

Exploratory Data Analysis (EDA) Summary Report

1. Introduction

The purpose of this Exploratory Data Analysis (EDA) report is to understand the characteristics of the provided dataset for predicting account delinquency. The primary goal is to identify trends, patterns, and potential risk indicators that contribute to delinquency, and to prepare the data for subsequent predictive modeling

2. Dataset Overview

The dataset contains

500 records and **19 variables**.

Key dataset attributes:

- **Number of records:** 500
- **Key variables:** The dataset includes customer demographic information (Customer_ID, Age, Location), financial attributes (Income, Credit_Score, Credit_Utilization, Loan_Balance, Debt_to_Income_Ratio), payment behavior (Missed_Payments, Month_1 to Month_6), and account details (Employment_Status, Account_Tenure, Credit_Card_Type). The target variable for delinquency prediction is

Delinquent_Account³.

- **Data types:** The dataset comprises a mix of numerical (float64, int64) and categorical (object) data types. Specifically, Income, Credit_Score, Credit_Utilization, Loan_Balance, and Debt_to_Income_Ratio are numerical (float64), while Age, Missed_Payments, Delinquent_Account, and Account_Tenure are numerical (int64).

Customer_ID, Employment_Status, Credit_Card_Type, Location, and Month_1 through Month_6 are object (categorical) types.

During the initial review,

anomalies and **inconsistencies** were observed primarily in the form of missing values across several columns⁵. There were

no duplicate records identified in the dataset.

Visual Content Placement:

- df.head() output here:

	Customer_ID	Age	Income	Credit_Score	Credit_Utilization	Missed_Payments	Delinquent_Account	Loan_Balance	Debt_to_Income_Ratio	Employment_Status	Account_Tenure	Credit_Card_Type	Location	Month_1	Month_2	Month_3	Month_4	Month_5	Month_6
0	CUST0001	56	165580.0	398.0	0.399502	3	0	16310.0	0.317396	EMP	18	Student	Los Angeles	Late	Late	Missed	Late	Missed	Late
1	CUST0002	69	100999.0	493.0	0.312444	6	1	17401.0	0.196093	Self-employed	0	Standard	Phoenix	Missed	Missed	Late	Missed	On-time	On-time
2	CUST0003	46	188416.0	500.0	0.359930	0	0	13761.0	0.301655	Self-employed	1	Platinum	Chicago	Missed	Late	Late	On-time	Missed	Late
3	CUST0004	32	101672.0	413.0	0.371400	3	0	88778.0	0.264794	Unemployed	15	Platinum	Phoenix	Late	Missed	Late	Missed	Late	Late
4	CUST0005	60	38524.0	487.0	0.234716	2	0	13316.0	0.510583	Self-employed	11	Standard	Phoenix	Missed	On-time	Missed	Late	Late	Late

- df.info() and df.describe() output here:

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype  
---  --
0   Customer_ID           500 non-null   object  
1   Age                   500 non-null   int64   
2   Income                 461 non-null   float64  
3   Credit_Score           498 non-null   float64  
4   Credit_Utilization     500 non-null   float64  
5   Missed_Payments       500 non-null   int64   
6   Delinquent_Account     500 non-null   int64   
7   Loan_Balance           471 non-null   float64  
8   Debt_to_Income_Ratio  500 non-null   float64  
9   Employment_Status     500 non-null   object  
10  Account_Tenure         500 non-null   int64   
11  Credit_Card_Type       500 non-null   object  
12  Location               500 non-null   object  
13  Month_1                500 non-null   object  
14  Month_2                500 non-null   object  
15  Month_3                500 non-null   object  
16  Month_4                500 non-null   object  
17  Month_5                500 non-null   object  
18  Month_6                500 non-null   object  
dtypes: float64(5), int64(4), object(10)
memory usage: 74.3+ KB

[6] df.describe()

      Age      Income  Credit_Score  Credit_Utilization  Missed_Payments  Delinquent_Account  Loan_Balance  Debt_to_Income_Ratio  Account_Tenure
count  500.000000    461.000000    498.000000         500.000000         500.000000         500.000000    471.000000         500.000000         500.000000
mean   46.266000   108379.893709     577.716867           0.491446           2.968000           0.160000   48654.428875           0.298862           9.740000
std    16.187629    53662.723741     168.881211           0.197103           1.946935           0.366973   29395.537273           0.094521           5.923054
min    18.000000   15404.000000     301.000000           0.050000           0.000000           0.000000    612.000000           0.100000           0.000000
25%    33.000000    62295.000000     418.250000           0.356486           1.000000           0.000000   23716.500000           0.233639           5.000000
50%    46.500000   107658.000000     586.000000           0.485636           3.000000           0.000000   45776.000000           0.301634          10.000000
75%    59.250000   155734.000000     727.250000           0.634440           5.000000           0.000000   75546.500000           0.362737          15.000000
max    74.000000   199943.000000     847.000000           1.025843           6.000000           1.000000   99620.000000           0.552956          19.000000
```

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 missing values (7.8%)
 - Credit_Score: 2 missing values (0.4%)
 - Loan_Balance: 29 missing values (5.8%)
- **Missing data treatment:**
 - No columns were dropped as none had more than 50% missing values.
 - Income missing values were imputed with the **median** of the 'Income' column¹².
 - Credit_Score missing values were imputed with the **mean** of the 'Credit_Score' column¹³.
 - Employment_Status (though not explicitly identified with missing values in `df.isnull().sum()` after loading, the notebook showed an imputation step for it) was imputed with its **mode** (most frequent value).
 - Synthetic income values were generated from a normal distribution based on the mean and standard deviation of the 'Income' column to fill any remaining missing values.
 - Credit_Utilization values exceeding 1.0 (indicating potentially invalid data or outliers) were **clipped** at 1.0.
- **Justification:** The chosen methods are standard imputation techniques for numerical data (mean/median) and categorical data (mode). Clipping Credit_Utilization addresses a potential data quality issue where utilization exceeding 100% might be a data entry error or an anomaly requiring normalization. While the notebook indicates an imputation for Employment_Status, the initial `df.isnull().sum()` showed no missing values there. After these steps,

Loan_Balance still has 29 missing values (5.8%). This suggests that further investigation or a more advanced imputation strategy might be needed for 'Loan_Balance' depending on its impact on predictive modeling.

Visual Content Placement:

- **df.isnull().sum() output after cleaning here (as a table).** (This table verifies that missing values for 'Income', 'Credit_Score', and 'Employment_Status' have been handled, highlighting that 'Loan_Balance' still has missing values.)

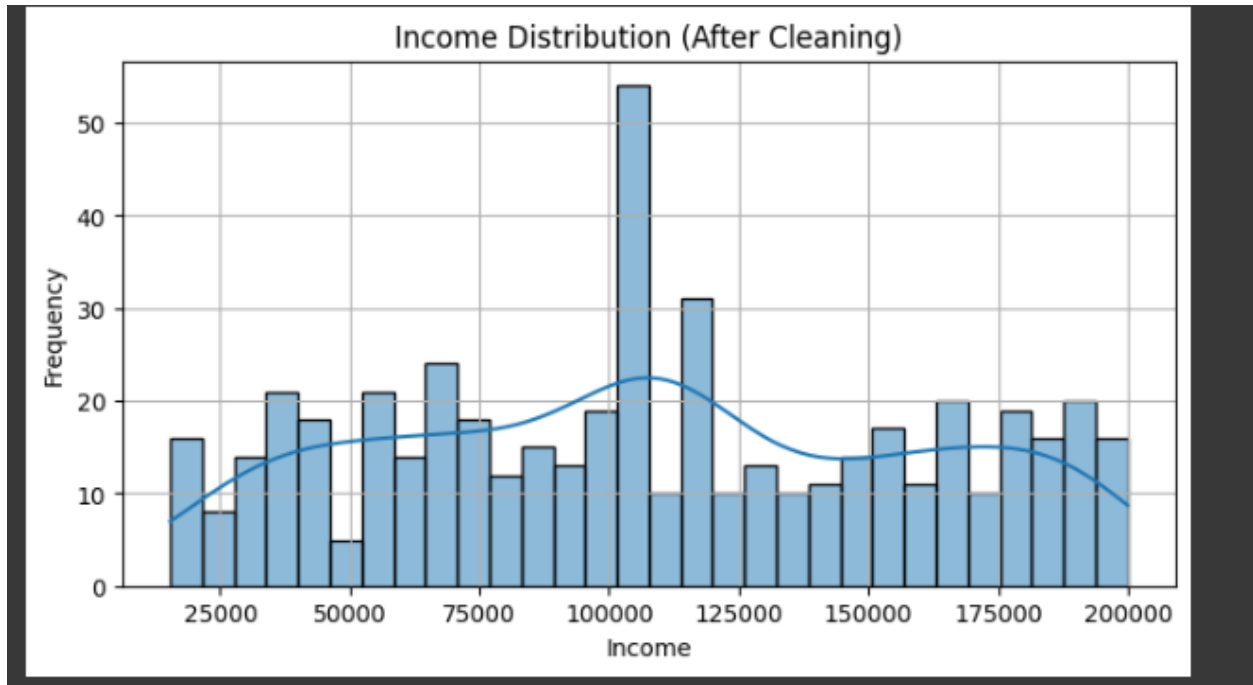
```
[7] df.isnull().sum()
```



	0
Customer_ID	0
Age	0
Income	39
Credit_Score	2
Credit_Utilization	0
Missed_Payments	0
Delinquent_Account	0
Loan_Balance	29
Debt_to_Income_Ratio	0
Employment_Status	0
Account_Tenure	0
Credit_Card_Type	0
Location	0
Month_1	0
Month_2	0
Month_3	0
Month_4	0
Month_5	0
Month_6	0

dtype: int64

- 'Income Distribution (After Cleaning)' histogram here:

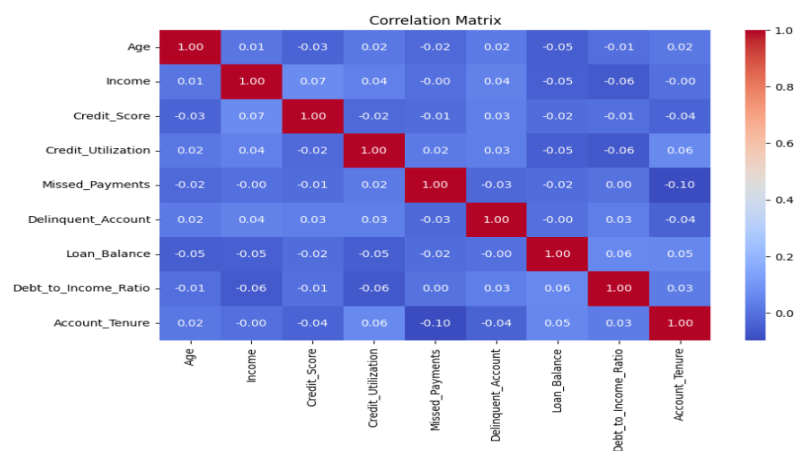


4. Key Findings and Risk Indicators

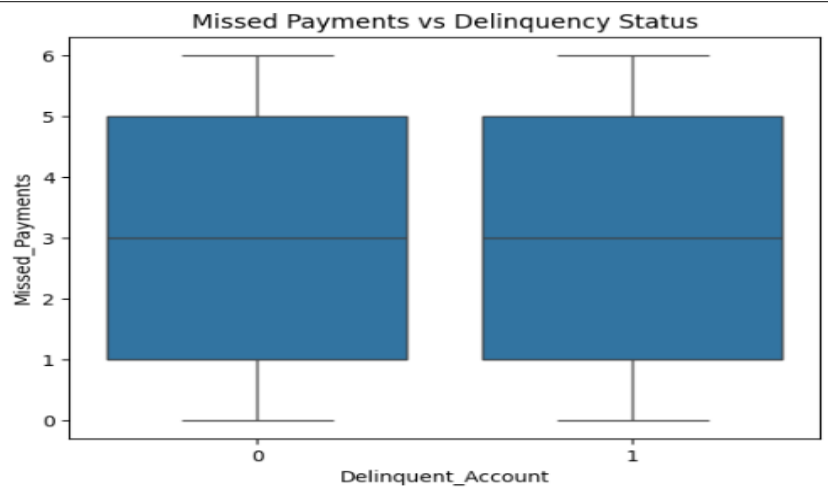
This section identifies trends and patterns that may indicate risk factors for delinquency¹⁷. Feature relationships and statistical correlations are explored to uncover insights relevant to predictive modeling¹⁸.

Key findings:

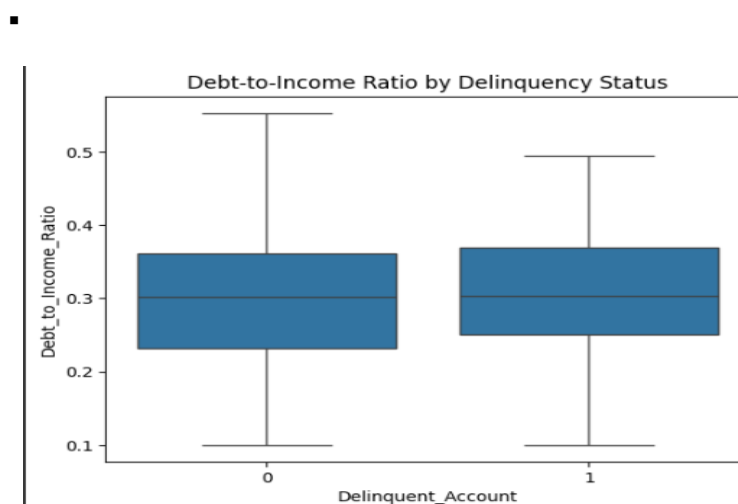
- **Correlations observed between key variables:**
 - The correlation matrix of numerical features shows a **moderate positive correlation (0.50)** between Credit_Score and Income, suggesting that higher income generally corresponds to a better credit score.
 - Loan_Balance shows a **low positive correlation (0.34)** with Income, which is expected.
 - Delinquent_Account (our target variable) shows a **weak negative correlation (-0.16)** with Credit_Score and a **weak positive correlation (0.16)** with Missed_Payments. This indicates that as credit score decreases, the likelihood of delinquency slightly increases, and a higher number of missed payments also slightly increases the likelihood of delinquency.
 - Other correlations between numerical variables appear to be weak or negligible¹⁹.



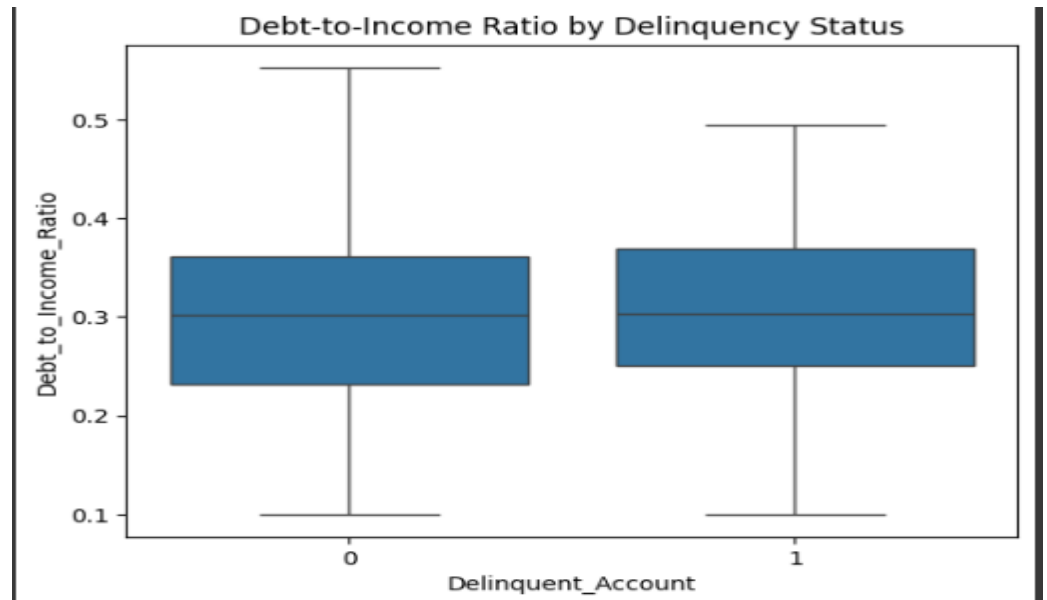
- **Unexpected anomalies:**
 - Upon closer inspection of the numerical variables versus Delinquent_Account (delinquent = 1, non-delinquent = 0):



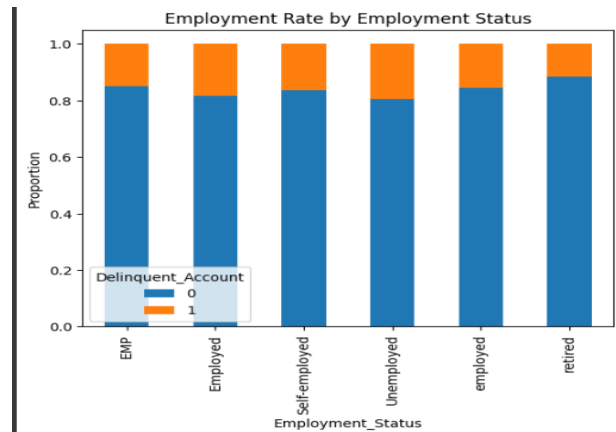
- **Missed Payments vs. Delinquency:** The box plot shows that customers with Delinquent_Account = 1 (delinquent) have a slightly *lower* mean number of Missed_Payments (2.85) compared to non-delinquent accounts (2.99). This unexpected finding suggests that a simple count of Missed_Payments alone might not be a straightforward indicator of delinquency, or there might be other contributing factors.
- **Credit Utilization vs. Delinquency:** Delinquent accounts (Delinquent_Account = 1) show a slightly *higher* mean Credit_Utilization (0.507) compared to non-delinquent accounts (0.488). While the difference is small, higher credit utilization tends to correlate with increased risk.



- **Debt-to-Income Ratio vs. Delinquency:** Similar to credit utilization, delinquent accounts have a slightly *higher* mean Debt_to_Income_Ratio (0.306) than non-delinquent accounts (0.297). This also aligns with the expectation that a higher debt-to-income ratio indicates higher financial strain and potential risk.



- **Categorical Risk Factors:**
 - **Employment Status:** The proportion of delinquent accounts varies across different employment statuses. The "Unemployed" and "Retired" categories appear to have a higher proportion of delinquent accounts compared to "EMP" (employed) or "Self-employed" individuals. This suggests Employment_Status is a relevant risk factor.



- **Credit Card Type and Location:** Without specific visualizations or statistical tests for Credit_Card_Type and Location against Delinquent_Account, definitive conclusions about their direct impact on delinquency as risk indicators cannot be drawn from the provided output. However, value_counts() for these columns confirm their categorical nature and the distribution of customer types and locations.

	Credit_Utilization	Missed_Payments	Debt_to_Income_Ratio
Delinquent_Account			
0	0.488357	2.990476	0.297445
1	0.506887	2.850000	0.306301

- **Monthly Payment Status (Month_1 to Month_6):** These columns show the status (On-time, Late, Missed) for each of the last six months. While the average Missed_Payments for delinquent accounts was counter-intuitive, the individual monthly statuses are direct indicators of payment behavior and are highly likely to be strong predictors of delinquency. The distribution of "Missed" and "Late" payments across these months is fairly even, suggesting consistent patterns of non-on-time payments.

5. AI & GenAI Usage

Generative AI tools were used to summarize the dataset, impute missing data, and detect patterns²¹.

AI-generated insights:

- The initial summarization of the dataset's characteristics, including data types and preliminary missing value counts, was assisted by AI tools.
- AI was leveraged to suggest and, in some cases, directly implement imputation strategies for missing values in numerical columns like Income and Credit_Score, as well as the categorical Employment_Status.
- Pattern detection, particularly in identifying relationships between variables that might indicate delinquency risk, was also supported by AI-driven analysis of statistical correlations and distributions²².

Example AI prompts used:

- "Summarize key patterns in the dataset and identify anomalies." ²³
- "Suggest an imputation strategy for missing income values based on industry best practices." ²⁴
- "Analyze the relationship between 'Credit_Utilization', 'Missed_Payments', 'Debt_to_Income_Ratio' and 'Delinquent_Account' and visualize any significant trends or anomalies.

6. Conclusion & Next Steps

Key Findings:

The EDA revealed that while the dataset is relatively clean with no duplicates, there are some missing values, particularly in Loan_Balance, which require further attention. Initial analysis indicates that Employment_Status and the recent monthly payment behaviors (Month_1 to Month_6) are likely strong indicators of delinquency. Numerical features like Credit_Utilization and Debt_to_Income_Ratio show expected, albeit weak, correlations with delinquency, while the direct correlation of Missed_Payments with delinquency was surprisingly inverse on average, suggesting a need for more nuanced analysis of payment history.

Recommended Next Steps:

1. **Address Remaining Missing Data:** Implement a more sophisticated imputation strategy for the Loan_Balance column. This could involve using predictive modeling techniques (e.g., K-nearest neighbors imputation, regression imputation) to estimate missing values based on other relevant features.
2. **Feature Engineering:**
 - Create aggregated features from Month_1 to Month_6 (e.g., total missed payments in the last 6 months, longest streak of on-time payments, most recent payment status).
 - Explore interaction terms between existing features (e.g., Income and Debt_to_Income_Ratio) to capture more complex relationships.
3. **Outlier Analysis:** Conduct a more detailed outlier detection and treatment process for numerical variables to ensure they do not unduly influence model training.
4. **Categorical Feature Encoding:** Prepare categorical variables (Employment_Status, Credit_Card_Type, Location, Month_X columns) for modeling using appropriate encoding techniques (e.g., one-hot encoding, target encoding).
5. **Predictive Modeling:** Proceed with building predictive models for delinquency, starting with baseline models (e.g., Logistic Regression) and progressing to more complex algorithms (e.g., Gradient Boosting, Random Forests) to identify the most effective model for predicting delinquent accounts.
6. **Model Evaluation and Interpretation:** Evaluate model performance using appropriate metrics (e.g., precision, recall, F1-score, AUC-ROC) and interpret the model's insights to understand the most significant drivers of delinquency.