

# Text Classification- Loan Data Analysis Use Case

Promothesh Chatterjee\*

---

\*Copyright© by Promothesh Chatterjee. All rights reserved. No part of this note may be reproduced, distributed, or transmitted in any form or by any means, including photocopying, recording, or other electronic or mechanical methods, without the prior written permission of the author.

## Understanding Loan Default Analysis

**Overview** Loan default analysis is crucial for financial institutions to minimize risk and maximize profitability. By examining different variables that influence loan defaults, we can predict and mitigate potential risks. This guide explains the importance of various variables in loan default analysis using the provided synthetic dataset.

**1. Variables in Loan Default Analysis** Understanding the impact of different variables helps in building robust predictive models. I have created a sythetic dataset (code posted in Canvas). Here are the key variables considered in the dataset:

1. **CustomerID**: Unique identifier for each customer.
2. **LoanAmount**: The amount of money borrowed by the customer.
3. **InterestRate**: The interest rate applied to the loan.
4. **LoanTerm**: The duration over which the loan must be repaid.
5. **CreditScore**: A numerical expression representing the creditworthiness of the customer.
6. **EmploymentStatus**: The employment status of the customer (Employed, Unemployed, Self-Employed).
7. **AnnualIncome**: The yearly income of the customer.
8. **LoanPurpose**: The reason for taking the loan (Home, Car, Education, Business, Other).
9. **Default**: Indicates whether the customer defaulted on the loan (Yes, No).
10. **ReviewText**: Textual data describing the loan purpose, generated based on writing style.
11. **WritingStyle**: The style in which the review text is written (Formal, Informal, Descriptive, Concise, Narrative).

## 2. Importance of Each Variable

- **CreditScore**: This is a critical variable. Lower credit scores often indicate a higher risk of default. The model uses a coefficient of -0.05 for CreditScore, suggesting that higher scores reduce the likelihood of default.
- **LoanAmount**: Higher loan amounts can increase the risk of default, as reflected by a positive coefficient (0.0065) in the logistic regression model.
- **AnnualIncome**: Higher income generally reduces the likelihood of default, with a negative coefficient (-0.005). It indicates that individuals with higher incomes are less likely to default on loans.
- **EmploymentStatus**: Employment status is significant, especially being unemployed, which drastically increases the likelihood of default (coefficient of 1.0).
- **InterestRate**: Higher interest rates can increase the financial burden on borrowers, leading to a higher probability of default (coefficient of 0.22).
- **LoanPurpose**: Different loan purposes have varying risks associated with them. For example, home loans might have a different default risk compared to car loans.
- **LoanTerm**: The duration of the loan affects repayment capacity. Shorter terms mean higher monthly payments, which could increase default risk.

**3. Generating the Default Variable** The default variable is generated using logistic regression probabilities. The logistic function converts the linear predictor into a probability. A set threshold determines whether a loan defaults or not.

**4. Textual Data Analysis** Textual data such as **ReviewText** and **WritingStyle** can provide additional insights. Analyzing sentiments, tone, and style can help in understanding the borrower's mindset and potential risk.

**5. Synthetic Data Generation** The dataset is generated synthetically with meaningful associations between variables: - **LoanPurpose** and **WritingStyle** are randomly assigned. - **CreditScore** and **InterestRate** have realistic constraints to reflect actual conditions. - **ReviewText** is generated based on a helper function using the loan purpose and writing style. - The default probability is calculated using a logistic regression model, ensuring a realistic default rate (~5%).

**6. Importance of Data Quality** High-quality, accurate data is essential for reliable loan default predictions. Synthetic data, while useful for educational purposes, must be validated against real-world data to ensure model robustness.

**7. Practical Applications** Understanding these variables and their impact on loan defaults allows financial institutions to: - Develop better risk assessment models. - Customize loan products based on risk profiles. - Implement targeted strategies for different customer segments.

**Conclusion** Analyzing loan defaults involves understanding multiple factors and their interplay. By comprehensively evaluating variables such as credit scores, loan amounts, income levels, and more, business analytics professionals can develop effective predictive models to manage and mitigate risks associated with loan defaults.