

Title: Credit Risk Modeling Notebook

1. Introduction:

- This notebook is centered around risk analytics, particularly within the banking and financial services sector using data from LendingClub, the largest peer-to-peer lending platform globally.
- The goal is to explore how data can be used to mitigate potential financial losses associated with lending to customers.

2. Data Understanding and Exploration:

- The dataset comprises information on past loan applicants and their default status.
- The primary objective is to identify patterns that indicate the probability of an applicant defaulting. This information can aid in decisions like loan denial, adjusting loan amounts, offering higher interest rates to risky applicants, and more.

3. Loading Libraries and Datasets:

- Libraries such as pandas, numpy, seaborn, matplotlib, and others are loaded to handle data manipulation, statistical analysis, and visualization tasks.
- Datasets `accepted_2007_to_2018Q4.csv` and `rejected_2007_to_2018Q4.csv` are loaded which contain information on accepted and rejected loans respectively.

4. Data Understanding:

- Initial examination of the data is performed by checking the first and last five rows of both datasets, and by checking the shape and unique values in the datasets.

5. Data Cleaning:

- Features from the accepted loans dataset are trimmed down to 20 significant features for easier analysis.
- Null values are identified and dropped from both datasets.
- Outliers are identified using the Interquartile Range (IQR) method and are handled.

6. Data Transformation:

- Functions are created to extract meaningful information from string values in the 'term' and 'Employment Length' columns.
- The Debt-To-Income Ratio in the rejected loans dataset is converted to a float type for numerical analysis.

7. Exploratory Data Analysis (EDA):

- Various visualizations are created to understand the distribution and relationship of different variables in the dataset.
 - Loan status count distribution shows the count of loans in different statuses.
 - Distribution of loan amounts and interest rates are visualized using histograms.

- Loan grade count distribution depicts the count of loans in different grades.
- Boxplots are used to show the distribution of loan amounts and interest rates across different loan grades.

8. Insights Derived:

- A significant number of loans are fully paid, with loans having a term of 36 months being the most common.
- Loan grades B and C have a higher frequency, indicating a moderate credit risk according to LendingClub's grading system.
- Employees with a tenure exceeding 10 years have better loan accessibility compared to those with shorter durations of employment.
- Loans around the 10,000 currency unit range are more prevalent, and a 15% interest rate is the most common amongst the loans.
- This analysis provides a preliminary understanding of the credit risk associated with lending to customers on the LendingClub platform.
- The insights derived could be used as a foundation for building predictive models to estimate the likelihood of a borrower defaulting on a loan.

-

9. **Target Variable Creation:** A new target variable `loan_status_log` was created by categorizing loan statuses into 'Current or Good' (1) and 'Default or Bad' (0) based on the values in the `loan_status` column.

10. Feature Engineering:

- Employment length (`emp_length`) was transformed into numerical values.
- The `purpose` variable was simplified into fewer categories.
- One-hot encoding was applied to categorical variables like `home_ownership`, `purpose`, `grade`, and `term`.

11. **Data Scaling and Transformation:** `StandardScaler` was utilized to scale the numerical features.

12. **Dataset Splitting:** The dataset was split into training and testing sets using a 70-30 split ratio.

Machine Learning Models

Four machine learning models were proposed for this analysis: Logistic Regression, Random Forest Classifier, Support Vector Classifier, and K Nearest Neighbors. However, only Logistic Regression and Random Forest Classifier were implemented and evaluated.

1. Logistic Regression:

- **Accuracy:** 87.01%
- **Precision:** 87.01%
- **Recall:** 100%
- **F1 Score:** 93.05%

- The model failed to identify any loans with credit risk, possibly due to class imbalance in the dataset.

2. Random Forest Classifier:

- **Accuracy:** 87.04%
- **Precision:** 75.75%
- **Recall:** 87.04%
- **F1 Score:** 81.00%
- Similar to Logistic Regression, this model too failed to identify any loans with credit risk.

Evaluation Metrics

- **Precision:** Proportion of true positive predictions among all positive predictions.
- **Recall:** Proportion of true positive predictions among all actual positives.
- **F1 Score:** Harmonic mean of precision and recall, providing a balance between the two metrics.
- **Confusion Matrix:** Provides a visualization of the model's performance on the test set.

Observations and Recommendations

- Both models exhibit a high accuracy rate, but fail to identify any loans with credit risk, suggesting a potential issue of class imbalance.
- More sophisticated techniques such as class balancing, using different algorithms, or feature engineering might help in addressing the shortcomings of the current models.
- Further exploration with the Support Vector Classifier, K Nearest Neighbours, and potentially other models like Gradient Boosting or Neural Networks could provide more insights.
- Fine-tuning the models with a broader range of hyperparameters or employing techniques like cross-validation could also contribute to improving the model performance.

This documentation provides a structured summary of the analysis workflow, evaluation of machine learning models, and suggestions for potential improvements in predicting the credit risk associated with loans in the dataset.