

Goals and Project Structure

✓ Setting up the Project

✓ Mount Drive

```
1 # mount to google drive
2 # from google.colab import drive
3 # drive.mount('/content/drive')
```

✓ Import Libraries

- will be importing more libraries as they become relevant in the project

```
1 pip install --upgrade numba ydata-profiling visions
```

```
Requirement already satisfied: numba in /usr/local/lib/python3.11/dist-packages (0.61.0)
Requirement already satisfied: ydata-profiling in /usr/local/lib/python3.11/dist-packages (4.12.2)
Requirement already satisfied: visions in /usr/local/lib/python3.11/dist-packages (0.7.6)
Collecting visions
  Using cached visions-0.8.1-py3-none-any.whl.metadata (11 kB)
Requirement already satisfied: llvmlite<0.45,>=0.44.0dev0 in /usr/local/lib/python3.11/dist-packages (from numba) (0.44.0)
Requirement already satisfied: numpy<2.2,>=1.24 in /usr/local/lib/python3.11/dist-packages (from numba) (1.26.4)
Requirement already satisfied: scipy<1.16,>=1.4.1 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (1.13.1)
Requirement already satisfied: pandas!=1.4.0,<3,>1.1 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (2.2.2)
Requirement already satisfied: matplotlib>=3.5 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (3.10.0)
Requirement already satisfied: pydantic>=2 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (2.10.6)
Requirement already satisfied: PyYAML<6.1,>=5.0.0 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (6.0.2)
Requirement already satisfied: jinja2<3.2,>=2.11.1 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (3.1.5)
Requirement already satisfied: htmlmin==0.1.12 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (0.1.12)
Requirement already satisfied: phik<0.13,>=0.11.1 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (0.12.4)
Requirement already satisfied: requests<3,>=2.24.0 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (2.32.3)
Requirement already satisfied: tqdm<5,>=4.48.2 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (4.67.1)
Requirement already satisfied: seaborn<0.14,>=0.10.1 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (0.13.2)
Requirement already satisfied: multimethod<2,>=1.4 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (1.12)
Requirement already satisfied: statsmodels<1,>=0.13.2 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (0.14.4)
Requirement already satisfied: typeguard<5,>=3 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (4.4.1)
Requirement already satisfied: imagehash==4.3.1 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (4.3.1)
Requirement already satisfied: wordcloud>=1.9.3 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (1.9.4)
Requirement already satisfied: dacite>=1.8 in /usr/local/lib/python3.11/dist-packages (from ydata-profiling) (1.9.2)
Requirement already satisfied: PyWavelets in /usr/local/lib/python3.11/dist-packages (from imagehash==4.3.1->ydata-profiling) (1.8.0)
Requirement already satisfied: pillow in /usr/local/lib/python3.11/dist-packages (from imagehash==4.3.1->ydata-profiling) (11.1.0)
Requirement already satisfied: attrs>=19.3.0 in /usr/local/lib/python3.11/dist-packages (from visions) (25.1.0)
Requirement already satisfied: networkx>=2.4 in /usr/local/lib/python3.11/dist-packages (from visions) (3.4.2)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from jinja2<3.2,>=2.11.1->ydata-profiling) (3
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (1.3
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (4.5
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (1.4
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (24.2
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (3.2.
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.11/dist-packages (from matplotlib>=3.5->ydata-profiling) (
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas!=1.4.0,<3,>1.1->ydata-profiling) (20
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas!=1.4.0,<3,>1.1->ydata-profiling) (
Requirement already satisfied: joblib>=0.14.1 in /usr/local/lib/python3.11/dist-packages (from phik<0.13,>=0.11.1->ydata-profiling) (1.4
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.11/dist-packages (from pydantic>=2->ydata-profiling) (0.
Requirement already satisfied: pydantic-core==2.27.2 in /usr/local/lib/python3.11/dist-packages (from pydantic>=2->ydata-profiling) (2.2
Requirement already satisfied: typing-extensions>=4.12.2 in /usr/local/lib/python3.11/dist-packages (from pydantic>=2->ydata-profiling)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.24.0->ydata-pro
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.24.0->ydata-profiling) (3.10
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.24.0->ydata-profiling)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.24.0->ydata-profiling)
Requirement already satisfied: patsy>=0.5.6 in /usr/local/lib/python3.11/dist-packages (from statsmodels<1,>=0.13.2->ydata-profiling) (1
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.7->matplotlib>=3.5->ydata-pr
```

```
1 import numpy as np, pandas as pd, matplotlib.pyplot as plt, seaborn as sns
2 from scipy import stats
3 from ydata_profiling import ProfileReport
```

Load Dataset

```
1 url = '/content/drive/MyDrive/Data_Science_Projects/Loan Default/Loan_default.csv'  
2 df = pd.read_csv(url)  
3 df_copy = df.copy()  
4 # quick glance at the data and ensure it is loaded properly  
5 df.head()
```

→

	LoanID	Age	Income	LoanAmount	CreditScore	MonthsEmployed	NumCreditLines	InterestRate	LoanTerm	DTIRatio	Education	EmploymentType	MaritalStatus	HasMortgage	HasDependents	LoanPurpose	HasCoSigner	Default
0	I38PQUQS96	56	85994	50587	520	80	4	15.23	36	0.44	Bachelor's	Employed	Married	Yes	No	Personal	No	
1	HPSK72WA7R	69	50432	124440	458	15	1	4.81	60	0.68	Master's	Employed	Married	Yes	No	Business/Finance	No	
2	C1OZ6DPJ8Y	46	84208	129188	451	26	3	21.17	24	0.31	Master's	Employed	Married	Yes	No	Business/Finance	No	
3	V2KKSFM3UN	32	31713	44799	743	0	3	7.07	24	0.23	High School	Employed	Married	Yes	No	Business/Finance	No	
4	EY08JDHTZP	60	20437	9139	633	8	4	6.51	48	0.73	Bachelor's	Employed	Married	Yes	No	Business/Finance	No	

Inspecting the Data

Shape

```
1 shape = df.shape  
2 print(f'The dataset has {shape[0]} rows and {shape[1]} columns.')
```

→ The dataset has 255347 rows and 18 columns.

Data Types

```
1 df.info()  
  
→ <class 'pandas.core.frame.DataFrame'>  
RangeIndex: 255347 entries, 0 to 255346  
Data columns (total 18 columns):  
 #   Column      Non-Null Count  Dtype     
---  --          --          --       --  
 0   LoanID      255347 non-null object    
 1   Age         255347 non-null int64     
 2   Income      255347 non-null int64     
 3   LoanAmount   255347 non-null int64     
 4   CreditScore  255347 non-null int64     
 5   MonthsEmployed  255347 non-null int64     
 6   NumCreditLines  255347 non-null int64     
 7   InterestRate  255347 non-null float64   
 8   LoanTerm     255347 non-null int64     
 9   DTIRatio     255347 non-null float64   
 10  Education    255347 non-null object    
 11  EmploymentType  255347 non-null object    
 12  MaritalStatus  255347 non-null object    
 13  HasMortgage   255347 non-null object    
 14  HasDependents  255347 non-null object    
 15  LoanPurpose   255347 non-null object    
 16  HasCoSigner   255347 non-null object    
 17  Default       255347 non-null int64     
dtypes: float64(2), int64(8), object(8)  
memory usage: 35.1+ MB
```

Missing Values

```
1 missing_vals = df.isna().sum()  
2  
3 if missing_vals.any():  
4   print('There are missing values in the dataset.')  
5 else:  
6   print('There are no missing values in the dataset.')
```

```
7  
8 missing_vals
```

→ There are no missing values in the dataset.

	0
LoanID	0
Age	0
Income	0
LoanAmount	0
CreditScore	0
MonthsEmployed	0
NumCreditLines	0
InterestRate	0
LoanTerm	0
DTIRatio	0
Education	0
EmploymentType	0
MaritalStatus	0
HasMortgage	0
HasDependents	0
LoanPurpose	0
HasCoSigner	0
Default	0

dtype: int64

▼ Duplicate Values

```
1 duplicate_vals = df.duplicated().sum()  
2  
3 if duplicate_vals > 0:  
4     print(f'There are {duplicate_vals} duplicate rows in the dataset.')  
5 else:  
6     print('There are no duplicate rows in the dataset.')
```

→ There are no duplicate rows in the dataset.

▼ Numerical Data

```
1 df.describe()
```

	Age	Income	LoanAmount	CreditScore	MonthsEmployed	NumCreditLines	InterestRate	LoanTerm	DTIRatio
count	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000	255347.000000
mean	43.498306	82499.304597	127578.865512	574.264346	59.541976	2.501036	13.492773	36.025894	0.500000
std	14.990258	38963.013729	70840.706142	158.903867	34.643376	1.117018	6.636443	16.969330	0.230000
min	18.000000	15000.000000	5000.000000	300.000000	0.000000	1.000000	2.000000	12.000000	0.100000
25%	31.000000	48825.500000	66156.000000	437.000000	30.000000	2.000000	7.770000	24.000000	0.300000
50%	43.000000	82466.000000	127556.000000	574.000000	60.000000	2.000000	13.460000	36.000000	0.500000
75%	56.000000	116219.000000	188985.000000	712.000000	90.000000	3.000000	19.250000	48.000000	0.700000
max	69.000000	149999.000000	249999.000000	849.000000	119.000000	4.000000	25.000000	60.000000	0.900000

Looking at our means we see on average:

- Most customers are ~ 44 years old and make ~\$82,500 a year
- The average loan amount per customer is \$127,578 for a term of ~17 months at an interest rate of approximately 13.5%
- The average credit score of a customer is 574 with a 0.5 DTI ratio

▼ Categorical Data

```
1 df.describe(include='object')
```

	LoanID	Education	EmploymentType	MaritalStatus	HasMortgage	HasDependents	LoanPurpose	HasCoSigner	
count	255347	255347	255347	255347	255347	255347	255347	255347	
unique	255347	4	4	3	2	2	5	2	
top	ZTH91CGL0B	Bachelor's	Part-time	Married	Yes	Yes	Business	Yes	
freq	1	64366	64161	85302	127677	127742	51298	127701	

From looking at our categorical data we see that:

- ~25% of customers have a Bachelor's degree, with this being the most common level of education achieved by the customer
- ~25% of customers also work part time
- ~50% of customers have a mortgage, a dependent, and/or a cosigner

▼ Detecting Outliers

- First we find the Interquartile Range (middle 50% of our data)
- Using the IQR we determine our upper and lower bounds
 - LowerBound = $Q1 - 1.5 \times IQR$
 - UpperBound = $Q3 + 1.5 \times IQR$
- Any value **below** the *lower bound* and **above** the *upper bound* is an outlier

```
1 # exclude the 'Default' column
2 # 'Default' column uses 1, 0 instead of yes, no
3 # not needed in outlier detection
4 df_no_default = df.drop('Default', axis=1)
5
6 def find_outliers(df):
7     outlier_count = {}
8     for col in df.select_dtypes(include='number'):
9         Q1 = df[col].quantile(0.25)
10        Q3 = df[col].quantile(0.75)
11        IQR = Q3 - Q1
12        lower_bound = Q1 - 1.5 * IQR
13        upper_bound = Q3 + 1.5 * IQR
14        outliers = df[(df[col] < lower_bound) | (df[col] > upper_bound)]
15        outlier_count[col] = outliers.shape[0]
16    return outlier_count
17
18 outliers_iqr = find_outliers(df_no_default)
19 outliers_df = pd.DataFrame(list(outliers_iqr.items()), columns=['Column', 'Outlier Count'])
20 outliers_df
```

	Column	Outlier Count	
0	Age	0	
1	Income	0	
2	LoanAmount	0	
3	CreditScore	0	
4	MonthsEmployed	0	
5	NumCreditLines	0	
6	InterestRate	0	
7	LoanTerm	0	
8	DTIRatio	0	

Next steps: [Generate code with outliers_df](#) [View recommended plots](#) [New interactive sheet](#)

Exploratory Data Analysis

```
1 eda = ProfileReport(df, title = 'EDA REPORT ')
2 eda
```

Summarize dataset: 100% 76/76 [00:21<00:00, 4.48it/s, Completed]

Generate report structure: 100% 1/1 [00:04<00:00, 4.65s/it]

Render HTML: 100% 1/1 [00:01<00:00, 1.51s/it]

EDA REPORT

Overview Variables Interactions Correlations Missing values Sample

Overview

Brought to you by [YData](#)

Overview	Alerts 1	Reproduction
Dataset statistics		Variable types
Number of variables 18		Text 1
Number of observations 255347		Numeric 7
Missing cells 0		Categorical 7
Missing cells (%) 0.0%		Boolean 3
Duplicate rows 0		
Duplicate rows (%) 0.0%		
Total size in memory 35.1 MiB		
Average record size in memory 144.0 B		

Variables

Select Columns

Some Observations:

- Our data is not normally distributed, rather most if it is roughly uniform, which means we need to use non-parametric models
- Our data also has low correlation, no one feature is affected by the other

Normality Check

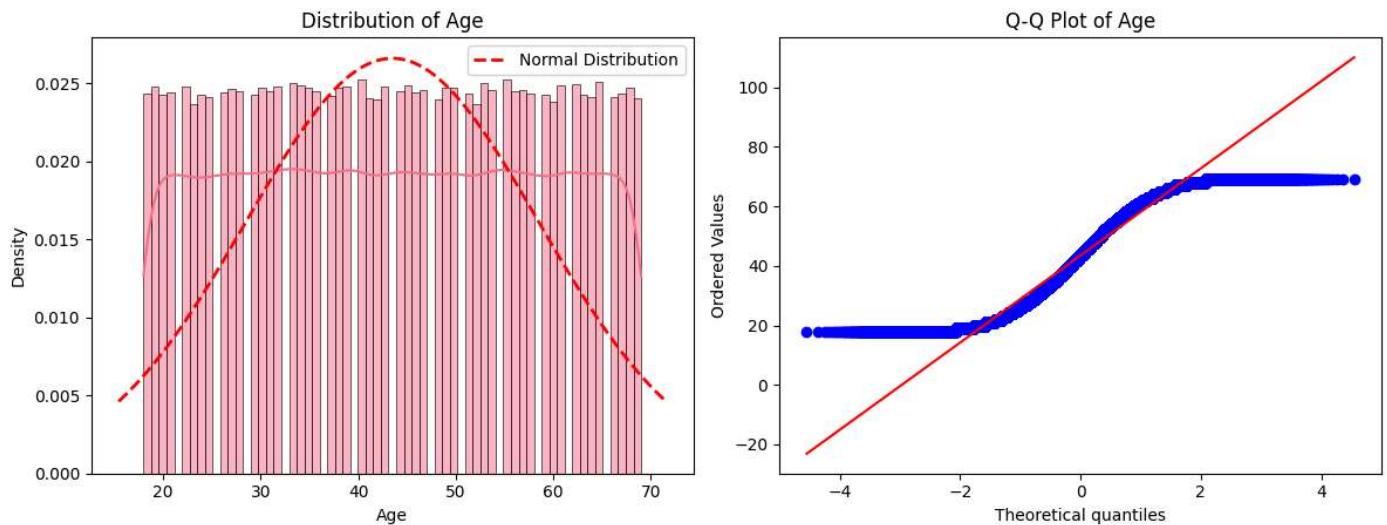
```
1 def plot_normality_check(data, feature, figsize=(12, 5)):  
2     """  
3         Create side-by-side plots to check normality of a numerical feature:  
4         1. Histogram with KDE and normal distribution overlay  
5         2. Q-Q plot  
6  
7     Parameters:  
8     -----  
9     data : pandas.DataFrame  
10        The dataframe containing the feature  
11    feature : str  
12        The name of the feature to analyze  
13    figsize : tuple, optional  
14        Figure size for the plot (width, height)  
15  
16    Returns:  
17    -----  
18    None: Displays the plots  
19    """  
20  
21    # Create figure and axes for side-by-side plots  
22    fig, (ax1, ax2) = plt.subplots(1, 2, figsize=figsize)  
23  
24    # Get the data for the feature  
25    feature_data = data[feature].dropna()  
26  
27    # Plot 1: Histogram with KDE  
28    sns.histplot(data=feature_data, stat='density', kde=True, ax=ax1)  
29  
30    # Add normal distribution curve  
31    xmin, xmax = ax1.get_xlim()  
32    x = np.linspace(xmin, xmax, 100)  
33    mu, std = feature_data.mean(), feature_data.std()  
34    p = stats.norm.pdf(x, mu, std)  
35    ax1.plot(x, p, 'r--', linewidth=2, label='Normal Distribution')  
36  
37    # Customize first plot  
38    ax1.set_title(f'Distribution of {feature}')  
39    ax1.set_xlabel(feature)  
40    ax1.set_ylabel('Density')  
41    ax1.legend()  
42  
43    # Plot 2: Q-Q plot  
44    stats.probplot(feature_data, dist="norm", plot=ax2)  
45    ax2.set_title(f'Q-Q Plot of {feature}')  
46  
47    # Add overall title and adjust layout  
48    plt.suptitle(f'Normality Check for {feature}', y=1.05)  
49    plt.tight_layout()  
50  
51    # Show plot  
52    plt.show()  
53  
54    # Print Shapiro-Wilk test results  
55    stat, p_value = stats.shapiro(feature_data.sample(min(5000, len(feature_data))))  
56    print(f'\nShapiro-Wilk Test Results for {feature}:')  
57    print(f'Staticstic: {stat:.4f}')  
58    print(f'p-value: {p_value:.4f}')  
59    if p_value < 0.05:  
60        print('The feature is not normally distributed (rejects null hypothesis)')  
61    else:  
62        print('The feature appears to be normally distributed (fails to reject null hypothesis)')  
63  
64    # Print skewness and kurtosis  
65    print(f'\nSkewness: {stats.skew(feature_data):.4f}')
```

```
66     print(f'Kurtosis: {stats.kurtosis(feature_data):.4f}')
67
68 # Example usage for multiple features:
69 def check_all_numeric_features(data, exclude_cols=None):
70     """
71     Check normality for all numeric features in the dataset
72
73     Parameters:
74     -----
75     data : pandas.DataFrame
76         The dataframe to analyze
77     exclude_cols : list, optional
78         List of column names to exclude from analysis
79
80     Returns:
81     -----
82     None: Displays plots for each numeric feature
83     """
84
85     # Get numeric columns
86     numeric_cols = data.select_dtypes(include=['int64', 'float64']).columns
87
88     # Exclude specified columns
89     if exclude_cols:
90         numeric_cols = [col for col in numeric_cols if col not in exclude_cols]
91
92     # Plot normality check for each numeric feature
93     for feature in numeric_cols:
94         plot_normality_check(data, feature)
95         print('\n' + '='*50 + '\n')

1 numeric_cols = df.select_dtypes(include=['int64', 'float64'])
2
3 check_all_numeric_features(numeric_cols, exclude_cols=['LoanID', 'Default', 'NumCreditLines', 'LoanTerm'])
```



Normality Check for Age

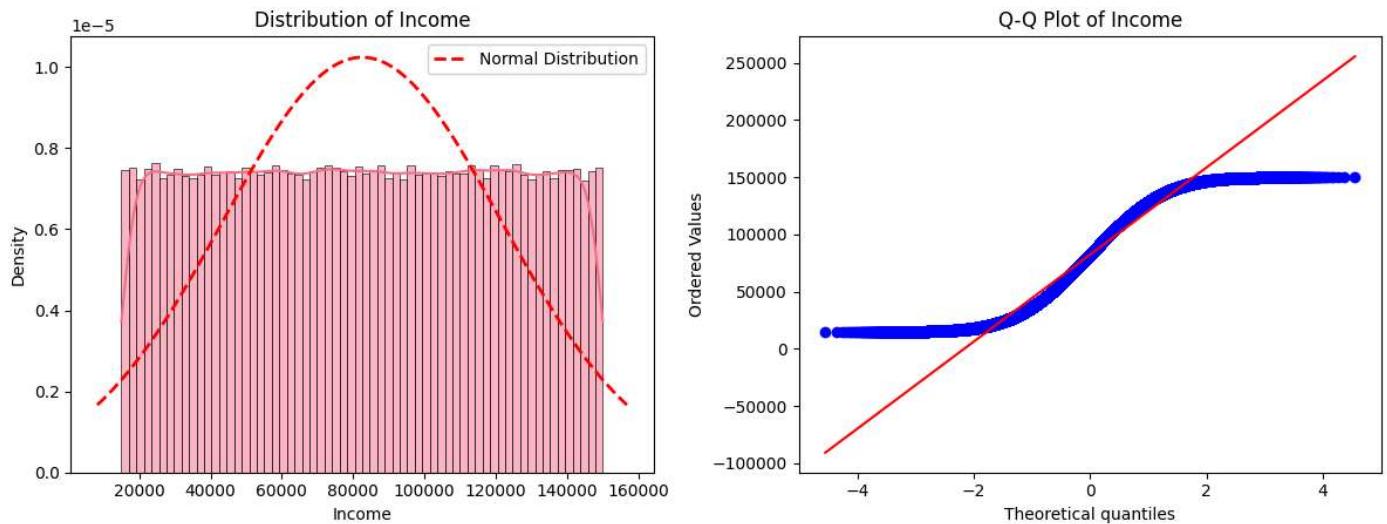


Shapiro-Wilk Test Results for Age:
Statistic: 0.9528
p-value: 0.0000
The feature is not normally distributed (rejects null hypothesis)

Skewness: 0.0007
Kurtosis: -1.1984

=====

Normality Check for Income

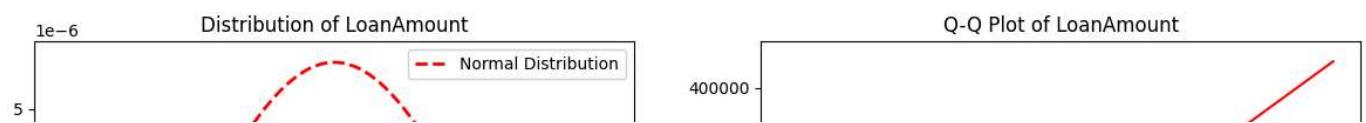


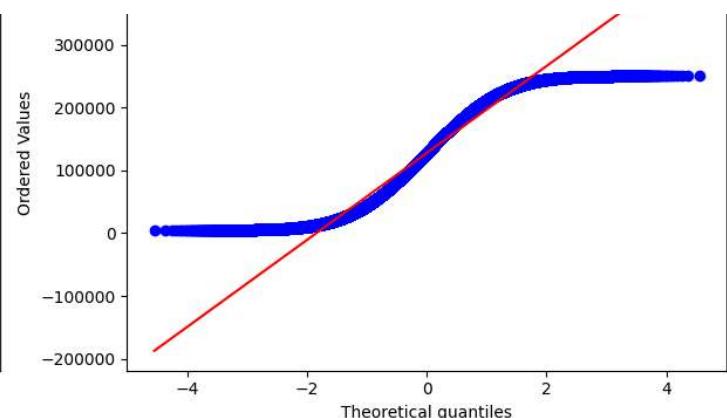
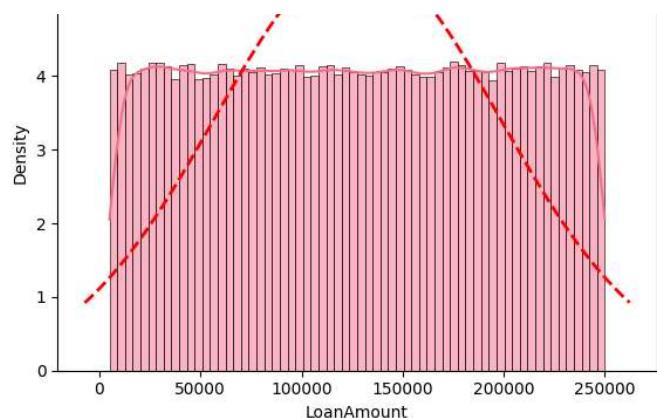
Shapiro-Wilk Test Results for Income:
Statistic: 0.9559
p-value: 0.0000
The feature is not normally distributed (rejects null hypothesis)

Skewness: -0.0004
Kurtosis: -1.1984

=====

Normality Check for LoanAmount

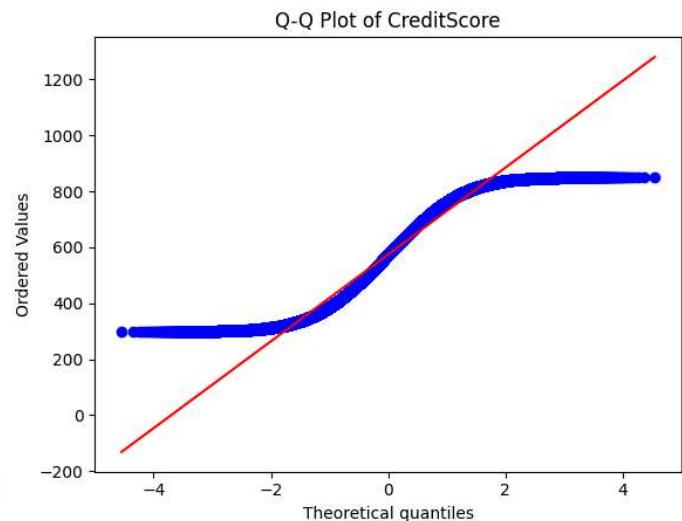
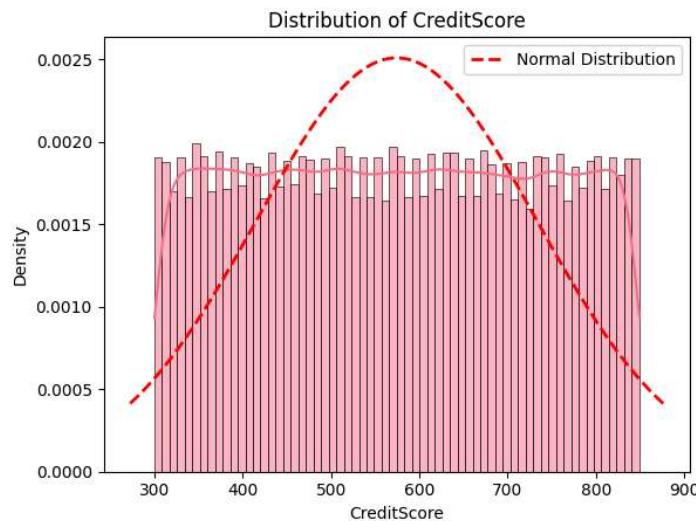




Shapiro-Wilk Test Results for LoanAmount:
 Statistic: 0.9550
 p-value: 0.0000
 The feature is not normally distributed (rejects null hypothesis)

Skewness: -0.0018
 Kurtosis: -1.2037

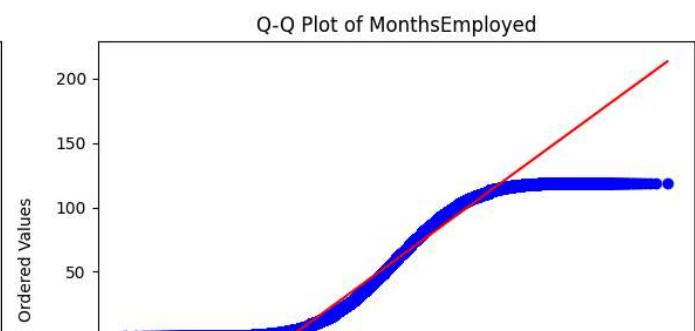
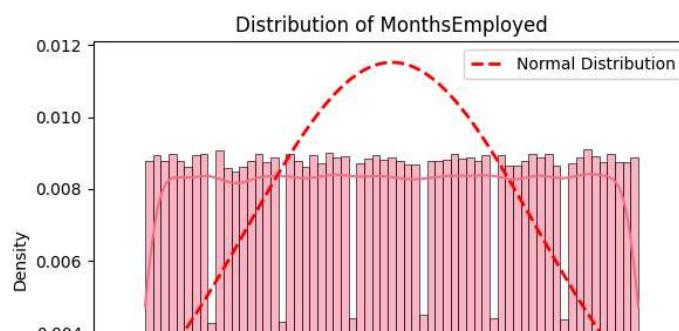
Normality Check for CreditScore

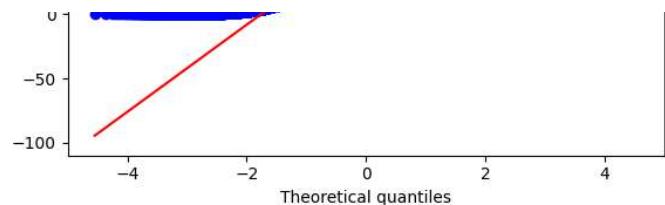
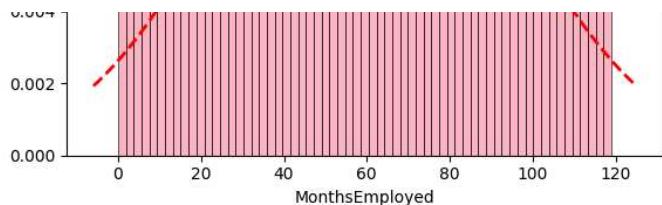


Shapiro-Wilk Test Results for CreditScore:
 Statistic: 0.9526
 p-value: 0.0000
 The feature is not normally distributed (rejects null hypothesis)

Skewness: 0.0047
 Kurtosis: -1.2003

Normality Check for MonthsEmployed

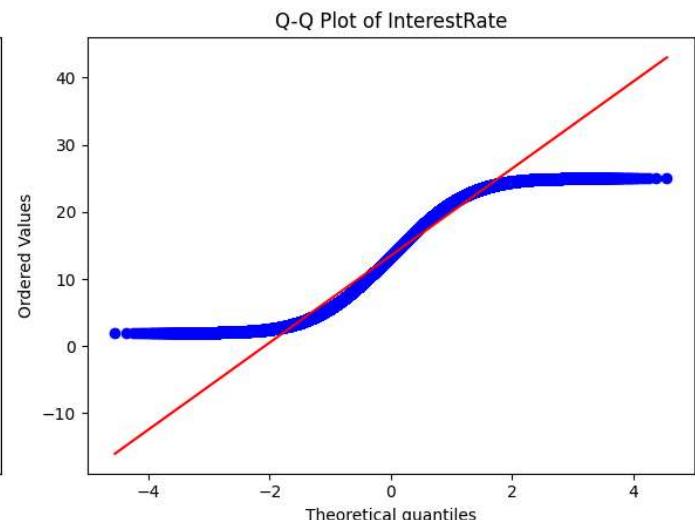
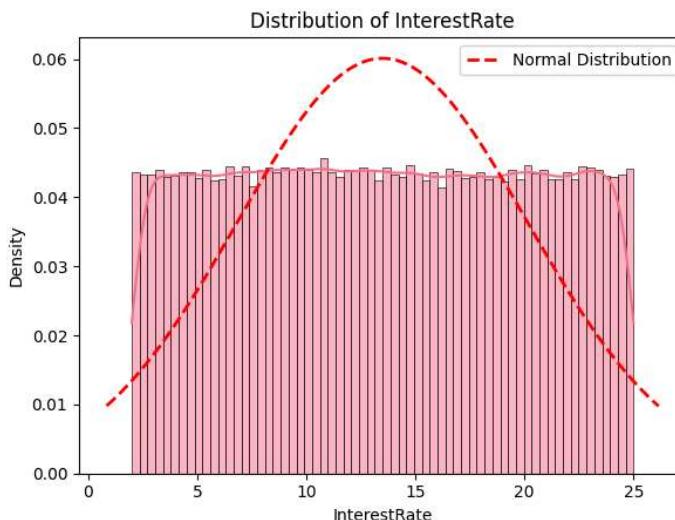




Shapiro-Wilk Test Results for MonthsEmployed:
 Statistic: 0.9531
 p-value: 0.0000
 The feature is not normally distributed (rejects null hypothesis)

Skewness: -0.0021
 Kurtosis: -1.1996

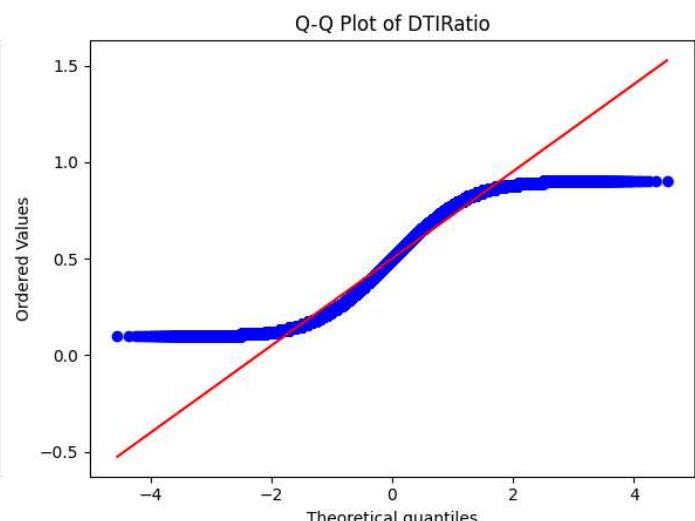
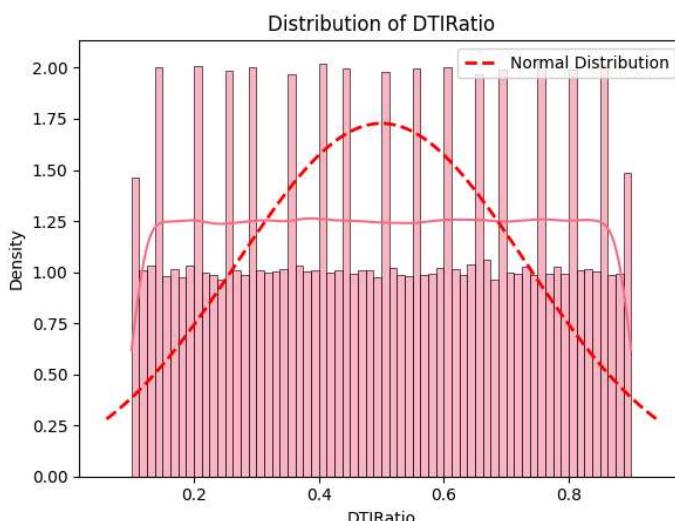
Normality Check for InterestRate



Shapiro-Wilk Test Results for InterestRate:
 Statistic: 0.9566
 p-value: 0.0000
 The feature is not normally distributed (rejects null hypothesis)

Skewness: 0.0046
 Kurtosis: -1.1972

Normality Check for DTIRatio



Shapiro-Wilk Test Results for DTIRatio:
Statistic: 0.9572
p-value: 0.0000
The feature is not normally distributed (rejects null hypothesis)

Skewness: -0.0015
Kurtosis: -1.1997

=====

We can see here that the p-value for each column is zero, right away this tells us that our data is not normal in any numerical column. Similarly looking at the graphs, we see our histogram data does not match the bell-curve of normalised data, and similarly our data does not sit neatly on the trend lines in our Q-Q plot, all signs that we are not working with normalised data.

From here we can apply transformations to our data to make it more optimized for models that rely on normal data. Since we do not need our data to be normal in the analysis phase, we can simply continue with the data in its current state, and later when choosing a model to train we can use one that does not rely on normalised data.

✓ Analyzing Risk

```
1 # Create new features
2 df['LoanToIncomeRatio'] = df['LoanAmount'] / df['Income']
3 df['RiskScore'] = (df['DTIRatio'] / df['CreditScore']) * 1000
4 df['MonthlyLoanBurden'] = (df['LoanAmount'] * (df['InterestRate']/1200)) / (1 - (1 + df['InterestRate']/1200)**(-df['LoanTerm']))
5 df['PaymentToIncome_Ratio'] = (df['MonthlyLoanBurden'] * 100) / (df['Income'] / 12)
```

✓ Risk Score and Loan-ToIncome Ratio

```
1 # Create bands for plot
2 df['RiskBand'] = pd.qcut(df['RiskScore'], q=5, labels=['Very Low', 'Low', 'Medium', 'High', 'Very High'])
3 df['CreditScoreBand'] = pd.cut(df['CreditScore'],
4                                bins=[300, 400, 500, 600, 700, 800, 850],
5                                labels=['Very Poor', 'Poor', 'Fair', 'Good', 'Very Good', 'Excellent'])
6
7 # Create DTI and LTI bins
8 dti_bins = pd.qcut(df['DTIRatio'], q=5)
9 df['LoanToIncomeRatio'] = df['LoanAmount'] / df['Income']
10 lti_bins = pd.qcut(df['LoanToIncomeRatio'], q=5)
11
12 print("\nDTI Ratio Ranges:")
13 print("-" * 50)
14 dti_ranges = [f"{interval.left:.1%}-{interval.right:.1%}"
15               for interval in dti_bins.unique().categories]
16 for i, interval in enumerate(dti_bins.unique().categories, 1):
17     print(f"Band {i}: {interval.left:.1%} to {interval.right:.1%}")
18
19 print("\nLoan-to-Income Ratio Ranges:")
20 print("-" * 50)
21 lti_ranges = [f"{interval.left:.2f}-{interval.right:.2f}x"
22                for interval in lti_bins.unique().categories]
23 for i, interval in enumerate(lti_bins.unique().categories, 1):
24     print(f"Band {i}: {interval.left:.2f}x to {interval.right:.2f}x annual income")
25
26 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 8))
27
28 # 1. Default Rate by Risk Band and DTI
29 risk_dti_default = df.groupby(['RiskBand', dti_bins], observed=False)[['Default']].mean().unstack()
30 risk_dti_default.columns = dti_ranges
31
32 sns.heatmap(risk_dti_default, annot=True, fmt='%.2%', cmap='YlOrRd', ax=ax1)
33 ax1.set_title('Default Rate by Risk Band and DTI Range', pad=20)
34 ax1.set_xlabel('DTI Ratio Range')
35 ax1.set_ylabel('Potential of Risk')
36 ax1.set_xticklabels(ax1.get_xticklabels(), rotation=45, ha='right')
37 ax1.set_yticklabels(ax1.get_yticklabels(), rotation=0)
38
39 # 2. Loan-to-Income Ratio vs Default Rate by Credit Score Band
40 lti_default = df.groupby(['CreditScoreBand', lti_bins], observed=False)[['Default']].mean().unstack()
41 lti_default.columns = lti_ranges
42
43 sns.heatmap(lti_default, annot=True, fmt='%.2%', cmap='YlOrRd', ax=ax2)
44 ax2.set_title('Default Rate by Credit Score and\nLoan-to-Income Ratio', pad=20)
45 ax2.set_xlabel('Loan-to-Income Ratio Range (x = Annual Income)')
46 ax2.set_ylabel('Credit Score Percentile')
47 ax2.set_xticklabels(ax2.get_xticklabels(), rotation=45, ha='right')
48 ax2.set_yticklabels(ax2.get_yticklabels(), rotation=0)
49
50 plt.tight_layout()
51 plt.show()
52
53 # Calculate Risk Band Statistics
```

```
54 risk_band_stats = df.groupby('RiskBand', observed=False)['Default'].agg([
55     ('count', 'count'),
56     ('default_rate', 'mean'),
57     ('num_defaults', lambda x: x.sum())
58 ]).round(4)
59
60 # Convert default rate to percentage
61 risk_band_stats['default_rate'] = (risk_band_stats['default_rate'] * 100).round(2)
62
63 # Add percentage of total loans column
64 risk_band_stats['percent_of_total'] = (risk_band_stats['count'] / risk_band_stats['count'].sum() * 100).round(2)
65
66 # Print Risk Band Statistics
67 print("\nRisk Band Statistics:")
68 print("-" * 50)
69 print(f"\nTotal number of loans: {df.shape[0]}")
70 print("\nBreakdown by Risk Band:")
71 print(risk_band_stats)
72 print("\nColumns:")
73 print("- count: Number of loans in this risk band")
74 print("- default_rate: Percentage of loans that defaulted in this band")
75 print("- num_defaults: Number of defaults in this band")
76 print("- percent_of_total: Percentage of total loans in this band")
```

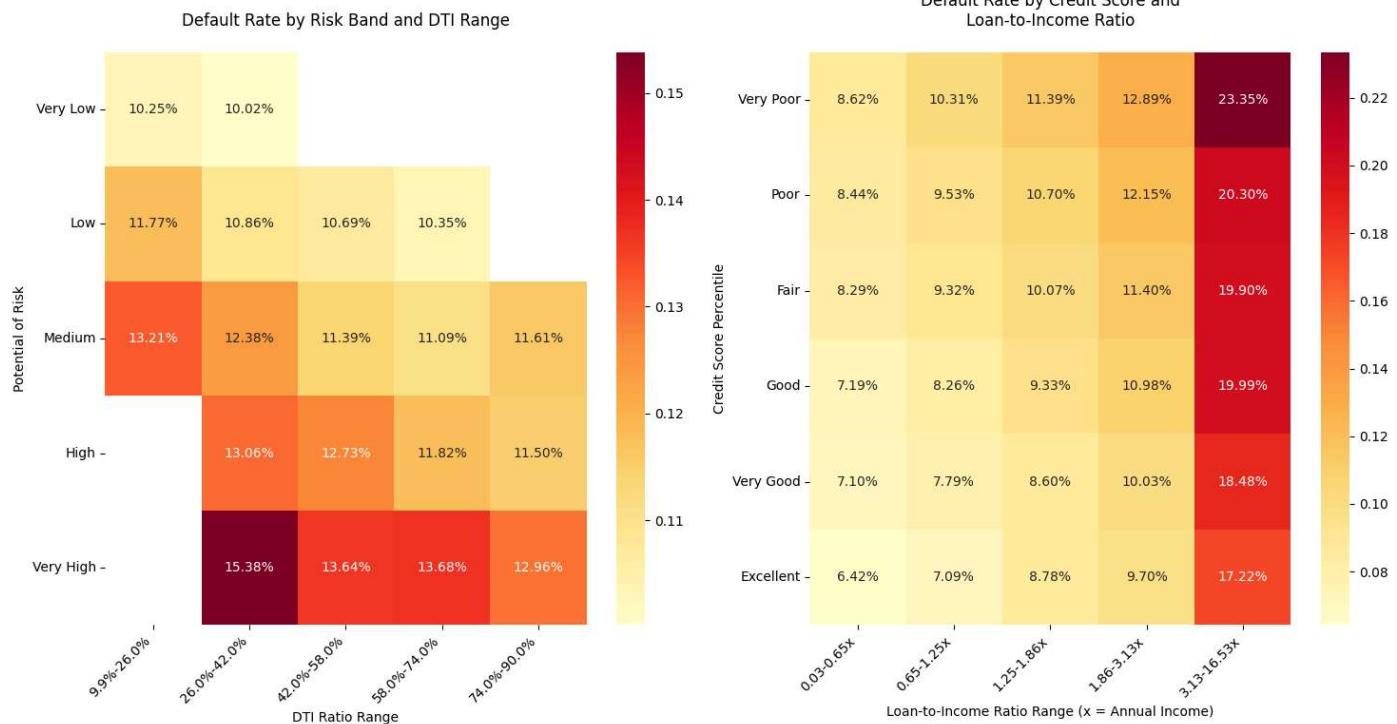


DTI Ratio Ranges:

Band 1: 9.9% to 26.0%
 Band 2: 26.0% to 42.0%
 Band 3: 42.0% to 58.0%
 Band 4: 58.0% to 74.0%
 Band 5: 74.0% to 90.0%

Loan-to-Income Ratio Ranges:

Band 1: 0.03x to 0.65x annual income
 Band 2: 0.65x to 1.25x annual income
 Band 3: 1.25x to 1.86x annual income
 Band 4: 1.86x to 3.13x annual income
 Band 5: 3.13x to 16.53x annual income



Risk Band Statistics:

Total number of loans: 255347

Breakdown by Risk Band:

RiskBand	count	default_rate	num_defaults	percent_of_total
Very Low	51075	10.21	5213	20.0
Low	51064	10.98	5607	20.0
Medium	51072	11.56	5905	20.0
High	51069	12.02	6140	20.0
Very High	51067	13.29	6788	20.0

Columns:

- count: Number of loans in this risk band
- default_rate: Percentage of loans that defaulted in this band
- num_defaults: Number of defaults in this band
- percent_of_total: Percentage of total loans in this band

Here we calculate the Potential of Risk as follows:

$$\text{Potential} = \frac{\text{DTI Ratio}}{\text{Credit Score}} \times 1000$$

We see that:

- Just like we predicted, customers with a higher risk score are way more likely to default
- Although customers with a high DTI ratio (58% +) are usually more likely to default, the customers within the range of (26%-42%) show the highest default rates across all risk potentials
- This tells us that using DTI ratio to predict which customers will default may not be as predictable as other factors, like credit score and Loan-To-Income Ratio that show clear trends
 - (Higher credit score and lower LTI ratio has less chance of defaulting and vice versa)

```

1 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))
2
3 # Plot 1: DTI Ratio vs Default Rate
4 # Create 20 bins of DTI Ratio and calculate mean default rate for each bin
5 dti_bins = pd.qcut(df['DTIRatio'], q=20)
6 default_by_dti = df.groupby(dti_bins, observed=False)[['Default']].mean() * 100 # Convert to percentage
7
8 # Plot the DTI ratio vs default rate
9 default_by_dti.plot(kind='line', marker='o', color='blue', ax=ax1)
10 ax1.set_title('DTI Ratio vs Default Rate', pad=20)
11 ax1.set_xlabel('DTI Ratio (binned)')
12 ax1.set_ylabel('Default Rate (%)')
13 ax1.tick_params(axis='x', rotation=45)
14 ax1.grid(True, alpha=0.3)
15
16 # Plot 2: Credit Score vs Default Rate
17 # Create 20 bins of Credit Score and calculate mean default rate for each bin
18 score_bins = pd.qcut(df['CreditScore'], q=20)
19 default_by_score = df.groupby(score_bins, observed=False)[['Default']].mean() * 100 # Convert to percentage
20
21 # Plot the Credit Score vs default rate
22 default_by_score.plot(kind='line', marker='o', color='red', ax=ax2)
23 ax2.set_title('Credit Score vs Default Rate', pad=20)
24 ax2.set_xlabel('Credit Score (binned)')
25 ax2.set_ylabel('Default Rate (%)')
26 ax2.tick_params(axis='x', rotation=45)
27 ax2.grid(True, alpha=0.3)
28
29 # Adjust layout to prevent overlap
30 plt.tight_layout()
31
32 # Print some summary statistics
33 print("\nDefault Rate Statistics:")
34 print("-" * 50)
35 print("\nDTI Ratio Quartiles and Default Rates:")
36 dti_quartiles = pd.qcut(df['DTIRatio'], q=4)
37 dti_default_rates = df.groupby(dti_quartiles, observed=False)[['Default']].agg(['mean', 'count'])
38 dti_default_rates['mean'] = dti_default_rates['mean'] * 100
39 print(dti_default_rates)
40
41 print("\nCredit Score Quartiles and Default Rates:")
42 score_quartiles = pd.qcut(df['CreditScore'], q=4)
43 score_default_rates = df.groupby(score_quartiles, observed=False)[['Default']].agg(['mean', 'count'])
44 score_default_rates['mean'] = score_default_rates['mean'] * 100
45 print(score_default_rates)

```



Default Rate Statistics:

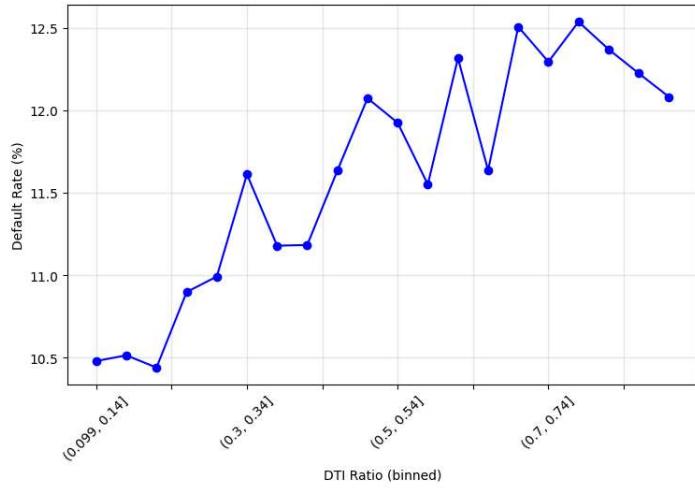
DTI Ratio Quartiles and Default Rates:

	mean	count
DTIRatio		
(0.099, 0.3]	10.659331	65248
(0.3, 0.5]	11.535095	63970
(0.5, 0.7]	11.987896	63781
(0.7, 0.9]	12.306730	62348

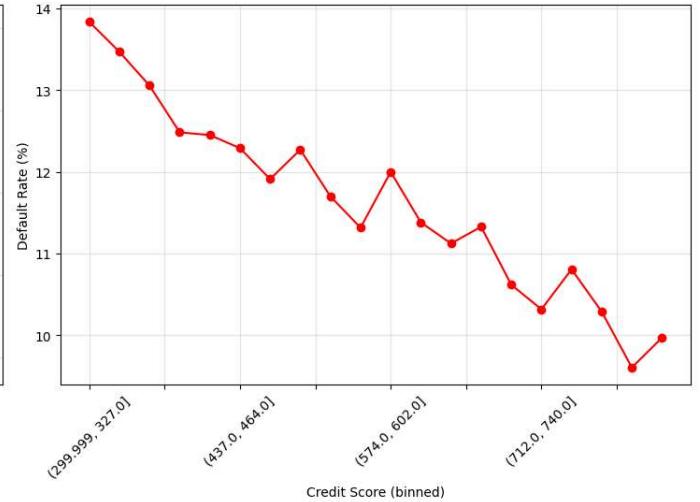
Credit Score Quartiles and Default Rates:

	mean	count
CreditScore		
(299.999, 437.0]	13.055815	64140
(437.0, 574.0]	11.894425	63803
(574.0, 712.0]	11.290222	63896
(712.0, 849.0]	10.197141	63508

DTI Ratio vs Default Rate



Credit Score vs Default Rate



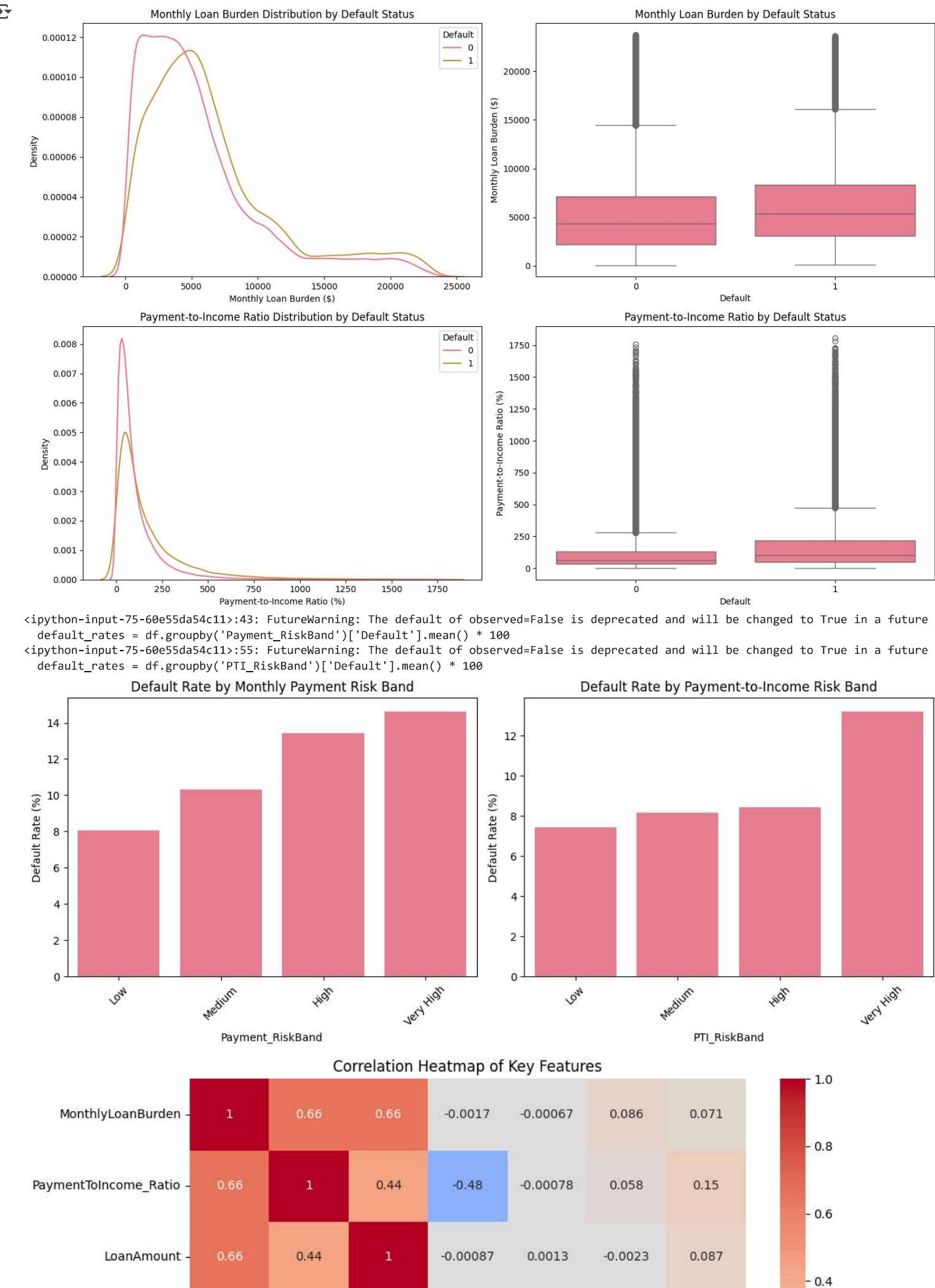
Plotting risk factors(DTI Ratio and Credit Score) by Default Rate we see clear trends indicating that a Higer DTI ratio and A lower credit score increases a customers chance of defaulting on their loan.

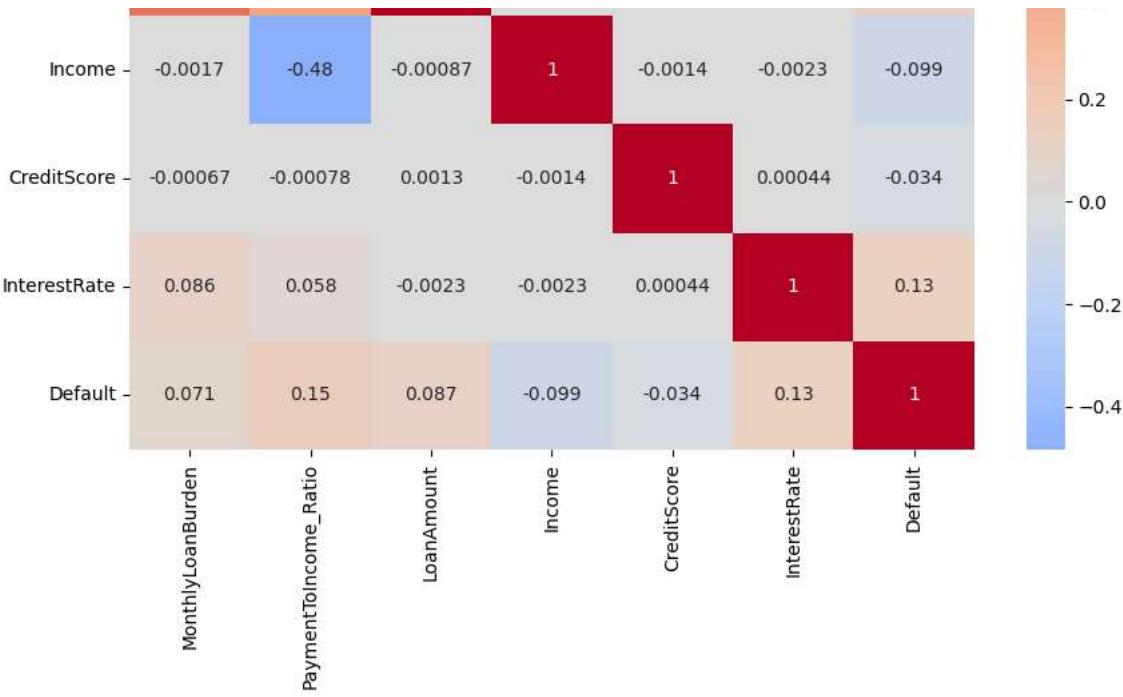
```
1 # Set up the plotting style
2 sns.set_palette("husl")
3
4 # Create figure with multiple subplots
5 fig = plt.figure(figsize=(15, 10))
6
7 # 1. Monthly Loan Burden Analysis
8 plt.subplot(2, 2, 1)
9 sns.kdeplot(data=df, x='MonthlyLoanBurden', hue='Default', common_norm=False)
10 plt.title('Monthly Loan Burden Distribution by Default Status')
11 plt.xlabel('Monthly Loan Burden ($)')
12 plt.ylabel('Density')
13
14 plt.subplot(2, 2, 2)
15 sns.boxplot(data=df, x='Default', y='MonthlyLoanBurden')
16 plt.title('Monthly Loan Burden by Default Status')
17 plt.ylabel('Monthly Loan Burden ($)')
18
19 # 2. Payment-to-Income Ratio Analysis
20 plt.subplot(2, 2, 3)
21 sns.kdeplot(data=df, x='PaymentToIncome_Ratio', hue='Default', common_norm=False)
22 plt.title('Payment-to-Income Ratio Distribution by Default Status')
23 plt.xlabel('Payment-to-Income Ratio (%)')
24 plt.ylabel('Density')
25
26 plt.subplot(2, 2, 4)
27 sns.boxplot(data=df, x='Default', y='PaymentToIncome_Ratio')
```

```

28 plt.title('Payment-to-Income Ratio by Default Status')
29 plt.ylabel('Payment-to-Income Ratio (%)')
30
31 plt.tight_layout()
32 plt.show()
33
34 # Create risk band analysis
35 plt.figure(figsize=(12, 5))
36
37 # Monthly Payment Risk Bands
38 payment_bins = df['MonthlyLoanBurden'].quantile([0, 0.25, 0.5, 0.75, 1.0])
39 df['Payment_RiskBand'] = pd.qcut(df['MonthlyLoanBurden'], q=4,
40                                     labels=['Low', 'Medium', 'High', 'Very High'])
41
42 plt.subplot(1, 2, 1)
43 default_rates = df.groupby('Payment_RiskBand')['Default'].mean() * 100
44 sns.barplot(x=default_rates.index, y=default_rates.values)
45 plt.title('Default Rate by Monthly Payment Risk Band')
46 plt.ylabel('Default Rate (%)')
47 plt.xticks(rotation=45)
48
49 # Payment-to-Income Risk Bands
50 pti_bins = [0, 20, 30, 40, float('inf')]
51 pti_labels = ['Low', 'Medium', 'High', 'Very High']
52 df['PTI_RiskBand'] = pd.cut(df['PaymentToIncome_Ratio'], bins=pti_bins, labels=pti_labels)
53
54 plt.subplot(1, 2, 2)
55 default_rates = df.groupby('PTI_RiskBand')['Default'].mean() * 100
56 sns.barplot(x=default_rates.index, y=default_rates.values)
57 plt.title('Default Rate by Payment-to-Income Risk Band')
58 plt.ylabel('Default Rate (%)')
59 plt.xticks(rotation=45)
60
61 plt.tight_layout()
62 plt.show()
63
64 # Create correlation heatmap
65 correlation_features = ['MonthlyLoanBurden', 'PaymentToIncome_Ratio',
66                         'LoanAmount', 'Income', 'CreditScore', 'InterestRate', 'Default']
67 correlation_matrix = df[correlation_features].corr()
68
69 plt.figure(figsize=(10, 8))
70 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', center=0)
71 plt.title('Correlation Heatmap of Key Features')
72 plt.tight_layout()
73 plt.show()
74
75 # Print summary statistics
76 print("\nMonthly Loan Burden Statistics:")
77 print("-" * 50)
78 print(df['MonthlyLoanBurden'].describe().round(2))
79
80 print("\nPayment-to-Income Ratio Statistics:")
81 print("-" * 50)
82 print(df['PaymentToIncome_Ratio'].describe().round(2))
83
84 # Print default rates by risk bands
85 print("\nDefault Rates by Payment Risk Band:")
86 print("-" * 50)
87 print(df.groupby('Payment_RiskBand')['Default'].agg(['count', 'mean']).round(4))
88
89 print("\nDefault Rates by Payment-to-Income Risk Band:")
90 print("-" * 50)
91 print(df.groupby('PTI_RiskBand')['Default'].agg(['count', 'mean']).round(4))

```





Monthly Loan Burden Statistics:

```

count    255347.00
mean     5649.06
std      4731.40
min      90.80
25%     2317.46
50%     4442.06
75%     7259.85
max     23735.67
Name: MonthlyLoanBurden, dtype: float64

```

Payment-to-Income Ratio Statistics:

```

count    255347.00
mean     115.73
std      146.64
min      0.82
25%     33.61
50%     68.25
75%     139.79
max     1808.61
Name: PaymentToIncome_Ratio, dtype: float64

```

Default Rates by Payment Risk Band:

Payment_RiskBand	count	mean
Low	63837	0.0806
Medium	63837	0.1032
High	63836	0.1343
Very High	63837	0.1464

Default Rates by Payment-to-Income Risk Band:

PTI_RiskBand	count	mean
Low	35993	0.0742
Medium	20325	0.0817
High	20750	0.0844
Very High	178279	0.1322

```

<ipython-input-75-60e55da54c11>:87: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future v
  print(df.groupby('Payment_RiskBand')['Default'].agg(['count', 'mean']).round(4))
<ipython-input-75-60e55da54c11>:91: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future v
  print(df.groupby('PTI_RiskBand')['Default'].agg(['count', 'mean']).round(4))

```


We calculate Monthly Loan Burden as follows:

$$\text{MonthlyLoanBurden} = \frac{\text{LoanAmount} \times \left(\frac{\text{InterestRate}}{1200} \right)}{1 - \left(1 + \frac{\text{InterestRate}}{1200} \right)^{-\text{LoanTerm}}}$$

- Using the standard amortization formula that ensures whether a loan is fully paid off by the end of its term or not based off of the loan amount, interest rate (converted to monthly), and loan term
- However looking at our kde plot, we can only gather that most customers have a monthly loan burden of ~\$5000, but this is true for both defaulted customers and those who paid off their loan

We calculate Payment to Income Ratio as:

$$\text{PaymentToIncome_Ratio} = \frac{\text{MonthlyLoanBurden} \times 100}{\frac{\text{Income}}{12}}$$

- Here we express monthly loan payment as a percentage of monthly income
- The higher the ratio is the more potential for risk
- Again our kde plot doesn't show much difference between defaulted customers and paid off customers

✓ ML Model Evaluation

- Here we will be testing the performance of various types of ML Models and choosing the best model for our data

✓ Pre Processing

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler, LabelEncoder
3 from sklearn.impute import SimpleImputer
4
5 def handle_missing_values(df):
6     """Impute missing values appropriately"""
7     data = df.copy()
8
9     # Separate numerical and categorical columns
10    numeric_columns = data.select_dtypes(include=['int64', 'float64']).columns
11    categorical_columns = data.select_dtypes(include=['object']).columns
12
13    # Impute numerical columns with median
14    numeric_imputer = SimpleImputer(strategy='median')
15    data[numeric_columns] = numeric_imputer.fit_transform(data[numeric_columns])
16
17    # Impute categorical columns with mode
18    categorical_imputer = SimpleImputer(strategy='most_frequent')
19    data[categorical_columns] = categorical_imputer.fit_transform(data[categorical_columns])
20
21    if data.isnull().values.any():
22        print("WARNING: There are still NaN values in the dataset after imputation!")
23        print("Investigate and fix the source of these missing values.")
24        print(data.isnull().sum()) # Print the columns with NaN values for debugging
25
26    return data
27
28    return data
29
30 def handle_outliers(df):
31     """Handle outliers in numerical columns"""
32     data = df.copy()
33
34     # List of columns to check for outliers
35     columns_to_check = ['Age', 'Income', 'LoanAmount', 'MonthsEmployed',
36                         'NumCreditLines', 'InterestRate', 'DTIRatio',
37                         'MonthlyLoanBurden', 'PaymentToIncome_Ratio',]
38
39     for column in columns_to_check:
40         Q1 = data[column].quantile(0.25)
41         Q3 = data[column].quantile(0.75)
42         IQR = Q3 - Q1
43         lower_bound = Q1 - 1.5 * IQR
```

```

44     upper_bound = Q3 + 1.5 * IQR
45     data[column] = data[column].clip(lower=lower_bound, upper=upper_bound)
46
47     return data
48
49 def encode_categorical_variables(df):
50     """Encode categorical variables appropriately"""
51
52     data = df.copy()
53
54     if not isinstance(df, pd.DataFrame):
55         raise TypeError("Input must be a pandas DataFrame")
56
57     # Binary encoding for yes/no variables
58     binary_columns = ['HasMortgage', 'HasDependents', 'HasCoSigner']
59
60     # Binary encoding
61     if binary_columns:
62         for col in binary_columns:
63             data[col] = data[col].map({'Yes': 1, 'No': 0})
64
65     # Direct mapping is the only strategy that works here
66
67     data['Education'] = data['Education'].map({"High School": 1,
68                                              "Bachelor's": 2,
69                                              "Master's": 3,
70                                              "PhD": 4})
71
72     data['EmploymentType'] = data['EmploymentType'].map({"Unemployed": 1,
73                                                       "Self-employed": 2,
74                                                       "Part-time": 3,
75                                                       "Full-time": 4})
76
77     data['MaritalStatus'] = data['MaritalStatus'].map({"Single": 1,
78                                                       "Married": 2,
79                                                       "Divorced": 3})
80
81     data['LoanPurpose'] = data['LoanPurpose'].map({"Auto": 1,
82                                                 "Business": 2,
83                                                 "Education": 3,
84                                                 "Home": 4,
85                                                 "Other": 5})
86
87     return data
88
89
90 # Step 5: Feature scaling
91 def scale_features(df):
92     """Scale numerical features"""
93     data = df.copy()
94
95     # Identify numerical columns to scale (exclude target variable and ID)
96     columns_to_scale = data.select_dtypes(include=['int64', 'float64']).columns
97     columns_to_scale = columns_to_scale.drop('Default') # Exclude target variable
98     if 'LoanID' in columns_to_scale:
99         columns_to_scale = columns_to_scale.drop('LoanID') # Exclude ID
100
101    # Scale the features
102    scaler = StandardScaler()
103    data[columns_to_scale] = scaler.fit_transform(data[columns_to_scale])
104
105    return data
106
107 # Main preprocessing pipeline
108 def preprocess_data(df, test_size=0.2, random_state=42):
109     """Complete preprocessing pipeline"""
110
111     print("Starting preprocessing pipeline...")
112
113     print("Handling missing values...")
114     df = handle_missing_values(df)
115
116     print("Handling outliers...")
117     df = handle_outliers(df)
118
119     print("Encoding categorical variables...")
120     df = encode_categorical_variables(df)

```

```

121
122     print("Scaling features...")
123     df = scale_features(df)
124
125     print("Splitting the data...")
126
127     # Remove LoanID if present
128     if 'LoanID' in df.columns:
129         df = df.drop('LoanID', axis=1)
130
131     # Remove Band columns
132     if 'RiskBand' in df.columns:
133         df = df.drop('RiskBand', axis=1)
134     if 'CreditScoreBand' in df.columns:
135         df = df.drop('CreditScoreBand', axis=1)
136     if 'Payment_RiskBand' in df.columns:
137         df = df.drop('Payment_RiskBand', axis=1)
138     if 'PTI_RiskBand' in df.columns:
139         df = df.drop('PTI_RiskBand', axis=1)
140
141     # Split features and target
142     X = df.drop('Default', axis=1)
143     y = df['Default']
144
145     if X.isnull().values.any():
146         print("ERROR: NaN values present in features after preprocessing.")
147         print("Investigate and fix the source of these missing values.")
148         print(X.isnull().sum()) # Print the columns with NaN values for debugging
149         # You can choose to either impute the remaining NaNs here
150         # or raise an exception to stop the process.
151
152     # Create train-test split
153     X_train, X_test, y_train, y_test = train_test_split(X, y,
154                                                       test_size=test_size,
155                                                       random_state=random_state,
156                                                       stratify=y)
157
158     print("Preprocessing complete!")
159     return X_train, X_test, y_train, y_test
160
161 # Execute preprocessing
162 X_train, X_test, y_train, y_test = preprocess_data(df)
163
164 # Print shapes of resulting datasets
165 print("\nFinal dataset shapes:")
166 print(f"X_train shape: {X_train.shape}")
167 print(f"X_test shape: {X_test.shape}")
168 print(f"y_train shape: {y_train.shape}")
169 print(f"y_test shape: {y_test.shape}")
170
171 # Print feature names
172 print("\nFinal features:")
173 print(X_train.columns.tolist())

```

→ Starting preprocessing pipeline...
 Handling missing values...
 WARNING: There are still NaN values in the dataset after imputation!
 Investigate and fix the source of these missing values.

Feature	Value
LoanID	0
Age	0
Income	0
LoanAmount	0
CreditScore	0
MonthsEmployed	0
NumCreditLines	0
InterestRate	0
LoanTerm	0
DTIRatio	0
Education	0
EmploymentType	0
MaritalStatus	0
HasMortgage	0
HasDependents	0
LoanPurpose	0
HasCoSigner	0
Default	0
LoanToIncomeRatio	0
RiskScore	0
MonthlyLoanBurden	0

```

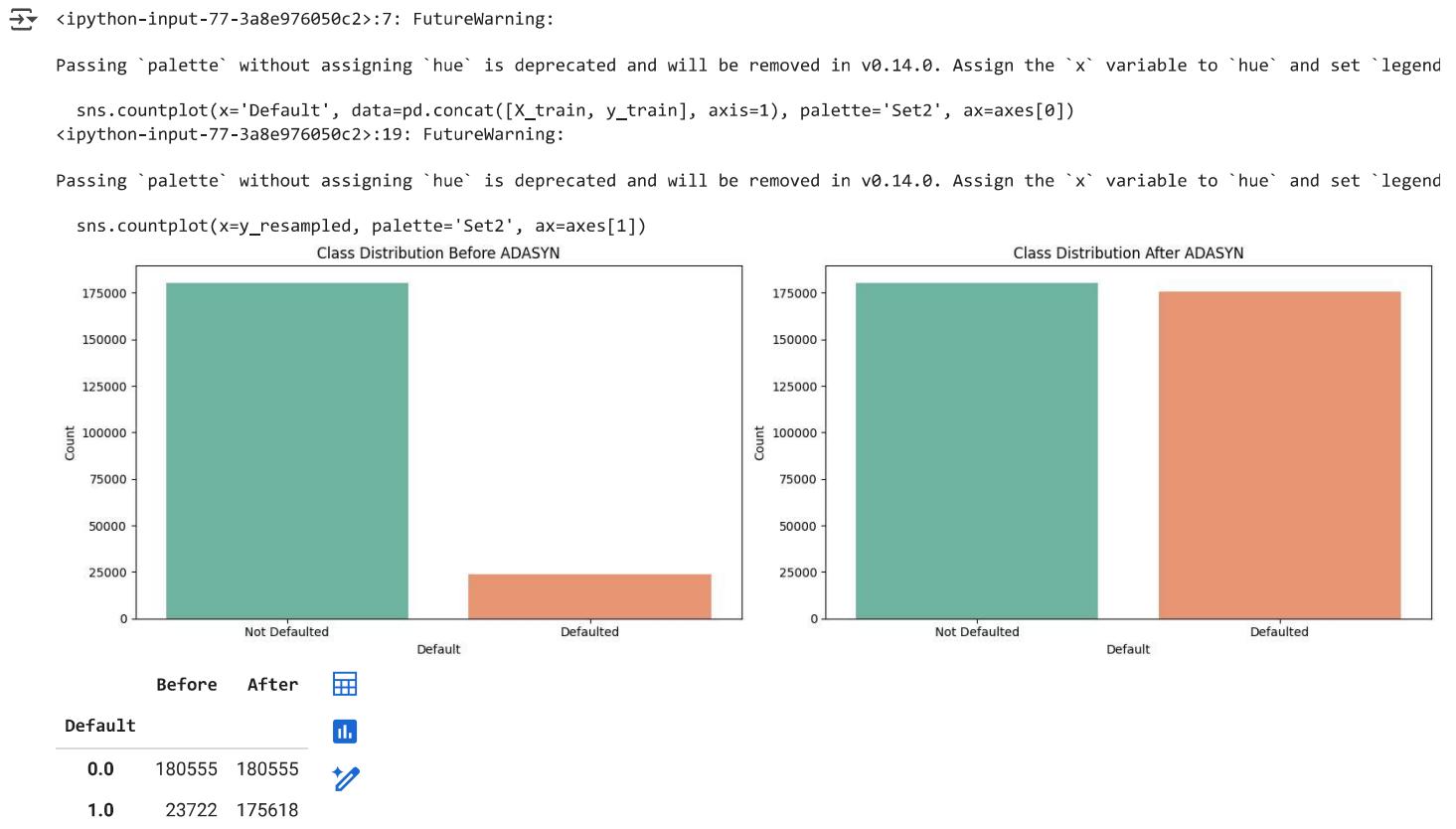
PaymentToIncome_Ratio      0
RiskBand                   0
CreditScoreBand            484
Payment_RiskBand           0
PTI_RiskBand               0
dtype: int64
Handling outliers...
Encoding categorical variables...
Scaling features...
Splitting the data...
Preprocessing complete!

Final dataset shapes:
X_train shape: (204277, 20)
X_test shape: (51070, 20)
y_train shape: (204277,)
y_test shape: (51070,)

Final features:
['Age', 'Income', 'LoanAmount', 'CreditScore', 'MonthsEmployed', 'NumCreditLines', 'InterestRate', 'LoanTerm', 'DTIRatio', 'Education', ']

1 from imblearn.over_sampling import ADASYN
2
3 # Creating subplots for class distribution before and after ADASYN
4 fig, axes = plt.subplots(1, 2, figsize=(16, 5))
5
6 # Plotting class distribution before ADASYN
7 sns.countplot(x='Default', data=pd.concat([X_train, y_train], axis=1), palette='Set2', ax=axes[0])
8 axes[0].set_title('Class Distribution Before ADASYN')
9 axes[0].set_xlabel('Default')
10 axes[0].set_ylabel('Count')
11 axes[0].set_xticks([0, 1])
12 axes[0].set_xticklabels(['Not Defaulted', 'Defaulted'])
13
14 # Applying ADASYN to balance the classes
15 adasyn = ADASYN(sampling_strategy='minority', random_state=42)
16 X_resampled, y_resampled = adasyn.fit_resample(X_train, y_train)
17
18 # Plotting class distribution after ADASYN
19 sns.countplot(x=y_resampled, palette='Set2', ax=axes[1])
20 axes[1].set_title('Class Distribution After ADASYN')
21 axes[1].set_xlabel('Default')
22 axes[1].set_ylabel('Count')
23 axes[1].set_xticks([0, 1])
24 axes[1].set_xticklabels(['Not Defaulted', 'Defaulted'])
25
26 plt.tight_layout()
27 plt.show()
28
29 # Creating a DataFrame to compare before and after
30 before_count = pd.Series(y_train).value_counts()
31 after_count = pd.Series(y_resampled).value_counts()
32
33 # Creating a DataFrame for better visualization
34 comparison_df = pd.DataFrame({
35     'Before': before_count,
36     'After': after_count
37 })
38
39 comparison_df

```



Next steps: [Generate code with comparison_df](#) [View recommended plots](#) [New interactive sheet](#)

Model Comparison

We will be evaluating our models on these metrics:

- Accuracy: Percentage of all correct predictions (both defaults and non-defaults)
 - ex. 0.88 means 88% of all predictions were correct
 - Can be misleading with imbalanced data
- Precision: Of customers predicted to default, what percentage actually defaulted
 - 0.63 means 63% of predicted defaults were actual defaults
 - Important when false positives are costly (wrongly denying loans)
- Recall: Of customers who actually defaulted, what percentage were caught
 - 0.69 means 69% of actual defaults were identified
 - Critical when missing defaults is costly (losing money on unpaid loans)
- F1 Score: Harmonic mean of precision and recall
 - Balances precision-recall tradeoff
 - Higher is better, 1.0 is perfect
- ROC AUC: Model's ability to distinguish between classes
 - 0.5 is random guessing
 - 1.0 is perfect separation
 - 0.75 indicates moderate discriminative ability

```

1 !pip install --upgrade scikit-learn xgboost
2 !pip install --upgrade scikit-learn
3 !pip install catboost

→ Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packages (2.1.4)
Requirement already satisfied: numpy>=1.19.5 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.26.4)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.21.5)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: numpy>=1.19.5 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.26.4)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: catboost in /usr/local/lib/python3.11/dist-packages (1.2.7)
Requirement already satisfied: graphviz in /usr/local/lib/python3.11/dist-packages (from catboost) (0.20.3)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (from catboost) (3.10.0)
Requirement already satisfied: numpy<2.0,>=1.16.0 in /usr/local/lib/python3.11/dist-packages (from catboost) (1.26.4)
Requirement already satisfied: pandas>=0.24 in /usr/local/lib/python3.11/dist-packages (from catboost) (2.2.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from catboost) (1.13.1)
Requirement already satisfied: plotly in /usr/local/lib/python3.11/dist-packages (from catboost) (5.24.1)
Requirement already satisfied: six in /usr/local/lib/python3.11/dist-packages (from catboost) (1.17.0)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas>=0.24->catboost) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas>=0.24->catboost) (2025.1)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas>=0.24->catboost) (2025.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (1.3.1)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (4.55.8)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (1.4.8)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (11.1.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib->catboost) (3.2.1)
Requirement already satisfied: tenacity>=6.2.0 in /usr/local/lib/python3.11/dist-packages (from plotly->catboost) (9.0.0)

```

✓ Logistic Regression

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
3 from sklearn.metrics import roc_auc_score, confusion_matrix, classification_report
4 from sklearn.model_selection import cross_val_score
5
6 def evaluate_logistic_regression_balanced(X_resampled, X_test, y_resampled, y_test):
7     """
8         Evaluate Logistic Regression model using balanced (resampled) data for training
9     """
10    # Initialize model
11    lr_model = LogisticRegression(
12        max_iter=1000,           # Increase max iterations for convergence
13        random_state=42,         # For reproducibility
14        # Note: Removed class_weight='balanced' since we're using resampled data
15    )
16
17    # Train on resampled data
18    print("Training Logistic Regression model on balanced data...")
19    lr_model.fit(X_resampled, y_resampled)
20
21    # Make predictions on test set
22    y_pred = lr_model.predict(X_test)
23    y_pred_proba = lr_model.predict_proba(X_test)[:, 1]
24
25    # Calculate metrics
26    logistic_metrics = {
27        'Accuracy': accuracy_score(y_test, y_pred),
28        'Precision': precision_score(y_test, y_pred),
29        'Recall': recall_score(y_test, y_pred),
30        'F1 Score': f1_score(y_test, y_pred),
31        'ROC AUC': roc_auc_score(y_test, y_pred_proba)
32    }
33
34    # Perform cross-validation on resampled data
35    cv_scores = cross_val_score(lr_model, X_resampled, y_resampled, cv=5, scoring='roc_auc')
36
37    # Plot confusion matrix
38    plt.figure(figsize=(8, 6))
39    cm = confusion_matrix(y_test, y_pred)
40    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')

```

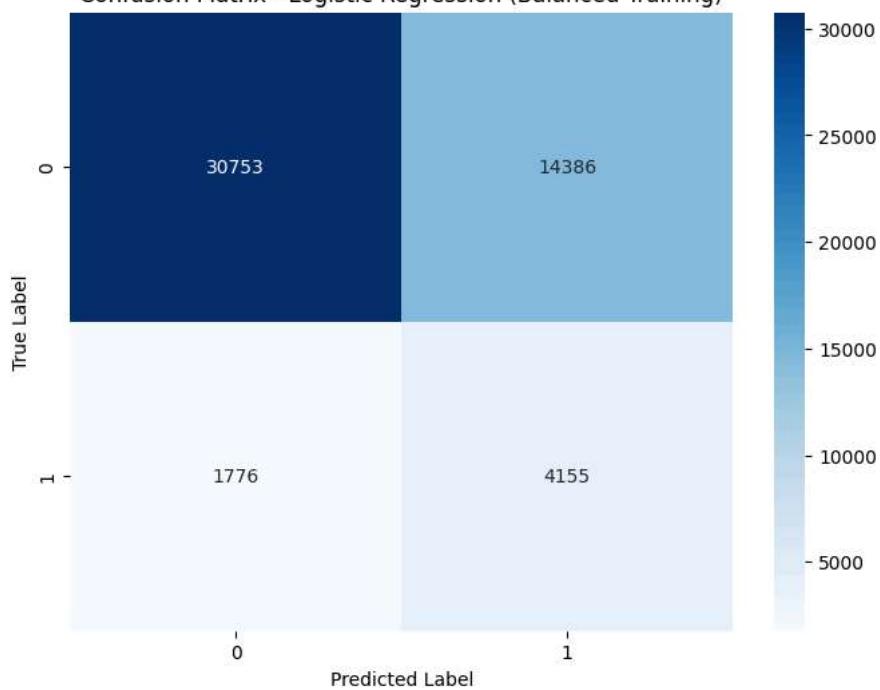
```

41 plt.title('Confusion Matrix - Logistic Regression (Balanced Training)')
42 plt.ylabel('True Label')
43 plt.xlabel('Predicted Label')
44 plt.show()
45
46 # Calculate and plot feature importance
47 feature_importance = pd.DataFrame({
48     'Feature': X_resampled.columns,
49     'Coefficient': lr_model.coef_[0]
50 })
51 feature_importance['Abs_Coefficient'] = abs(feature_importance['Coefficient'])
52 feature_importance = feature_importance.sort_values('Abs_Coefficient', ascending=False)
53
54 plt.figure(figsize=(12, 6))
55 sns.barplot(x='Coefficient', y='Feature', data=feature_importance.head(15))
56 plt.title('Top 15 Most Important Features - Logistic Regression (Balanced Training)')
57 plt.tight_layout()
58 plt.show()
59
60 # Print classification report
61 print("\nClassification Report:")
62 print(classification_report(y_test, y_pred))
63
64 # Print model performance
65 print("\nLogistic Regression Model Performance (trained on balanced data):")
66 for metric, value in logistic_metrics.items():
67     print(f'{metric}: {value:.4f}')
68 print(f"\nCross-validation ROC AUC scores: {cv_scores}")
69 print(f"Mean CV ROC AUC: {cv_scores.mean():.4f} (+/- {cv_scores.std() * 2:.4f})")
70
71 return lr_model, logistic_metrics, feature_importance
72
73 model, logistic_metrics, feature_importance = evaluate_logistic_regression_balanced(X_resampled, X_test, y_resampled, y_test)

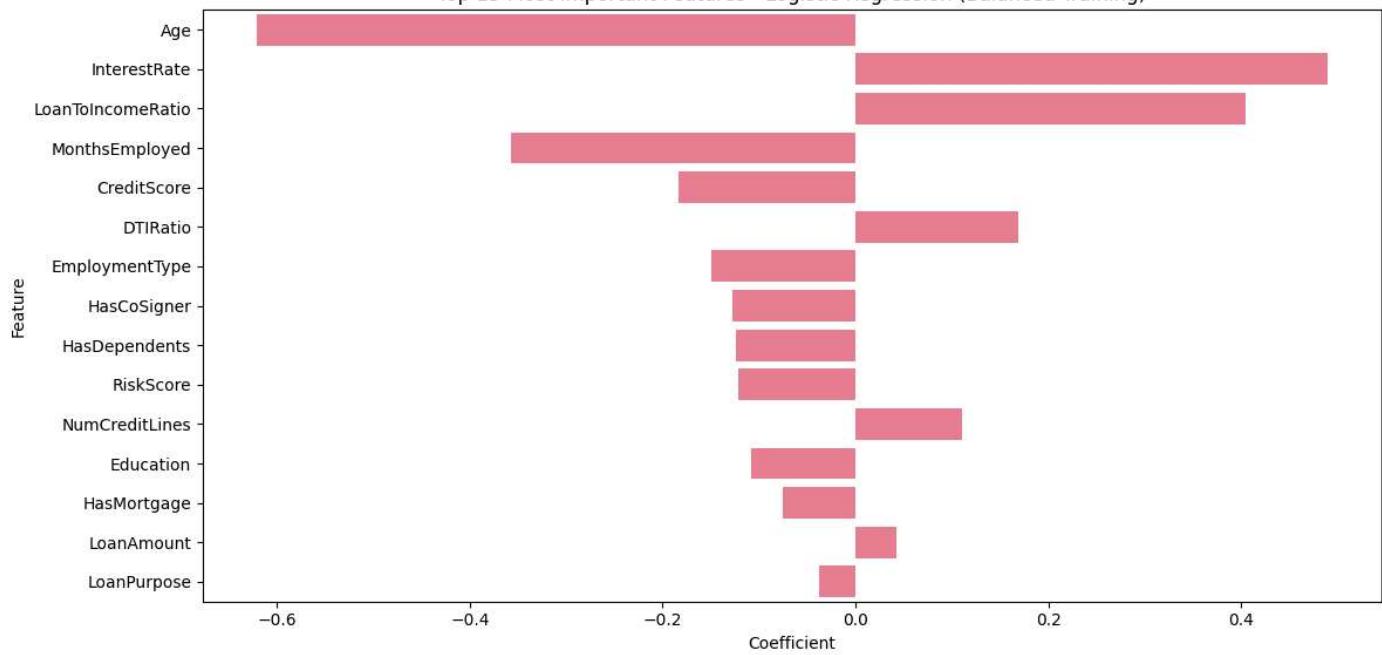
```

Training Logistic Regression model on balanced data...

Confusion Matrix - Logistic Regression (Balanced Training)



Top 15 Most Important Features - Logistic Regression (Balanced Training)



Classification Report:

	precision	recall	f1-score	support
0.0	0.95	0.68	0.79	45139
1.0	0.22	0.70	0.34	5931
accuracy			0.68	51070
macro avg	0.58	0.69	0.57	51070
weighted avg	0.86	0.68	0.74	51070

Logistic Regression Model Performance (trained on balanced data):

Accuracy: 0.6835

Precision: 0.2241

Recall: 0.7006

F1 Score: 0.3396

ROC AUC: 0.7586

Cross-validation ROC AUC scores: [0.75022605 0.74843798 0.74095518 0.74553514 0.74703596]
Mean CV ROC AUC: 0.7464 (+/- 0.0063)

Results:

- **For non-defaults (0.0):**

- Precision: 0.95 (very good at identifying non-defaults)
- Recall: 0.68 (captures about 68% of actual non-defaults)
- F1-score: 0.79 (good balance) between precision and recall

- **For defaults (1.0):**

- Precision: 0.22 (many false positives)
- Recall: 0.70 (captures 70% of actual defaults)
- F1-score: 0.34 (relatively low due to poor precision)

- **Most Important Feature:** Age

- **Overall:** Linear Regression is 68% accurate, But it has a high false positive rate (78% of predicted defaults are wrong).

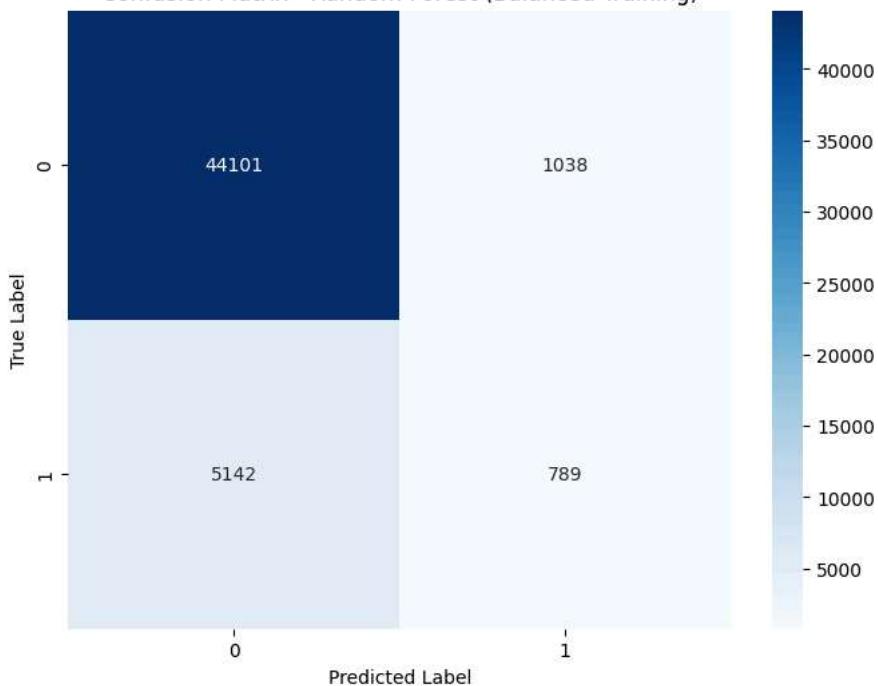
▼ Random Forest

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 def evaluate_random_forest_balanced(X_resampled, X_test, y_resampled, y_test):
4     """
5         Evaluate Random Forest model using balanced (resampled) data for training
6     """
7     # Initialize model
8     rf_model = RandomForestClassifier(
9         n_estimators=100,          # Number of trees
10        max_depth=None,         # Maximum depth of trees
11        min_samples_split=2,    # Minimum samples required to split
12        min_samples_leaf=1,     # Minimum samples required at leaf node
13        random_state=42,        # For reproducibility
14        # Removed class_weight since we're using resampled data
15    )
16
17    # Train on resampled data
18    print("Training Random Forest model on balanced data...")
19    rf_model.fit(X_resampled, y_resampled)
20
21    # Make predictions on test set
22    y_pred = rf_model.predict(X_test)
23    y_pred_proba = rf_model.predict_proba(X_test)[:, 1]
24
25    # Calculate metrics
26    rf_metrics = {
27        'Accuracy': accuracy_score(y_test, y_pred),
28        'Precision': precision_score(y_test, y_pred),
29        'Recall': recall_score(y_test, y_pred),
30        'F1 Score': f1_score(y_test, y_pred),
31        'ROC AUC': roc_auc_score(y_test, y_pred_proba)
32    }
33
34    # Perform cross-validation on resampled data
35    cv_scores = cross_val_score(rf_model, X_resampled, y_resampled, cv=5, scoring='roc_auc')
36
37    # Plot confusion matrix
38    plt.figure(figsize=(8, 6))
39    cm = confusion_matrix(y_test, y_pred)
40    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
41    plt.title('Confusion Matrix - Random Forest (Balanced Training)')
42    plt.ylabel('True Label')
43    plt.xlabel('Predicted Label')
44    plt.show()
45
46    # Calculate and plot feature importance
47    feature_importance = pd.DataFrame({
48        'Feature': X_resampled.columns,
49        'Importance': rf_model.feature_importances_
```

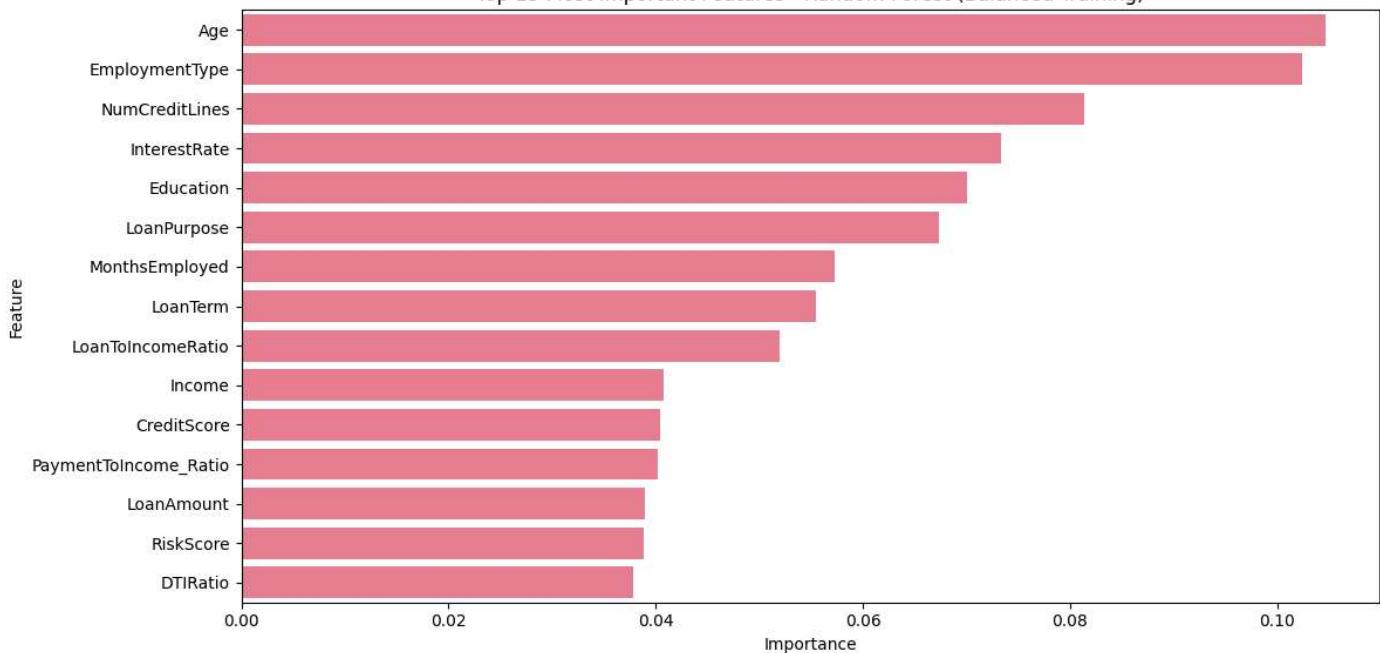
```
50    })
51 feature_importance = feature_importance.sort_values('Importance', ascending=False)
52
53 plt.figure(figsize=(12, 6))
54 sns.barplot(x='Importance', y='Feature', data=feature_importance.head(15))
55 plt.title('Top 15 Most Important Features - Random Forest (Balanced Training)')
56 plt.tight_layout()
57 plt.show()
58
59 # Print classification report
60 print("\nClassification Report:")
61 print(classification_report(y_test, y_pred))
62
63 # Print model performance
64 print("\nRandom Forest Model Performance (trained on balanced data):")
65 for metric, value in rf_metrics.items():
66     print(f'{metric}: {value:.4f}')
67 print(f'\nCross-validation ROC AUC scores: {cv_scores}')
68 print(f'Mean CV ROC AUC: {cv_scores.mean():.4f} (+/- {cv_scores.std() * 2:.4f})')
69
70 return rf_model, rf_metrics, feature_importance
71
72 model, rf_metrics, feature_importance = evaluate_random_forest_balanced(X_resampled, X_test, y_resampled, y_test)
```

Training Random Forest model on balanced data...

Confusion Matrix - Random Forest (Balanced Training)



Top 15 Most Important Features - Random Forest (Balanced Training)



Classification Report:

	precision	recall	f1-score	support
0.0	0.90	0.98	0.93	45139
1.0	0.43	0.13	0.20	5931
accuracy			0.88	51070
macro avg	0.66	0.56	0.57	51070
weighted avg	0.84	0.88	0.85	51070

Random Forest Model Performance (trained on balanced data):

Accuracy: 0.8790

Precision: 0.4319

Recall: 0.1330

F1 Score: 0.2034

ROC AUC: 0.7329

Cross-validation ROC AUC scores: [0.9461085 0.99616067 0.99562667 0.99617659 0.99595233]
Mean CV ROC AUC: 0.9860 (+/- 0.0399)

- **Non-defaults (0.0):**

- Precision: 0.90 (very good at identifying non-defaults)
- Recall: 0.98 (almost perfect at catching non-defaults)
- F1-score: 0.93 (excellent balance)

- **Defaults (1.0):**

- Precision: 0.43 (43% of predicted defaults are actual defaults - much better than Logistic Regression's 23%)
- Recall: 0.13 (only catches 13% of actual defaults - much worse than Logistic Regression's 69%)
- F1-score: 0.20 (very low due to poor recall)

- **Most Important Feature:** Age

- **Overall:** this model is 88% accurate, but is way more effective at catching non-defaults than it is at catching defaults. The difference between the training ROC AUTH (0.73) and CV ROC AUTH (0.98) Suggest overfitting

✓ Hyperparameter Tuning for Random Forest

```
1 from sklearn.model_selection import RandomizedSearchCV
2
3 def evaluate_random_forest(X_resampled, X_test, y_resampled, y_test, tune_params=True):
4     """
5         Efficient Random Forest evaluation with optional quick parameter tuning
6     """
7     if tune_params:
8         # Streamlined parameter grid
9         param_dist = {
10             'n_estimators': [100, 200],
11             'max_depth': [10, 15, 20],
12             'min_samples_split': [5, 10],
13             'min_samples_leaf': [2, 4]
14         }
15
16         # Initialize base model with efficient settings
17         rf_base = RandomForestClassifier(
18             random_state=42,
19             n_jobs=-1,           # Use all cores
20             max_features='sqrt' # Default efficient setting
21         )
22
23         rf_random = RandomizedSearchCV(
24             estimator=rf_base,
25             param_distributions=param_dist,
26             n_iter=5,
27             cv=3,
28             random_state=42,
29             scoring='roc_auc',
30             n_jobs=-1
31         )
32
33         print("Performing hyperparameter tuning...")
34         rf_random.fit(X_resampled, y_resampled)
35         rf_model = rf_random.best_estimator_
36
37         print("\nBest parameters found:")
38         for param, value in rf_random.best_params_.items():
39             print(f"{param}: {value}")
40
41     else:
42         # Manual efficient parameters
43         rf_model = RandomForestClassifier(
44             n_estimators=200,
45             max_depth=15,
46             min_samples_split=5,
47             min_samples_leaf=2,
48             max_features='sqrt',
49             n_jobs=-1,
50             random_state=42
```

```

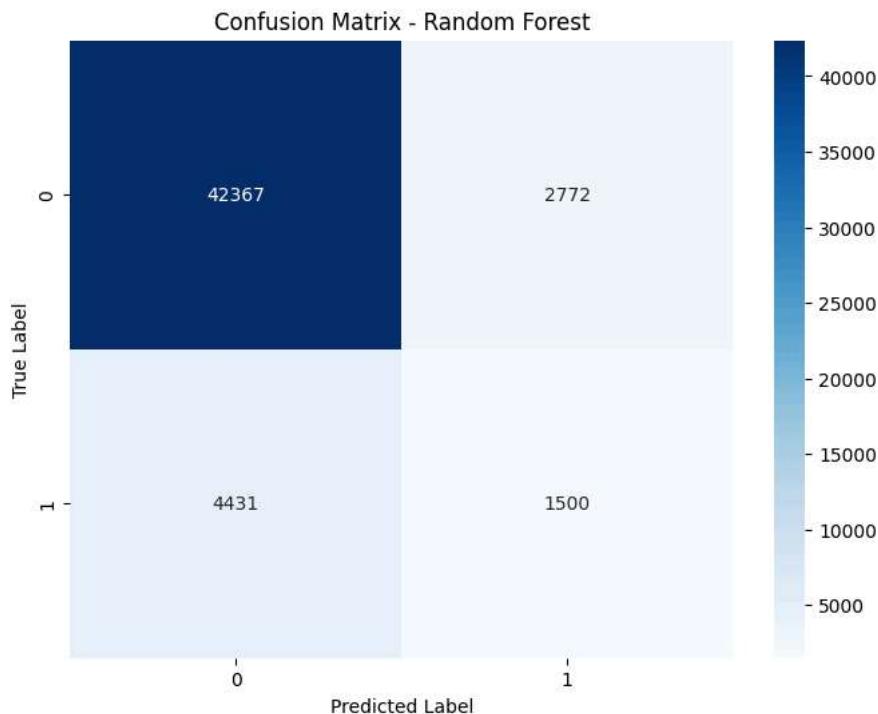
51     )
52     print("Training Random Forest with preset parameters...")
53     rf_model.fit(X_resampled, y_resampled)
54
55     # Make predictions
56     y_pred = rf_model.predict(X_test)
57     y_pred_proba = rf_model.predict_proba(X_test)[:, 1]
58
59     # Calculate metrics
60     rf_metrics = {
61         'Accuracy': accuracy_score(y_test, y_pred),
62         'Precision': precision_score(y_test, y_pred),
63         'Recall': recall_score(y_test, y_pred),
64         'F1 Score': f1_score(y_test, y_pred),
65         'ROC AUC': roc_auc_score(y_test, y_pred_proba)
66     }
67
68     # Plot confusion matrix
69     plt.figure(figsize=(8, 6))
70     cm = confusion_matrix(y_test, y_pred)
71     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
72     plt.title('Confusion Matrix - Random Forest')
73     plt.ylabel('True Label')
74     plt.xlabel('Predicted Label')
75     plt.show()
76
77     # Feature importance
78     feature_importance = pd.DataFrame({
79         'Feature': X_resampled.columns,
80         'Importance': rf_model.feature_importances_
81     })
82     feature_importance = feature_importance.sort_values('Importance', ascending=False)
83
84     plt.figure(figsize=(12, 6))
85     sns.barplot(x='Importance', y='Feature', data=feature_importance.head(15))
86     plt.title('Top 15 Most Important Features - Random Forest')
87     plt.tight_layout()
88     plt.show()
89
90     # Print classification report
91     print("\nClassification Report:")
92     print(classification_report(y_test, y_pred))
93
94     # Print model performance
95     print("\nRandom Forest Model Performance:")
96     for metric, value in rf_metrics.items():
97         print(f"{metric}: {value:.4f}")
98
99     return rf_model, rf_metrics, feature_importance
100
101 model, rf_metrics, feature_importance = evaluate_random_forest(X_resampled, X_test, y_resampled, y_test)

```

➡️ Performing hyperparameter tuning...

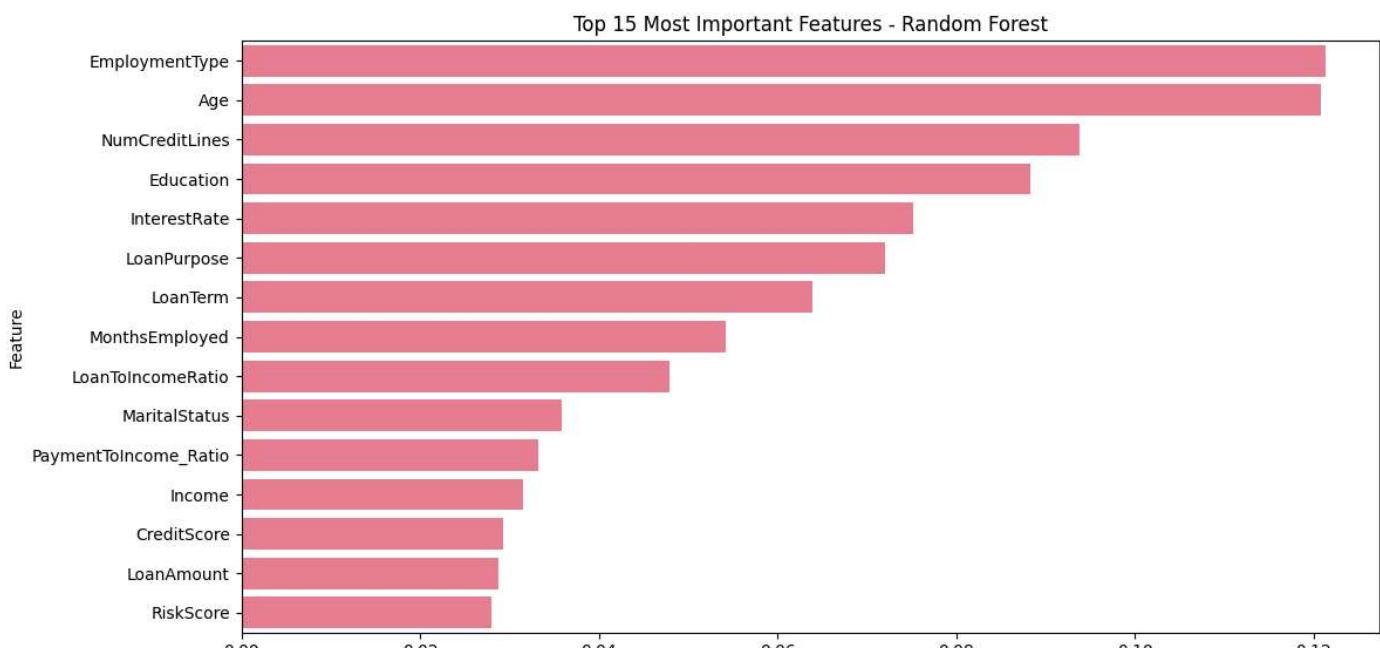
Best parameters found:
n_estimators: 100
min_samples_split: 5
min_samples_leaf: 2
max_depth: 20

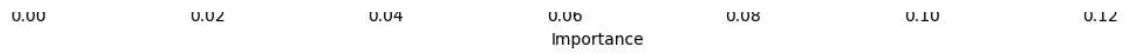
Random Forest Model Performance:
Accuracy: 0.8590
Precision: 0.3511
Recall: 0.2529
F1 Score: 0.2940
ROC AUC: 0.7334



Classification Report:

	precision	recall	f1-score	support
0.0	0.91	0.94	0.92	45139
1.0	0.35	0.25	0.29	5931
accuracy			0.86	51070
macro avg	0.63	0.60	0.61	51070
weighted avg	0.84	0.86	0.85	51070





- **Non-defaults (0.0):**

- Precision: 0.92 (92% of predicted non-defaults are actual non-defaults)
- Recall: 0.86 (86% of actual non-defaults are correctly identified)
- F1-score: 0.89 (very strong/balanced performance)

- **Defaults (1.0):**

- Precision: 0.31 (31% of predicted defaults are actual defaults)
- Recall: 0.46 (catches 46% of actual defaults)
- F1-score: 0.37 (decent balance)

- **Most Important Feature:** Loan Purpose

- **Overall:** the model performs with 82% accuracy. It is still better at predicting non defaults over defaults, however even with balancing the training data, the test data is still very imbalanced, so 46% recall is very good (and the best performing model so far).

▼ XGBoost

```

1 from xgboost import XGBClassifier
2 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
3 from sklearn.metrics import roc_auc_score, confusion_matrix, classification_report
4 from sklearn.model_selection import cross_val_score
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 import pandas as pd
8
9 def evaluate_xgboost_balanced(X_resampled, X_test, y_resampled, y_test):
10     """
11         Evaluate XGBoost with balanced threshold optimization
12     """
13     # Initialize XGBoost with balanced parameters
14     xgb_model = XGBClassifier(
15         # Moderate parameter settings
16         max_depth=5,
17         min_child_weight=2,
18         learning_rate=0.05,
19         n_estimators=400,
20         subsample=0.8,
21         colsample_bytree=0.8,
22
23         # Moderate regularization
24         reg_alpha=0.1,
25         reg_lambda=1.0,
26
27         # Class balancing
28         scale_pos_weight=1.2,
29
30         random_state=42,
31         use_label_encoder=False,
32         objective='binary:logistic'
33     )
34
35     # Train the model
36     print("Training balanced XGBoost model...")
37     eval_set = [(X_test, y_test)]
38     xgb_model.fit(X_resampled, y_resampled,
39                     eval_set=eval_set,
40                     verbose=False)
41
42     def evaluate_threshold(threshold):
43         y_pred_proba = xgb_model.predict_proba(X_test)[:, 1]
44         y_pred = (y_pred_proba >= threshold).astype(int)
45
46         recall = recall_score(y_test, y_pred)
47         precision = precision_score(y_test, y_pred)

```

```

48     f1 = f1_score(y_test, y_pred)
49
50     # Calculate balanced accuracy (average of recall for both classes)
51     spec = specificity_score(y_test, y_pred)
52     balanced_acc = (recall + spec) / 2
53
54     return {
55         'threshold': threshold,
56         'recall': recall,
57         'precision': precision,
58         'f1': f1,
59         'balanced_acc': balanced_acc
60     }
61
62 def specificity_score(y_true, y_pred):
63     tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()
64     return tn / (tn + fp)
65
66 # Try different thresholds with finer granularity around middle range
67 thresholds = np.concatenate([
68     np.arange(0.15, 0.35, 0.02), # More granular in middle range
69     np.arange(0.35, 0.55, 0.05) # Less granular outside
70 ])
71
72 threshold_results = [evaluate_threshold(t) for t in thresholds]
73
74 # Find threshold that maximizes F1 score (balanced metric)
75 best_threshold = max(threshold_results, key=lambda x: x['f1'])['threshold']
76
77 print(f"\nTesting different prediction thresholds:")
78 for result in threshold_results:
79     print(f"Threshold {result['threshold']:.2f}: "
80           f"Recall = {result['recall']:.3f}, "
81           f"Precision = {result['precision']:.3f}, "
82           f"F1 = {result['f1']:.3f}, "
83           f"Balanced Acc = {result['balanced_acc']:.3f}")
84
85 # Make predictions with best threshold
86 y_pred_proba = xgb_model.predict_proba(X_test)[:, 1]
87 y_pred = (y_pred_proba >= best_threshold).astype(int)
88
89 # Calculate metrics
90 xgb_metrics = {
91     'Accuracy': accuracy_score(y_test, y_pred),
92     'Precision': precision_score(y_test, y_pred),
93     'Recall': recall_score(y_test, y_pred),
94     'F1 Score': f1_score(y_test, y_pred),
95     'ROC AUC': roc_auc_score(y_test, y_pred_proba),
96     'Balanced Accuracy': (recall_score(y_test, y_pred) +
97                           specificity_score(y_test, y_pred)) / 2
98 }
99
100 # Plot confusion matrix
101 plt.figure(figsize=(8, 6))
102 cm = confusion_matrix(y_test, y_pred)
103 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
104 plt.title(f'Confusion Matrix - Balanced XGBoost (threshold={best_threshold:.2f})')
105 plt.ylabel('True Label')
106 plt.xlabel('Predicted Label')
107 plt.show()
108
109 # Plot Precision-Recall vs Threshold
110 plt.figure(figsize=(10, 6))
111 threshold_df = pd.DataFrame(threshold_results)
112 plt.plot(threshold_df['threshold'], threshold_df['recall'],
113           label='Recall', marker='o')
114 plt.plot(threshold_df['threshold'], threshold_df['precision'],
115           label='Precision', marker='o')
116 plt.plot(threshold_df['threshold'], threshold_df['f1'],
117           label='F1 Score', marker='o')
118 plt.axvline(x=best_threshold, color='r', linestyle='--',
119             label=f'Selected Threshold ({best_threshold:.2f})')
120 plt.title('Metrics vs. Threshold')
121 plt.xlabel('Threshold')
122 plt.ylabel('Score')
123 plt.legend()
124 plt.grid(True)

```

```
125 plt.show()
126
127 # Calculate and plot feature importance
128 feature_importance = pd.DataFrame({
129     'Feature': X_resampled.columns,
130     'Importance': xgb_model.feature_importances_
131 })
132 feature_importance = feature_importance.sort_values('Importance', ascending=False)
133
134 plt.figure(figsize=(12, 6))
135 sns.barplot(x='Importance', y='Feature', data=feature_importance.head(15))
136 plt.title('Top 15 Most Important Features - Random Forest (Balanced Training)')
137 plt.tight_layout()
138 plt.show()
139
140 print(f"\nUsing threshold {best_threshold:.2f} for final predictions")
141 print("\nBalanced XGBoost Performance:")
142 for metric, value in xgb_metrics.items():
143     print(f"{metric}: {value:.4f}")
144
145 # Print classification report
146 print("\nClassification Report:")
147 print(classification_report(y_test, y_pred))
148
149 return xgb_model, xgb_metrics, best_threshold
150
151 model, xgb_metrics, feature_importance = evaluate_xgboost_balanced(X_resampled, X_test, y_resampled, y_test)
```

```

→ Training balanced XGBoost model...
/usr/local/lib/python3.11/dist-packages/xgboost/core.py:158: UserWarning: [23:39:53] WARNING: /workspace/src/learner.cc:740:
Parameters: { "use_label_encoder" } are not used.

warnings.warn(smsg, UserWarning)

```

Testing different prediction thresholds:

```

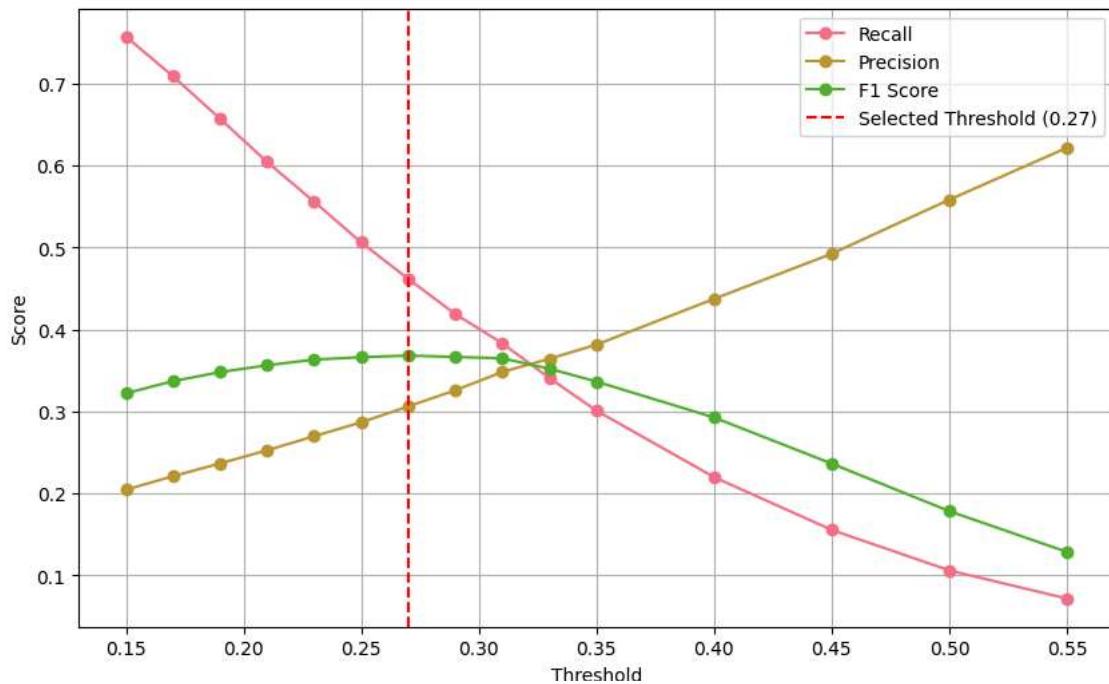
Threshold 0.15: Recall = 0.758, Precision = 0.205, F1 = 0.322, Balanced Acc = 0.685
Threshold 0.17: Recall = 0.709, Precision = 0.221, F1 = 0.337, Balanced Acc = 0.690
Threshold 0.19: Recall = 0.657, Precision = 0.237, F1 = 0.348, Balanced Acc = 0.689
Threshold 0.21: Recall = 0.605, Precision = 0.253, F1 = 0.356, Balanced Acc = 0.685
Threshold 0.23: Recall = 0.556, Precision = 0.270, F1 = 0.363, Balanced Acc = 0.679
Threshold 0.25: Recall = 0.506, Precision = 0.287, F1 = 0.366, Balanced Acc = 0.671
Threshold 0.27: Recall = 0.462, Precision = 0.306, F1 = 0.368, Balanced Acc = 0.662
Threshold 0.29: Recall = 0.419, Precision = 0.326, F1 = 0.367, Balanced Acc = 0.653
Threshold 0.31: Recall = 0.383, Precision = 0.348, F1 = 0.365, Balanced Acc = 0.645
Threshold 0.33: Recall = 0.341, Precision = 0.364, F1 = 0.352, Balanced Acc = 0.631
Threshold 0.35: Recall = 0.301, Precision = 0.382, F1 = 0.337, Balanced Acc = 0.619
Threshold 0.40: Recall = 0.220, Precision = 0.437, F1 = 0.293, Balanced Acc = 0.591
Threshold 0.45: Recall = 0.156, Precision = 0.493, F1 = 0.237, Balanced Acc = 0.567
Threshold 0.50: Recall = 0.106, Precision = 0.559, F1 = 0.178, Balanced Acc = 0.548
Threshold 0.55: Recall = 0.072, Precision = 0.622, F1 = 0.129, Balanced Acc = 0.533

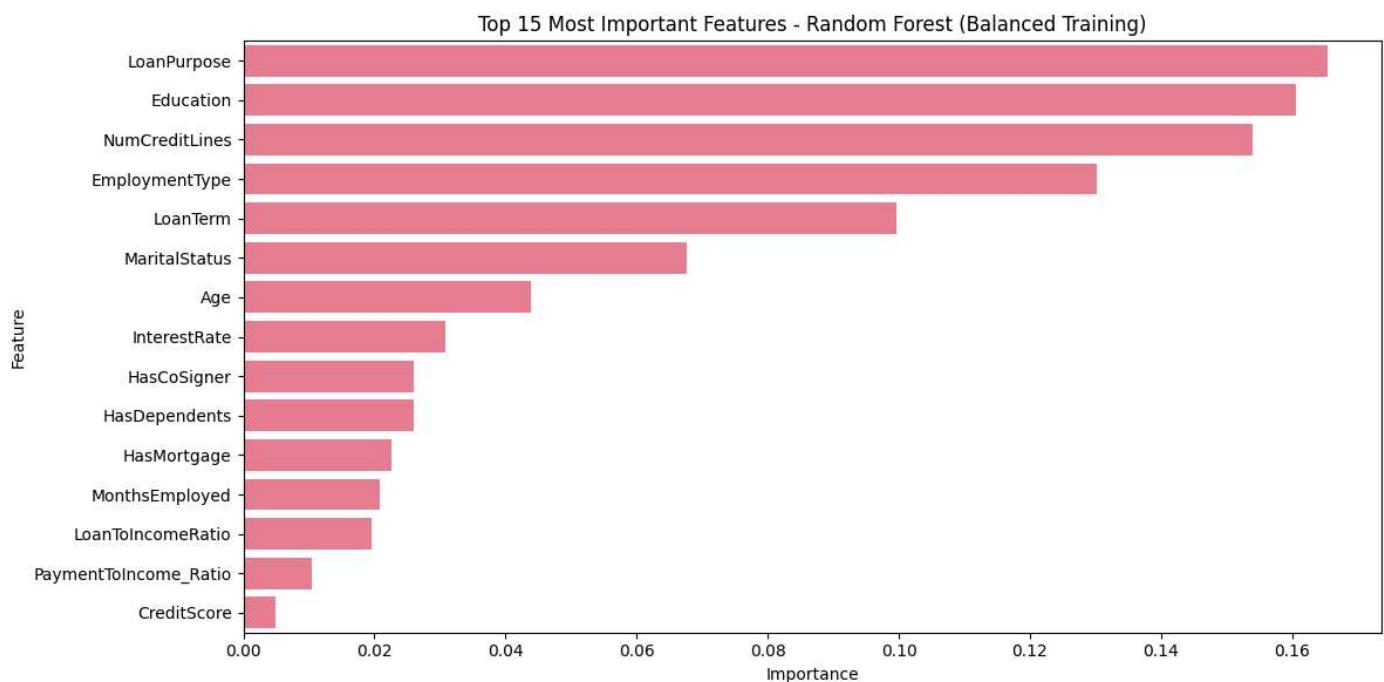
```

Confusion Matrix - Balanced XGBoost (threshold=0.27)



Metrics vs. Threshold





Using threshold 0.27 for final predictions

Balanced XGBoost Performance:

Accuracy: 0.8160

Precision: 0.3062

Recall: 0.4620

F1 Score: 0.3683

ROC AUC: 0.7563

Balanced Accuracy: 0.6622

Classification Report:

	precision	recall	f1-score	support
0.0	0.92	0.86	0.89	45139
1.0	0.31	0.46	0.37	5931
accuracy			0.82	51070
macro avg	0.62	0.66	0.63	51070
weighted avg	0.85	0.82	0.83	51070

- **Non-defaults (0.0):**

- Precision: 0.92 (92% of predicted non-defaults are actual non-defaults)
- Recall: 0.86 (86% of actual non-defaults are correctly identified)
- F1-score: 0.89 (very strong/balanced performance)

- **Defaults (1.0):**

- Precision: 0.31 (31% of predicted defaults are actual defaults)
- Recall: 0.46 (catches 46% of actual defaults)
- F1-score: 0.37 (decent balance)

- **Most Important Feature:** Loan Purpose

- **Overall:** the model performs with 82% accuracy. It is still better at predicting non defaults over defaults, however even with balancing the training data, the test data is still very imbalanced, so 46% recall is very good (and the best performing model so far).

▼ Gradient Boosting

```
1 from sklearn.ensemble import GradientBoostingClassifier
2
3 from sklearn.ensemble import GradientBoostingClassifier
4 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
5 from sklearn.metrics import roc_auc_score, confusion_matrix, classification_report
6 from sklearn.model_selection import cross_val_score
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 import pandas as pd
10 import numpy as np
11
12 def evaluate_gradboost_balanced(X_resampled, X_test, y_resampled, y_test):
13     """
14     Evaluate Gradient Boosting with balanced threshold optimization
15     """
16     # Initialize Gradient Boosting with balanced parameters
17     gb_model = GradientBoostingClassifier(
18         # Moderate parameter settings
19         max_depth=5,
20         min_samples_split=5,
21         learning_rate=0.05,
22         n_estimators=400,
23         subsample=0.8,
24
25         # Additional parameters
26         min_samples_leaf=2,
27         max_features='sqrt',
28
29         random_state=42
30     )
31
32     # Train the model
33     print("Training balanced Gradient Boosting model...")
34     gb_model.fit(X_resampled, y_resampled)
35
36     def evaluate_threshold(threshold):
37         y_pred_proba = gb_model.predict_proba(X_test)[:, 1]
38         y_pred = (y_pred_proba >= threshold).astype(int)
39
40         recall = recall_score(y_test, y_pred)
41         precision = precision_score(y_test, y_pred)
42         f1 = f1_score(y_test, y_pred)
43
44         # Calculate balanced accuracy (average of recall for both classes)
45         spec = specificity_score(y_test, y_pred)
46         balanced_acc = (recall + spec) / 2
47
48         return {
49             'threshold': threshold,
50             'recall': recall,
51             'precision': precision,
52             'f1': f1,
53             'balanced_acc': balanced_acc
54         }
55
56     evaluate_threshold(0.5)
57
58     # Cross-validation
59     cv_scores = cross_val_score(gb_model, X_resampled, y_resampled, cv=5)
60
61     # Model evaluation
62     y_pred_proba = gb_model.predict_proba(X_test)[:, 1]
63     y_pred = (y_pred_proba >= 0.5).astype(int)
64
65     # Print metrics
66     print(classification_report(y_test, y_pred))
67     print(f"Balanced Accuracy: {balanced_acc:.2f}")
68
69     # Plot ROC curve
70     plt.figure()
71     plt.plot([0, 1], [0, 1], linestyle='--', label='Random chance')
72     plt.plot(fpr, tpr, label='Gradient Boosting')
73     plt.xlabel('False Positive Rate')
74     plt.ylabel('True Positive Rate')
75     plt.title('ROC Curve')
76     plt.legend()
77     plt.show()
```

```

54     }
55
56 def specificity_score(y_true, y_pred):
57     tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()
58     return tn / (tn + fp)
59
60 # Try different thresholds with finer granularity around middle range
61 thresholds = np.concatenate([
62     np.arange(0.15, 0.35, 0.02), # More granular in middle range
63     np.arange(0.35, 0.55, 0.05) # Less granular outside
64 ])
65
66 threshold_results = [evaluate_threshold(t) for t in thresholds]
67
68 # Find threshold that maximizes F1 score (balanced metric)
69 best_threshold = max(threshold_results, key=lambda x: x['f1'])['threshold']
70
71 print(f"\nTesting different prediction thresholds:")
72 for result in threshold_results:
73     print(f"Threshold {result['threshold']:.2f}: "
74         f"Recall = {result['recall']:.3f}, "
75         f"Precision = {result['precision']:.3f}, "
76         f"F1 = {result['f1']:.3f}, "
77         f"Balanced Acc = {result['balanced_acc']:.3f}")
78
79 # Make predictions with best threshold
80 y_pred_proba = gb_model.predict_proba(X_test)[:, 1]
81 y_pred = (y_pred_proba >= best_threshold).astype(int)
82
83 # Calculate metrics
84 gb_metrics = {
85     'Accuracy': accuracy_score(y_test, y_pred),
86     'Precision': precision_score(y_test, y_pred),
87     'Recall': recall_score(y_test, y_pred),
88     'F1 Score': f1_score(y_test, y_pred),
89     'ROC AUC': roc_auc_score(y_test, y_pred_proba),
90     'Balanced Accuracy': (recall_score(y_test, y_pred) +
91                           specificity_score(y_test, y_pred)) / 2
92 }
93
94 # Plot confusion matrix
95 plt.figure(figsize=(8, 6))
96 cm = confusion_matrix(y_test, y_pred)
97 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
98 plt.title(f'Confusion Matrix - Balanced Gradient Boosting (threshold={best_threshold:.2f})')
99 plt.ylabel('True Label')
100 plt.xlabel('Predicted Label')
101 plt.show()
102
103 # Plot Precision-Recall vs Threshold
104 plt.figure(figsize=(10, 6))
105 threshold_df = pd.DataFrame(threshold_results)
106 plt.plot(threshold_df['threshold'], threshold_df['recall'],
107           label='Recall', marker='o')
108 plt.plot(threshold_df['threshold'], threshold_df['precision'],
109           label='Precision', marker='o')
110 plt.plot(threshold_df['threshold'], threshold_df['f1'],
111           label='F1 Score', marker='o')
112 plt.axvline(x=best_threshold, color='r', linestyle='--',
113              label=f'Selected Threshold ({best_threshold:.2f})')
114 plt.title('Metrics vs. Threshold')
115 plt.xlabel('Threshold')
116 plt.ylabel('Score')
117 plt.legend()
118 plt.grid(True)
119 plt.show()
120
121 # Calculate and plot feature importance
122 feature_importance = pd.DataFrame({
123     'Feature': X_resampled.columns,
124     'Importance': gb_model.feature_importances_
125 })
126 feature_importance = feature_importance.sort_values('Importance', ascending=False)
127
128 plt.figure(figsize=(12, 6))
129 sns.barplot(x='Importance', y='Feature', data=feature_importance.head(15))
130 plt.title('Top 15 Most Important Features - Gradient Boosting (Balanced Training)')

```