

Financial Risk Analysis Report

Yaw Jantuah Boateng

February 2026

1 Introduction

This report summarizes the cleaning and analysis process on a loan default dataset. A publicly available financial dataset was extracted from Kaggle for this exercise. The data was then cleaned and analysed with Python using Google Colab notebook. The process was highlighted on GitHub repositories to show a reproducible process. A report was also prepared with LaTeX using Overleaf editor.

2 Data Cleaning

Missing values were explored in the dataset. Fortunately, there were no missing values in any of the columns, implying no need to perform any missing data imputation. Duplicates were then successfully removed from the data.

```
1 # Check missing values
2 df.isnull().sum()
3
4 # Remove duplicates
5 df = df.drop_duplicates()
```

3 Initial Findings

Summary statistics indicate a total of 255,347 participants in the data, with an average age of 43 and a standard deviation of 14.99. Income averages about 82,500, and the loan amount is around 127,500. Loan default shows a remarkable average of 0.12 over an average loan term of 36 months. The average credit score was also recorded to be 574.26.

Table 1: Summary Statistics After Cleaning

	Age	Income	LoanAmount	CreditScore	MonthsEmployed	NumCreditLines
count	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000
mean	43.498306	82499.304597	127578.865512	574.264346	59.541976	2.500000
std	14.990258	38963.013729	70840.706142	158.903867	34.643376	1.118034
min	18.000000	15000.000000	5000.000000	300.000000	0.000000	1.000000
25%	31.000000	48825.500000	66156.000000	437.000000	30.000000	2.000000
50%	43.000000	82466.000000	127556.000000	574.000000	60.000000	2.000000
75%	56.000000	116219.000000	188985.000000	712.000000	90.000000	3.000000
max	69.000000	149999.000000	249999.000000	849.000000	119.000000	4.000000

Basic distribution analysis highlights that income and loan amounts appear to be right-skewed, suggesting the presence of high-value outliers or a concentration of lower values with a tail extending to higher values.

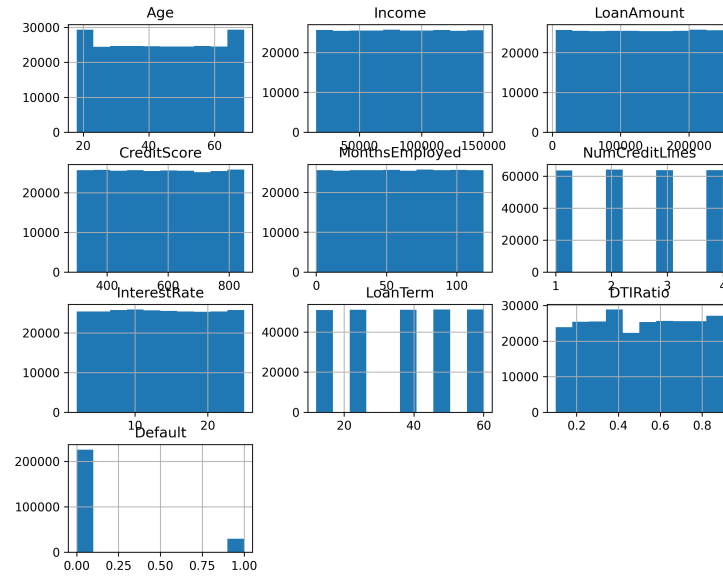


Figure 1: Basic Distribution Analysis

4 Conclusion

The dataset is now clean and ready for further modelling. The outputs of the basic analysis provide insights into the dataset's characteristics, such as the absence of missing values and the skewed distributions of income and loan amount.

A Full Python Code

```
1  # upload data
2  from google.colab import files
3  uploaded = files.upload()
4
5  import pandas as pd
6
7  # Load the dataset
8  df = pd.read_csv("raw_data.csv.csv")
9
10 # Show first 5 rows
11 df.head()
12
13 # Check structure
14 df.info()
15
16 # Check summary statistics
17 summary_table = df.describe()
18 print(summary_table.to_latex())
19
20 # Check missing values
21 df.isnull().sum()
22
23 # Remove duplicates
24 df = df.drop_duplicates()
25 print("done")
26
27 # Save cleaned dataset:
28 df.to_csv('cleaned_data.csv', index=False)
29 files.download("cleaned_data.csv")
30 print("Saved")
31
32 # Basic Visualization
33 import matplotlib.pyplot as plt
34 df.hist(figsize=(10, 8))
35 plt.title("Distribution Analysis")
36 plt.savefig("distribution_analysis.png", dpi=300,
37             bbox_inches='tight')
38 plt.show()
39 files.download("distribution_analysis.png")
```