

Exploratory Data Analysis (EDA) Summary

Report Template

1. Introduction

The goal of this project report is to design an AI-powered solution that enables Geldium, a digital lending provider, to significantly improve its ability to predict and manage credit card delinquency. Faced with rising numbers of customers missing payments beyond 30 days, the collections team currently relies on historical trends and manual case handling, limiting efficiency and decision-making speed. The purpose of the report is to demonstrate how artificial intelligence, particularly GenAI techniques, can identify at-risk customer's early, forecast delinquency likelihood, and guide targeted intervention strategies. The report will include an exploratory analysis of customer data to detect patterns and key risk factors, a predictive model design blueprint, and proposed intervention approaches. It will also emphasize ethical considerations and industry best practices to ensure fairness and transparency in decision-making. Ultimately, the report aims to help Geldium reduce delinquencies, optimize collection efficiency, and enhance customer experience by supporting proactive credit risk management with AI-based insights.

2. Dataset Overview

This section summarizes the dataset, including the number of records, key variables, and data types. It also highlights any anomalies, duplicates, or inconsistencies observed during the initial review.

Key dataset attributes:

- Number of records: **500**

Column Name	Description	Data Type
-------------	-------------	-----------

Customer_ID	Unique Identifier for each customer	Categorical
Age	Customers age	Numerical
Income	Annual income of the customer	Numerical
Credit_Score	Credit score rating	Numerical
Credit_Utilization	Ratio of credit used to available credit	Numerical
Missed_Payments	Count of missed payments	Numerical
Delinquent_Account	Binary indicator of delinquency(1= delinquent)	categorical
Loan_Balance	Outstanding loan balance	Numerical
Debt_to_Income_Ratio	Ratio of debt to income	Numerical
Employment_Status	Employment category('Employed', 'Unemployed', 'Self-employed')	Categorical
Account_Tenure	Duration of account ownership in years	Numerical
Credit_Card_Type	Type of credit card held('Standard', 'Gold', 'Platinum')	Categorical
Location	City of residence	Categorical
Month_1 to Month_6	Monthly payment status('On-time', 'Late', 'Missed')	Categorical

Descriptive highlights

Variable	Mean	Std	Min	Max	Observations
Age	46.3	16.2	18	74	Broad age range, evenly distributed.

Variable	Mean	Std	Min	Max	Observations
Income	108,380	53,663	15,404	199,943	High variability (possible income inequality).
Credit Score	578	169	301	847	Covers poor to excellent credit ranges.
Credit Utilization	0.49	0.20	0.05	1.03	Some over-utilization (>1) indicates risk.
Debt-to-Income Ratio	0.30	0.09	0.10	0.55	Mostly moderate but some risky levels (>0.5).
Loan Balance	48,654	29,396	612	99,620	Large variability in debt loads.

Outliers

- **Credit_Utilization** slightly exceeding 1.0 which may indicate data entry errors or unusual financial behavior.
- **Debt-to-Income Ratio** near 0.55 suggest financial stress cases that may correlate with delinquency.
- **Loan_Balance**, several entries are missing which may distort debt-related metrics.

Missing values and inconsistencies

- **Income:** has 39 missing cells which could affect financial ratio modelling thus should be carefully imputed, i.e. Median or regression imputation.
- **Credit_Score:** has 2 missing cells which is minimal hence can be safely imputed (mean/median).
- **Loan_Balance:** has 29 missing cells which might indicate customers without loans (could encode as 0 or "No Loan").

Top 3 predictive variables for delinquency

- **Missed_Payments :** Strong direct correlation with delinquency; higher counts often precede delinquent status.

- **Credit_Utilization :** High utilization suggests financial stress, which is a known risk factor for delinquency.
- **Debt_to_Income_Ratio:**
Indicates repayment capacity; higher ratios often align with delinquent behavior

Fields that might cause problems

- Missing or invalid entries were noted in Income, Credit_Score, and Loan_Balance fields due to incomplete data submission or unmatched records.
- Some variables showed extremely high leverage or utilization (>1), indicating possible data entry errors or risky credit usage.
- Customer_ID is a unique identifier which makes it not useful for modeling.
- Employment_Status is non-numeric, contains mixed categories (e.g., “EMP”, “Self-employed”, “Unemployed”), and requires encoding and consistency checks.
- Month_1–Month_6 are Categorical ("On-time", "Late", "Missed") suggests aggregating into counts or weighted scores (e.g., ordinal or frequency encoding).

3. Missing Data Analysis

Identifying and addressing missing data is critical to ensuring model accuracy. This section outlines missing values in the dataset, the approach taken to handle them, and justifications for the chosen method.

Key missing data findings:

Variables with missing values:

- Income – 39
- Credit_Score – 2
- Loan_Balance – 29

Other columns were complete with no missing entries i.e. Age Credit_Utilization, Debt_to_Income_Ratio and Missed_Payments.

Missing data treatment:

- **Imputation (continuous variables):** Missing values in **Income** and **Loan_Balance** were imputed using the median value to retain distribution characteristics and reduce the effect of extreme values.
- **Mean Substitution (small gaps):** Missing **Credit_Score** entries were replaced using the mean score value since the missing count was very low.
- **No deletion:** Deleting records was avoided to prevent potential bias and loss of valuable information, especially since missing values were below 10% for most variables.

Justification:

- Median imputation helps manage skewed distributions often seen in financial datasets.
- Mean substitution is suitable for variables with minimal missingness.
- Retaining all data points ensures sufficient sample size for training predictive models without introducing sampling bias.

4. Key Findings and Risk Indicators

This section identifies trends and patterns that may indicate risk factors for delinquency. Feature relationships and statistical correlations are explored to uncover insights relevant to predictive modeling.

Key findings:

Correlations observed between key variables:

- Strong positive correlation was observed between **credit utilization** and delinquency rates in that, customers with utilization above 0.6 are more likely to miss payments.

- **Debt-to-income ratio** vs Delinquent_Account shows moderate correlation ($r \approx 0.35$) higher ratios above 0.35 indicate repayment stress and greater risk.
- **Credit Score** is negatively correlated with Delinquent_Account ($r \approx -0.5$) lower scores reflect higher risk of missed payments.
- **Missed payments history** (MSP1–MSP6) also shows significant correlation to Delinquent_Account ($r \approx 0.5-0.7$) quality of previous on-time payment behavior is a strong indicator of future late repayment.
- **Monthly_Status (Month_1 to Month_6)** correlate with Delinquent_Account due to frequent “Missed” and “Late” entries.

Unexpected anomalies that require further investigations:

- Monthly_Status column need transformation, some rows show all “Missed” or all “On-time” which may be outliers.
- Employment_Status variants has inconsistent labels like “EMP”, “employed” and “Employed” need standardization.
- A small cluster with high income (>200k) but consistently high delinquency suggests behavioral risk rather than financial constraints. Few customers report extremely low income which may be misreported.
- Some customers with long tenure (>120 months) and high utilization (>1.0) never miss payment, may indicate preferential borrowing or secured collateral cases.
- Some entries show questionably low balances despite high income or tenure under the Loan_Balance column.

5. AI & GenAI Usage

Generative AI tools were used to summarize the dataset, impute missing data, and detect patterns. This section documents AI-generated insights and the prompts used to obtain results.

Example AI prompts used:

- "Analyze this dataset and summarize key patterns, outliers, and missing values."
- "Highlight any fields that might present problems for modeling delinquency."

- "Identify the top 3 variables most likely to predict delinquency based on this dataset."
- "Dataset Overview: number of records, key variables, and data types."
- "Missing Data Analysis: variables with missing values and treatment approach."
- "Identify trends and patterns that may indicate risk factors for delinquency."
- "Summarize correlations observed between key variables."
- "Highlight unexpected anomalies requiring further investigation."
- "Summarize key findings and outline the recommended next steps."

6. Conclusion & Next Steps

Summary of key findings

- **Strong Predictors of Delinquency** – Variables like **credit utilization**, **debt-to-income ratio**, and **past missed payments** consistently show the strongest correlation with future late payments. Higher utilization (>0.75) and ratios above ~ 0.4 significantly increase bad-debt risk.
- **Data Quality Insights** – Missing values were most common in **income**, **loan balance**, and **credit score**, requiring targeted imputation. Applying median imputation and AI-based methods helped maintain distribution integrity without reducing sample size.
- **Behavioral Patterns** – AI analysis highlighted critical patterns such as increasing utilization over consecutive months, sudden income drops, and repeated partial payments, all linked to higher default likelihood.
- **Anomalies Detected** – Some high-income customers exhibited delinquency, suggesting that lending behavior, not just finances, is a key predictor. Long-tenure accounts with new delinquency spikes may indicate emerging stress.

Recommended Next Steps

- **Build the Predictive Model:** Use identified key predictors to train machine-learning classifiers (e.g., logistic regression with regularization, tree-based models, or neural networks), validating with k-fold cross-validation.
- **Create Risk Engine Logic:** Develop a scoring mechanism that classifies customers into low, medium, and high risk based on feature thresholds and predicted probabilities.
- **Implement AI-Assisted Notifications:** Integrate the model into the collections workflow to trigger proactive interventions for at-risk accounts.
- **Continuous Monitoring:** Periodically retrain models with new data, monitor key performance metrics (ROC-AUC, precision-recall), and audit for bias or drift.
- **Feature Expansion:** Consider advanced features like cash-flow trends, transaction frequency, and external credit bureau data to further improve predictive accuracy.