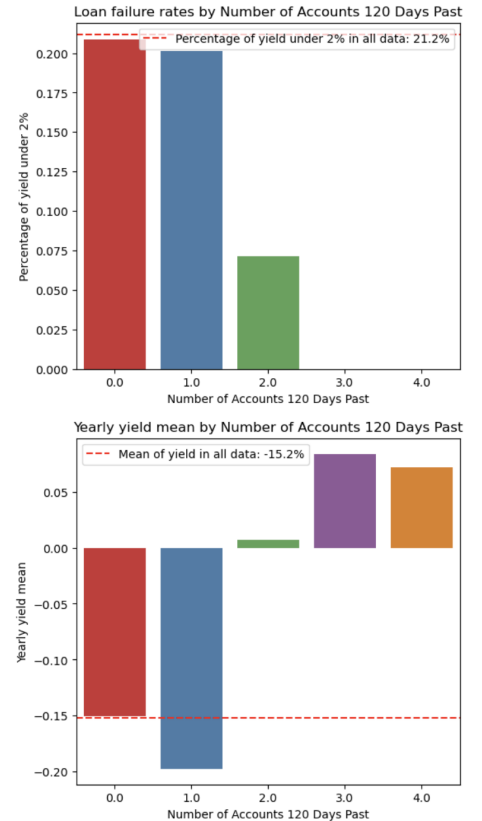| | delinq_amnt | Count | normalize count | yield over percentages | yearly yield mean |
|---|---|---|---|---|---|
| 0 | 0 | 326746 | 0.994546 | 0.211804 | -0.151707 |
| 1 | 65000 | 40 | 0.000122 | 0.325000 | -0.310498 |
| 2 | 25 | 24 | 0.000073 | 0.208333 | -0.005655 |
| 3 | 56 | 19 | 0.000058 | 0.263158 | -0.242170 |
| 4 | 30 | 19 | 0.000058 | 0.263158 | -0.625535 |
| ... | ... | ... | ... | ... | ... |

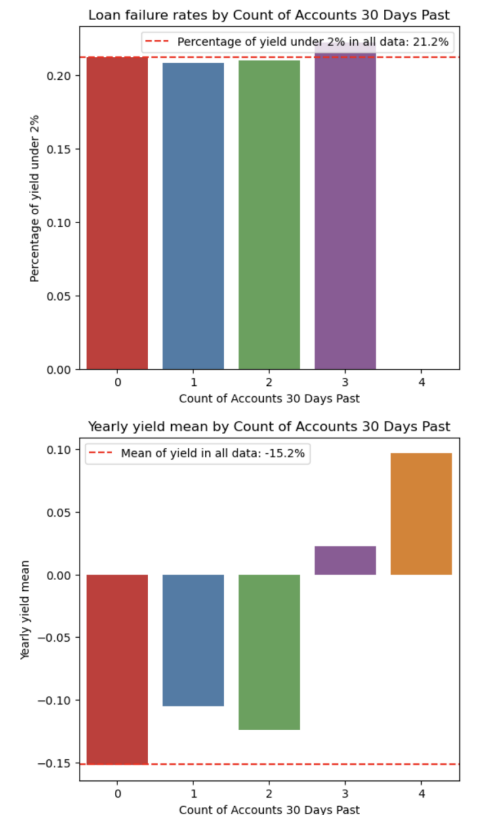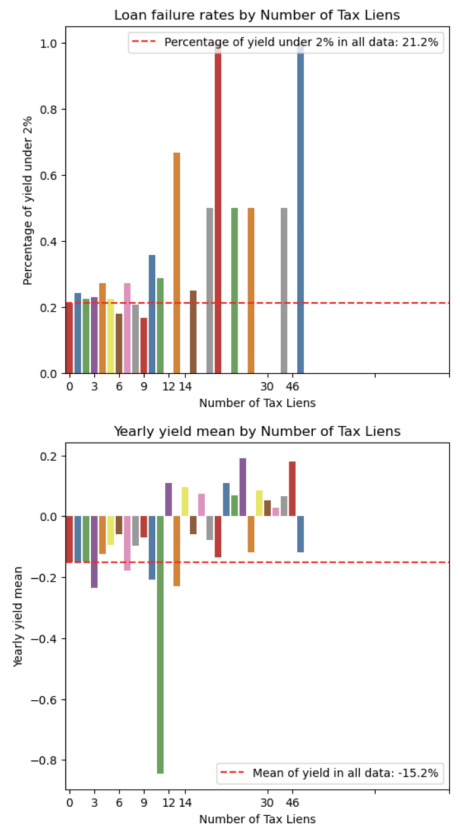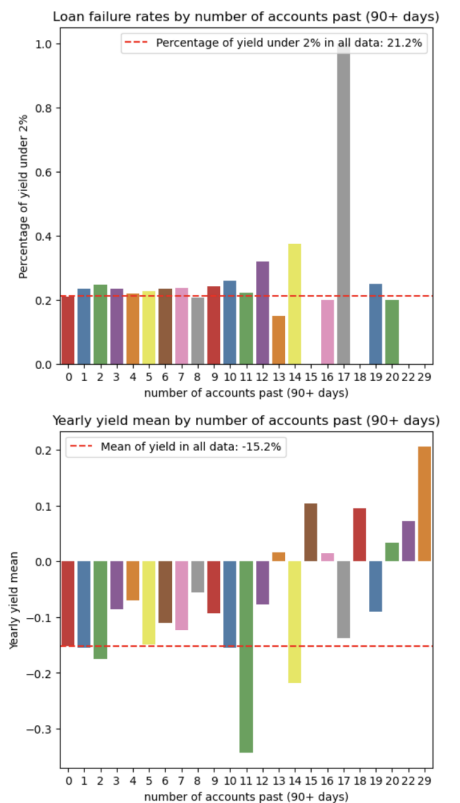| | | |
|---|---|---|
| **acc_now_delinq** | This variable contains 99.36% of the data in one category, where both the failure and return rates are close to the overall average. Therefore, we will remove this variable and thus prevent potential biases or inaccuracies in our analysis that could arise from including a variable with minimal variation and limited impact on the target variable. |  |
| **chargeoff_within_12_mths** | This variable contains 99.17% of the data in one category, where both the failure and return rates are close to the overall average. Therefore, we will remove this variable and thus prevent potential biases or inaccuracies in our analysis that could arise from including a variable with minimal variation and limited impact on the target variable. |  |

| **num_tl_120dpd_2m** | This variable contains 99.89% of the data in one category, where both the failure and return rates are close to the overall average. Therefore, we will remove this variable and thus prevent potential biases or inaccuracies in our analysis that could arise from including a variable with minimal variation and limited impact on the target variable. |  |
|---|---|---|
| **num_tl_30dpd** | This variable contains 99.58% of the data in one category, where both the failure and return rates are close to the overall average. Therefore, we will remove this variable and thus prevent potential biases or inaccuracies in our analysis that could arise from including a variable with minimal variation and limited impact on the target variable. |  |

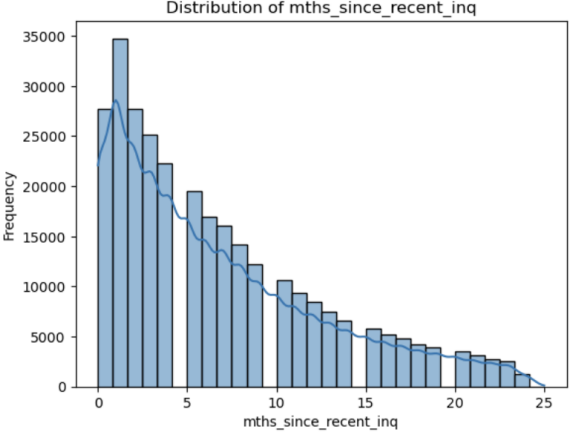| | | |
|---|---|---|
| **tax_liens** | In this variable, 95% of the values are 0 tax liens that have similar rates of failure and average returns as the population, while higher values indicate a higher risk of default with lower average returns. This variable is relevant for predicting loan defaults because it indicates a borrower's history of tax liens, which may suggest difficulty in paying taxes and a higher risk of default. We converted to a binary variable where 0 represents no tax liens and 1 represents one or more tax liens. |  |
| **num_tl_90g_dpd_24m** | This variable has low entropy but is relevant for predicting loan defaults because a high number of accounts past due may indicate a history of not repaying debts on time, which increases the risk of default. While high numbers of accounts past due generally lead to higher failure rates, for certain cases with a high number of accounts past due, they do not represent a large percentage of the data. Therefore, we can make this variable binary, where 0 represents no accounts past due and 1 represents a history of having accounts past due for 90 or more days in the last 24 months. |  |

| pub_rec_bankruptcies | This variable has low entropy but is relevant for predicting loan defaults because having a history of bankruptcy may indicate that the borrower has experienced financial difficulties in the past and may be at a higher risk of defaulting on the loan. Therefore, higher values of "pub_rec_bankruptcies" may be associated with a higher risk of default. Therefore, we can make this variable binary, where 0 represents no public record bankruptcies and 1 represents a more then one public record bankruptcies |  |
|---|---|---|

**Appendix B:**

| Column Name | Missing Values | % Missing | Imputation Method | |
|---|---|---|---|---|
| il_util | 43357 | 13.2% | LinearRegression with ['all_util', 'open_il_12m', 'total_bal_il', 'open_il_24m', 'open_act_il'] The "il_util" feature reflects a borrower's account utilization, which is important in assessing ability to return a loan. We can use other related columns such as payment and credit utilization to predict this variable, we can see that they have high correlation. |  |

| | | | | |
|---|---|---|---|---|
| | | | We observed that the distribution of "il_util" has a median greater than the mean, indicating a non-normal distribution. Hence, we will use the median to fill in the missing values, which is a more conservative approach. A lower value of "il_util" suggests that the borrower is utilizing a smaller portion of their available credit, potentially reducing the risk of default and improving their creditworthiness. | <br>Distribution of il_util |
| mths_since_recent_inq | 32274 | 9.82% | LinearRegression with ['mo_sin_rcnt_tl', 'mo_sin_rcnt_rev_tl_op', 'fico_range_high', 'fico_range_low']<br>The feature reflects the number of months since the borrower's most recent installment accounts opened, which is important in assessing ability to return a loan. We can use other related columns , we can see that there is high correlation. |  |
| | | | The "mths_since_recent_inq" distribution shows a higher mean than the median, indicating some outliers in the data. To address missing values, we will use the median as a conservative approach. Higher values of "mths_since_rcnt_il" suggest that the borrower has not recently applied for new credit, which could reduce the risk of default. | <br>Distribution of mths_since_recent_inq |

| | | | | |
|---|---|---|---|---|
| emp_length | 21394 | 6.51% | We divided the annual_inc to bins and then generated all possible combinations of home ownership and income ranges, and computed the mode value of the available data for each combination. | |
| | | | As you can see, there is no clear distribution for the variable, but it seems that most of the values are 10, so we will use mode |  Distribution of emp_length |
| mths_since_rcn | 8843 | 2.69% | The mean is larger than the median. This indicates that there are a few very large values (outliers) that are pulling up the mean. To be more conservative in our approach, we will use the average to impute missing values. |  Distribution of mths_since_rcnt_il |
| | | | However, because the difference between the average and the median is relatively large due to the presence of outliers, we will also consider using the median to impute missing values to have a more robust and conservative approach. | |

| | | | | |
|---|---|---|---|---|
| mo_sin_old_il_acct | 8792 | 2.67% | The median is larger than the mean. This indicates that there are a few very large values (outliers) that are pulling up the mean. To be more conservative in our approach, we will use the median to impute missing values. |  |
| | | | However, because the difference between the average and the median is relatively large due to the presence of outliers, we will also consider using the mean to impute missing values to have a more robust and conservative approach. | |
| bc_util | 3707 | 1.12% | The mean is larger than the median. This indicates that there are a few very large values (outliers) that are pulling up the mean. To be more conservative in our approach, we will use the average to impute missing values. |  |
| | | | However, because the difference between the average and the median is relatively large due to the presence of outliers, we will also consider using the median to impute missing values to have a more robust and conservative approach. | |

| percent_bc_gt_75 | 3567 | 1.08% | Because of the unclear distribution, we will complete both by mean and median. |  Distribution of percent_bc_gt_75 |
|---|---|---|---|---|
| bc_open_to_buy | 3542 | 1.07% | Because the mean and the median are very far from each other due to many outliers we will use both the mean and the median. |  Distribution of bc_open_to_buy |
| mths_since_recent_bc | 3328 | 1.01% | Because the mean and the median are very far from each other due to many outliers we will use both the mean and the median. |  Distribution of mths_since_recent_bc |

**Appendix C:**

| Column Name | Transformation Details | | New Column Name |
|---|---|---|---|
| num_tl_90g_dpd_24m | Binary(0- no accounts, 1- a history of having accounts) | detailed in Appendix A | has_past_due_accounts_90g_24m |
| pub_rec_bankruptcies | Binary(0- no public record bankruptcies, 1- more than one public record bankruptcies) | detailed in Appendix A | has_pub_rec_bankruptcies |
| purpose | Binary(1- yield over 2% percentage greater than population mean,0- otherwise) | detailed in the step B report | purpose_danger |
| addr_state | Binary(1- yield over 2% percentage greater than population mean,0- otherwise) | detailed in the step B report | addr_state_danger |
| earliest_cr_line | data['issue_d'] - data["earliest_cr_line"] | detailed in the step B report | Age_of_credit_history |
| inq_last_6mths | Ordinal encoding (high-3,medium-2,low-1 | detailed in Appendix D | inq_last_6mths_cat |
| inq_last_12m | data['inq_last_12m']-data['inq_last_6mths'] | We want to capture the change in credit-seeking behavior of the borrower. If the difference is positive, it suggests an increase in credit-seeking behavior, which may be indicative of financial instability. Conversely, a negative difference could suggest a decrease in credit-seeking behavior, which may be indicative of greater financial stability. This new feature could potentially improve the model's ability to predict loan default risk. | recent_inquiry_difference |

| Mo_sin_rcnt_rev_tl_op and mo_sin_rcnt_tl | data['mo_sin_rcnt_rev_tl_op']/data['mo_sin_rcnt_tl'] | We discovered that the correlation between the columns "mo_sin_rcnt_rev_tl_op" and "mo_sin_rcnt_tl"  is 0.61- Thus, we created a new feature called "revolving_account_opened_ratio " by dividing "mo_sin_rcnt_rev_tl_op" by "mo_sin_rcnt_tl", aiming to capture the ratio of months since the most recent revolving account opened to months since the most recent account opened. The resulting feature may provide additional information about a borrower's credit behavior, which can be useful in predicting their ability to repay a loan. | revolving_account_opened_ratio |

**Appendix D:**

Note that many categorical variables already undergone transformations so reducing the number of variables requires further processing. You can refer to Appendix C.

| Column Name | Imputation Method | Suitable Model | |
|---|---|---|---|
| home_owners hip | We will one-hot encode the categories RENT and MORTGAGE since they have a loan percentage above 2% and the average yield is above/below the mean in the data. | For both |  |

| purpose | Leave One Out Encoding was used for the "purpose" variable in both regression and classification models due to its high cardinality and multiple categories. It helps prevent overfitting and improves model performance by using a more robust estimate of the mean for each category, making the models more generalizable. | For both |  |
|---|---|---|---|
| inq_last_6mths | We are transforming this variable into an ordinal categorical variable with three levels based on the number of credit inquiries made in the past 6 months. This is done to capture the relationship between inquiry levels and default rates, which can be used as a predictor in a machine learning model to predict loan defaults. Ordinal encoding is appropriate because the categories have a natural order and relative values. | For both |  |

Yearly yield mean by Loan Purpose Category
Mean of yield in all data: -9.2%

Loan failure rates by Number of inquiries in last 6 months
Percentage of yield under 2% in all data: 21.2%

Yearly yield mean by Number of inquiries in last 6 months
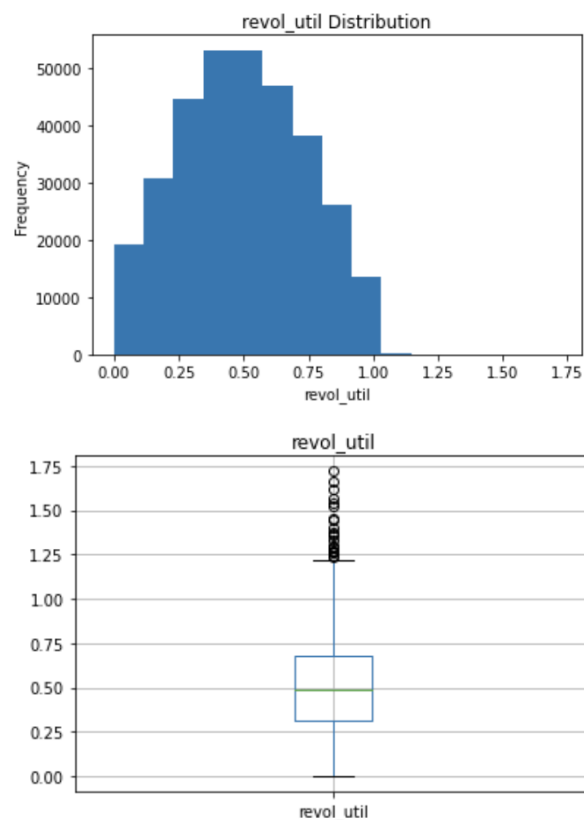Mean of yield in all data: -9.2%

**Appendix E:**

| Column Name | Treatment | |
|---|---|---|
| DTI | After observing that the distribution of our data appears to be approximately normal, we used a boxplot to identify any potential outliers. By examining the boxplot, we looked for data points that fell outside the whiskers, which are typically set at a distance of 1.5 times the interquartile range (IQR). However, in our case, we decided to be more conservative and considered data points that were more than 3 standard deviations away from the mean as outliers. |  |
| DTI | |  |
| annual_inc | when the distribution of a variable appears to be exponential or highly skewed, applying a logarithmic transformation can sometimes help to make the data more normally distributed. Taking the logarithm of the variable can compress the range of values, |  |

| | reduce the skewness, and make the distribution more symmetric. | |
|---|---|---|
| Total rev limit | when the distribution of a variable appears to be exponential or highly skewed, applying a logarithmic transformation can sometimes help to make the data more normally distributed. Taking the logarithm of the variable can compress the range of values, reduce the skewness, and make the distribution more symmetric. |  |
| revol_bal | |  |
| Revol_bal- num of outliers over 2000000 | 916 + when the distribution of a variable appears to be exponential or highly skewed, applying a logarithmic transformation can sometimes help to make the data more normally distributed. Taking the logarithm of the variable can compress the range of values, reduce the skewness, and make the distribution more symmetric. |  |

| | | |
|---|---|---|
| revol_util | After observing that the distribution of our data appears to be approximately normal, we used a boxplot to identify any potential outliers. By examining the boxplot, we looked for data points that fell outside the whiskers, which are typically set at a distance of 1.5 times the interquartile range (IQR). However, in our case, we decided to be more conservative and considered data points that were more than 3 standard deviations away from the mean as outliers. |  |

revol_util Distribution

revol_util

**Appendix F:**

| Column Name | Treatment | |
|---|---|---|
| installment | To reduce multicollinearity due to a very strong correlation (0.97, Spearman correlation) between 'installment' and 'loan_amnt' we decided to drop installment. | Scatter Plot: Relationship between installment and loan_amnt: strong positive association |
| tot_hi_cred_lim | To reduce the number of features and also reduce multicollinearity due to a very strong correlation (0.92, Spearman correlation) between 'tot_cur_bal' and 'tot_hi_cred_lim' we decided to drop installment. | Scatter Plot: Relationship between "tot_hi_cred_lim" and "tot_cur_bal": strong positive association |
| Bc_util_med - Ratio of total current balance to high credit/credit limit for all bankcard accounts.<br>Revol_util - Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit. | Both measures essentially reflect the borrower's credit utilization behavior, but at different levels of aggregation. The ratio of total current balance to high credit/credit limit focuses specifically on bankcard accounts, while the revolving line utilization rate considers all available revolving credit.<br>We created a Credit_Utilization_Score feature that represents the combined credit utilization across bankcard accounts. | Scatter Plot: Relationship between "bc_util_med" and "revol_util"<br>Polynomial Regression |