# Exploratory Data Analysis (EDA) Report

## 1. Project & Data Overview

- Dataset Name: Loan Approval Prediction Dataset
- ML Task: Binary Classification (Predicting Loan_Status).
- Total Observations : 614
- Total Features ): 13

---

## 2. Data Quality & Missing Values

- Crucial Finding: There are NO missing values in the dataset.
  - Action for ML Engineer: This is ideal. No imputation step is required. The data is clean in terms of completeness.

---

## 3. Target Variable Analysis

- Target Variable: Loan_Status
- Distribution:
  - Y (Approved): 68.73%
  - N (Denied): 31.27%
- ML Impact: The target variable exhibits a mild class imbalance (roughly 2:1 ratio).
  - Recommendation:
    - Standard classification algorithms (e.g., Logistic Regression, Tree-based models) should perform reasonably well, but the imbalance must be monitored.
    - Evaluation Metrics must prioritize F1 score and AUC-ROC over simple Accuracy.
    - If initial models struggle, consider techniques like SMOTE (Oversampling) or adjusting class weights (e.g., in XGBoost/Random Forest).

---

## 4. Numerical Feature Analysis (Outliers & Skewness)

The summary statistics indicate potential issues with skewness and outliers in the income features, typical for financial data.

| Feature | Min | Max | Mean | 75th Percentile | Observation |
|---|---|---|---|---|---|
| ApplicantIncome | 150 | 10171 | 4617 | 5795 | High Range/Outliers: Max (10171) is significantly higher than the 75th percentile (5795), suggesting a right skew and potential outliers (high-income applicants). |
| CoapplicantIncome | 0 | 5743 | 1420 | 2297 | Zero Values: The 25th percentile is 0.00, indicating a substantial number of applications where the applicant has no co applicant income. |
| LoanAmount | 9 | 262 | 137 | 165 | Skewness: The Max (262) is distant from the 75th percentile (165), suggesting some larger loans act as outliers. |

- 
  ML Preprocessing Recommendation:
  1. Transformation: Apply a Log Transformation to ApplicantIncome, CoapplicantIncome (handle 0.0 values first, e.g., log(1+X), and LoanAmount to normalize the distributions and mitigate the impact of outliers.
  2. Feature Engineering: Create a single, more stable predictor: Total_Income = ApplicantIncome + CoapplicantIncome.

## 5. Bivariate Analysis: Feature-Target Relationship

The relationship between Credit_History and Loan_Status is extremely strong and highly critical for modeling.

| Credit_History | Denied (N) Rate | Approved (Y) Rate | Observation |
|---|---|---|---|
| 0.0 (No/Bad History) | 92.13% | 7.87% | Almost all applicants with a bad credit history are DENIED the loan. |
| 1.0 (Good History) | 20.95% | 79.05% | The majority of applicants with a good credit history are APPROVED. |

- 

  ML Impact: Credit_History is likely the single most predictive feature. Model interpretability (XAI) should validate its high importance. Any model that fails to leverage this feature's power will perform poorly.

---

## 6. Summary of ML Preparation Steps

1. Data Cleaning: No missing value imputation or duplicate handling is required (data is clean).
2. Feature Engineering:
   - Drop the non-predictive Loan_ID.
   - Create Total_Income = ApplicantIncome + CoapplicantIncome.
3. Numerical Processing:
   - Apply Log Transformation to Total_Income and LoanAmount.
   - Apply Standard Scaling or MinMaxScaler to the transformed numerical features.
4. Categorical Processing:
   - Apply One-Hot Encoding to the remaining categorical features (e.g., Gender, Married, Education, Property_Area, etc.).

5. Modeling: Start with robust classifiers like Logistic Regression (due to the strong Credit_History factor) and Tree-based models (Random Forest/XGBoost), and focus evaluation on AUC-ROC.