# BANKING –CAPSTONE PROJECT

Jupyter **Banking project** Last Checkpoint: 11 hours ago   (autosaved)

File    Edit    View    Insert    Cell    Kernel    Widgets    Help

▶ Run    ■    C    ▶▶    Code    ▾    ⌨

## Data Preparation and Preliminary Analysis

```
In [ ]:    1  import pandas as pd
           2
           3  # Load your dataset
           4  df = pd.read_csv('loan_CSV.csv')
           5
```

```
In [9]:    1  # Check the structure of the data
           2  print(df.shape)  # To see the dimensions (rows, columns)
           3  print(df.columns)  # To see all columns
           4  print(df.info())  # Summary of columns and their types
           5
```

```
(39717, 21)
Index(['id', 'member_id', 'loan_amnt', 'funded_amnt_inv', 'term', 'int_rate',
       'installment', 'grade', 'sub_grade', 'emp_title', 'emp_length',
       'home_ownership', 'annual_inc', 'verification_status', 'issue_d',
       'loan_status', 'purpose', 'title', 'revol_bal', 'revol_util',
       'application_type'],
      dtype='object')
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 39717 entries, 0 to 39716
Data columns (total 21 columns):
 #   Column               Non-Null Count   Dtype
---  ------               --------------   -----
 0   id                   39717 non-null   int64
 1   member_id            39717 non-null   int64
 2   loan_amnt            39717 non-null   int64
 3   funded_amnt_inv      39717 non-null   float64
 4   term                 39717 non-null   object
 5   int_rate             39717 non-null   object
 6   installment          39717 non-null   float64
 7   grade                39557 non-null   object
 8   sub_grade            39717 non-null   object
 9   emp_title            37258 non-null   object
 10  emp_length           38642 non-null   object
 11  home_ownership       39717 non-null   object
 12  annual_inc           39717 non-null   float64
```

```
17  title                39706 non-null  object
18  revol_bal            39717 non-null  int64
19  revol_util           39667 non-null  object
20  application_type     39717 non-null  object
dtypes: float64(3), int64(4), object(14)
memory usage: 6.4+ MB
None
```

In [29]:
```python
1  #Handling Missing Values and Duplicates
2
3  # Check for missing values
4  print(df.isnull().sum())  # Number of missing values in each column
5
6  # Remove duplicates
7  df.drop_duplicates(inplace=True)
8
```

```
numerical_column    0
dtype: int64
```

In [27]:
```python
1   import pandas as pd
2
3   # Example DataFrame
4   data = {'numerical_column': [1, 2, float('nan'), 4, 5, float('nan'), 7, 8]}
5   df = pd.DataFrame(data)
6
7   # Calculate mean and fill NaNs
8   mean_value = df['numerical_column'].mean()
9   df['numerical_column'].fillna(mean_value, inplace=True)
10
11  print(df)
12
```

```
   numerical_column
0               1.0
1               2.0
2               4.5
3               4.0
4               5.0
5               4.5
6               7.0
7               8.0
```

In [32]:
```python
1  #Variable Renaming:
2
3  # Rename columns if needed
4  df.rename(columns={'old_column_name': 'new_column_name'}, inplace=True)
5
```
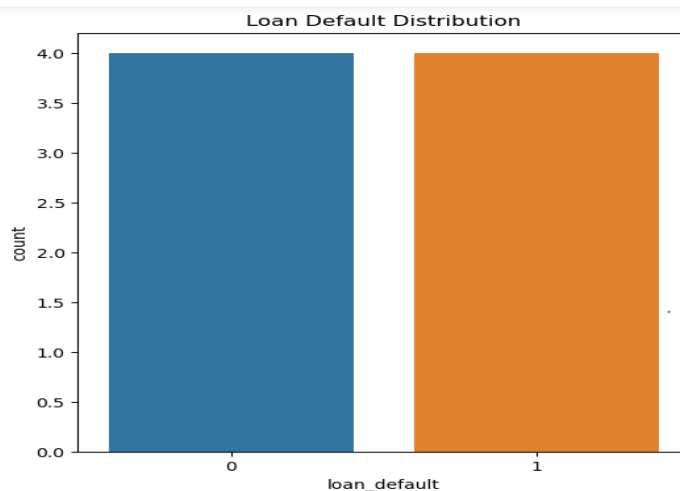
# Exploratory Data Analysis (EDA)

In [33]:
```
1  #Statistical Description:
2
3  # Summary statistics
4  print(df.describe())
5
```

```
       numerical_column
count                7.0
mean                 4.5
std                  2.5
min                  1.0
25%                  3.0
50%                  4.5
75%                  6.0
max                  8.0
```
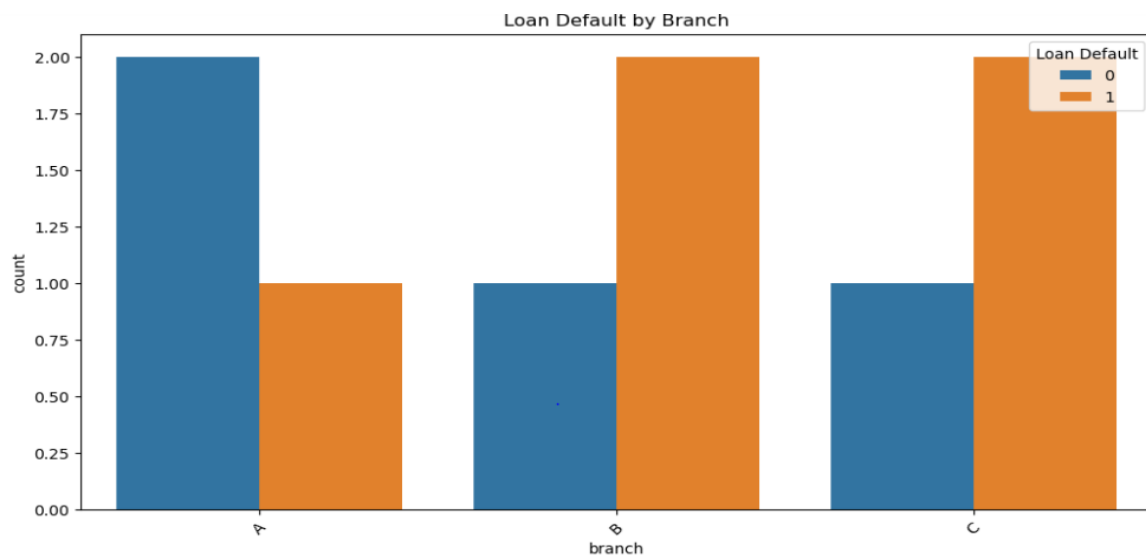
In [46]:
```
1  import matplotlib.pyplot as plt
2  import seaborn as sns
3  import pandas as pd
4
5  # Example DataFrame
6  data = {'loan_default': [0, 1, 1, 0, 0, 1, 0, 1]}
7  df = pd.DataFrame(data)
8
9  # Convert loan_default to categorical if it's numeric
10 df['loan_default'] = df['loan_default'].astype('category')
11
12 # Plotting with Seaborn countplot
13 plt.figure(figsize=(6, 6))
14 sns.countplot(x='loan_default', data=df)
15 plt.title('Loan Default Distribution')
16 plt.show()
17
```

```python
1  #Relationships Across Categories:
2
3  # Example: Loan default across different branches
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  import pandas as pd
7
8
9  data = {
10     'branch': ['A', 'B', 'A', 'C', 'B', 'C', 'A', 'B', 'C'],
11     'loan_default': [0, 1, 1, 0, 0, 1, 0, 1, 1]
12 }
13 df = pd.DataFrame(data)
14
15 # Convert 'branch' to categorical if needed
16 df['branch'] = df['branch'].astype('category')
17
18 # Plotting with Seaborn countplot
19 plt.figure(figsize=(12, 6))
20 sns.countplot(x='branch', hue='loan_default', data=df)
21 plt.title('Loan Default by Branch')
22 plt.xticks(rotation=45)
23 plt.legend(title='Loan Default', loc='upper right')
24 plt.show()
```
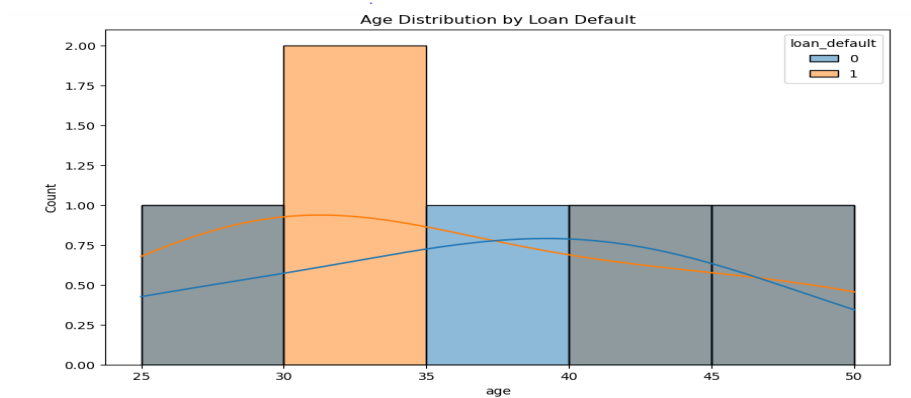
```python
1  #Age and Defaulting:
2
3  # Age distribution for defaulters vs. non-defaulters
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  import pandas as pd
7
8  # Example DataFrame
9  data = {
10     'age': [25, 30, 28, 35, 40, 32, 45, 50, 42],
11     'loan_default': [0, 1, 1, 0, 0, 1, 0, 1, 1]
12 }
13 df = pd.DataFrame(data)
14
15 # Plotting with Seaborn histplot
16 plt.figure(figsize=(10, 6))
17 sns.histplot(x='age', hue='loan_default', data=df, kde=True)
18 plt.title('Age Distribution by Loan Default')
19 plt.show()
20
```
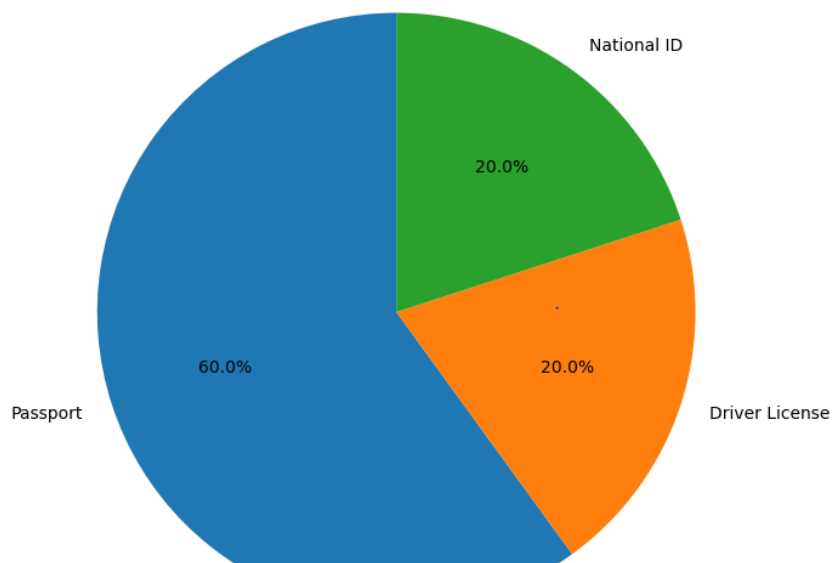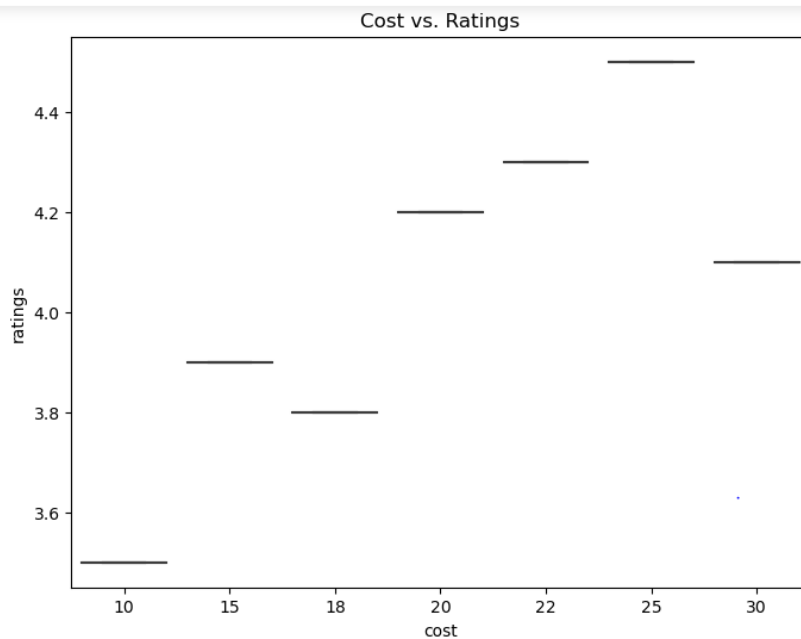
Age Distribution by Loan Default

```python
#Type of ID Presented:

# Count of different ID types
import matplotlib.pyplot as plt
import pandas as pd

# Example DataFrame (assuming 'id_type' is a valid column)
data = {
    'id_type': ['Passport', 'Driver License', 'Passport', 'Passport', 'National ID'],
}
df = pd.DataFrame(data)

# Print value counts to verify the column exists
print(df['id_type'].value_counts())

# Visualize with a pie chart
plt.figure(figsize=(8, 8))
df['id_type'].value_counts().plot(kind='pie', autopct='%1.1f%%', startangle=90)
plt.title('ID Types Presented by Customers')
plt.ylabel('')
plt.show()
```

```
id_type
Passport          3
Driver License    1
National ID       1
Name: count, dtype: int64
```



ID Types Presented by Customers

```
In [56]:   1  #Factors Affecting Ratings
           2
           3  # Example: Relationship between cost and ratings
           4  import matplotlib.pyplot as plt
           5  import seaborn as sns            .
           6  import pandas as pd
           7
           8  # Example DataFrame
           9  data = {
          10      'cost': [10, 20, 15, 25, 30, 18, 22],
          11      'ratings': [3.5, 4.2, 3.9, 4.5, 4.1, 3.8, 4.3]
          12  }
          13  df = pd.DataFrame(data)
          14
          15  # Plotting with Seaborn boxplot
          16  plt.figure(figsize=(8, 6))
          17  sns.boxplot(x='cost', y='ratings', data=df)
          18  plt.title('Cost vs. Ratings')
          19  plt.show()
          20
```



Cost vs. Ratings

# Logistic Regression Modeling

In [58]:
```python
#Data Preparation for Modeling:

import pandas as pd
from sklearn.model_selection import train_test_split

# Example DataFrame
data = {
    'age': [25, 30, 35, 40, 28],
    'employment_type': ['Full-time', 'Part-time', 'Full-time', 'Self-employed', 'Full-time'],
    'credit_score': [700, 650, 720, 690, 710],
    'loan_default': [0, 1, 0, 1, 0]
}
df = pd.DataFrame(data)

# Select relevant features and target variable
X = df[['age', 'employment_type', 'credit_score']]  # Select relevant features
y = df['loan_default']  # Target variable

# Encode categorical variables (if necessary)
X = pd.get_dummies(X, columns=['employment_type'], drop_first=True)

# Split data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Verify the selected features and target variable
print("Selected Features:")
print(X.columns)
print("\nTarget Variable:")
print(y.head())
```

```
Selected Features:
Index(['age', 'credit_score', 'employment_type_Part-time',
       'employment_type_Self-employed'],
      dtype='object')

Target Variable:
0    0
1    1
2    0
3    1
4    0
Name: loan_default, dtype: int64
```

In [59]:
```python
#Logistic Regression Model:


from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix, classification_report

# Initialize and fit logistic regression model
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

# Predictions
y_pred = model.predict(X_test)

# Evaluate model
print(confusion_matrix(y_test, y_pred))
print(classification_report(y_test, y_pred))
```

```
[[1]]
              precision    recall  f1-score   support

           1       1.00      1.00      1.00         1

    accuracy                           1.00         1
   macro avg       1.00      1.00      1.00         1
weighted avg       1.00      1.00      1.00         1
```

```python
In [62]:  1  import pandas as pd
          2  import numpy as np
          3  import matplotlib.pyplot as plt
          4  import seaborn as sns
          5  from sklearn.model_selection import train_test_split
          6  from sklearn.linear_model import LogisticRegression
          7  from sklearn.metrics import confusion_matrix, classification_report
          8
          9  # Set seed for reproducibility
         10  np.random.seed(42)
         11
         12  # Generate simulated data
         13  n_samples = 1000
         14
         15  # Credit bureau score (assuming normally distributed scores)
         16  credit_bureau_scores = np.random.normal(700, 50, n_samples)
         17
         18  # Default status (0 for non-default, 1 for default)
         19  default_status = np.random.choice([0, 1], size=n_samples, p=[0.8, 0.2])
         20
         21  # Primary and secondary account details (random amounts)
         22  primary_loan_amounts = np.random.uniform(1000, 5000, n_samples)
         23  secondary_loan_amounts = np.random.uniform(500, 3000, n_samples)
         24
         25  # Sanctioned and disbursed amounts (assuming some variability)
         26  sanctioned_amounts = primary_loan_amounts + np.random.uniform(-500, 500, n_samples)
         27  disbursed_amounts = primary_loan_amounts - np.random.uniform(0, 500, n_samples)
         28
         29  # Inquiry count (number of inquiries)
         30  inquiry_count = np.random.poisson(3, n_samples)
         31
```

```python
         32  # Credit history details
         33  new_loans_last_six_months = np.random.randint(0, 4, n_samples)
         34  loans_defaulted_last_six_months = np.random.randint(0, 2, n_samples)
         35  time_since_first_loan = np.random.randint(1, 10, n_samples)
         36
         37  # Create DataFrame
         38  data = pd.DataFrame({
         39      'credit_bureau_score': credit_bureau_scores,
         40      'default_status': default_status,
         41      'primary_loan_amount': primary_loan_amounts,
         42      'secondary_loan_amount': secondary_loan_amounts,
         43      'sanctioned_amount': sanctioned_amounts,
         44      'disbursed_amount': disbursed_amounts,
         45      'inquiry_count': inquiry_count,
         46      'new_loans_last_six_months': new_loans_last_six_months,
         47      'loans_defaulted_last_six_months': loans_defaulted_last_six_months,
         48      'time_since_first_loan': time_since_first_loan
         49  })
         50
```
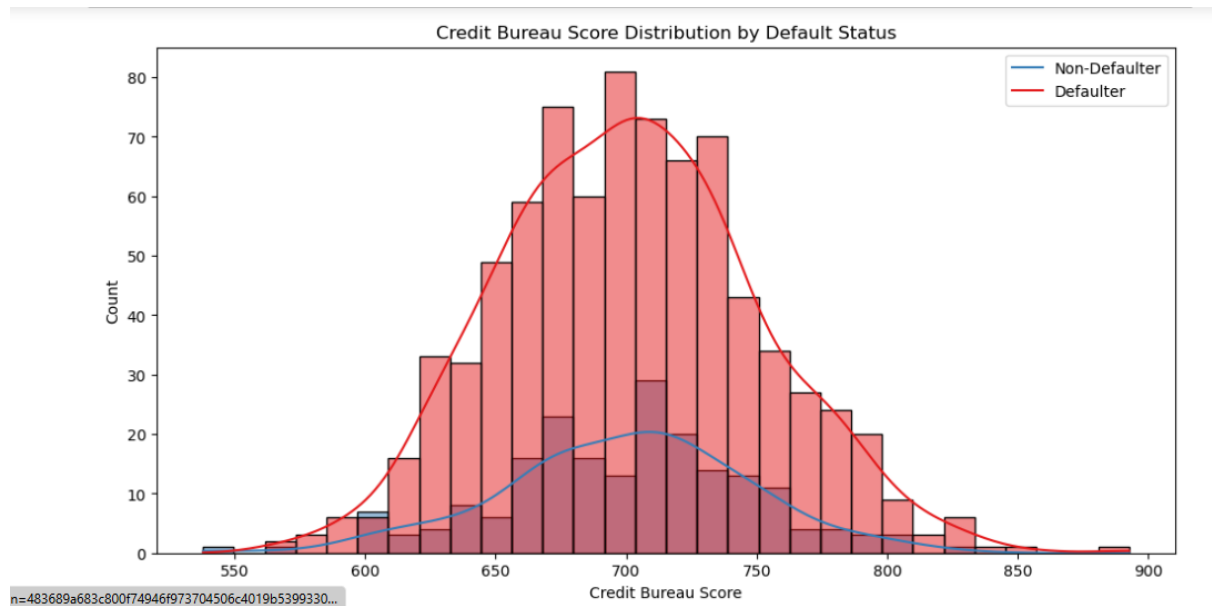
```python
          1  # Plotting credit bureau score distribution for defaulters vs. non-defaulters
          2  plt.figure(figsize=(12, 6))
          3  sns.histplot(data=data, x='credit_bureau_score', hue='default_status', kde=True, bins=30, palette='Set1')
          4  plt.title('Credit Bureau Score Distribution by Default Status')
          5  plt.xlabel('Credit Bureau Score')
          6  plt.ylabel('Count')
          7  plt.legend(['Non-Defaulter', 'Defaulter'])
          8  plt.show()
```
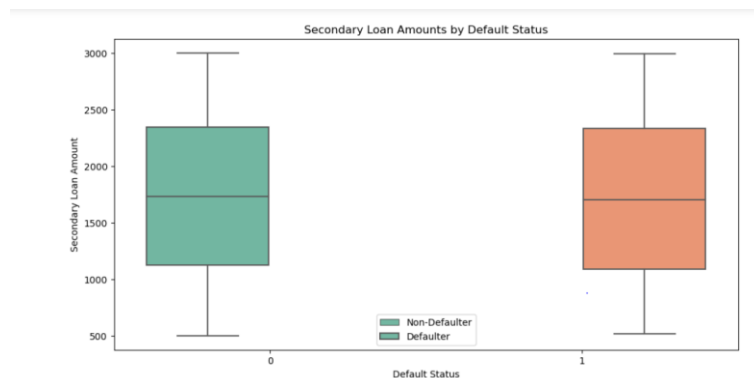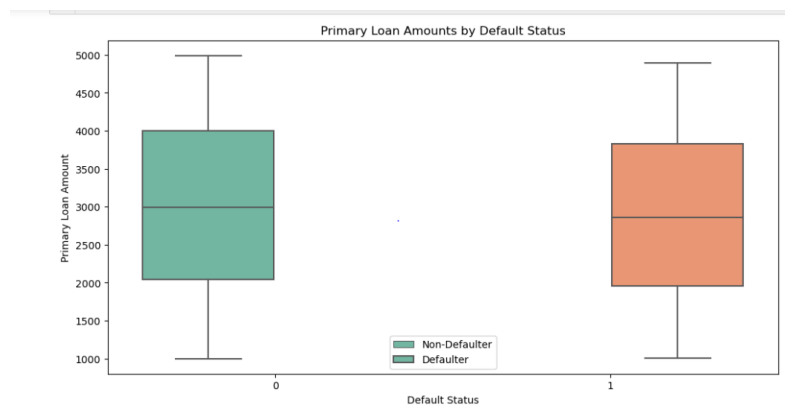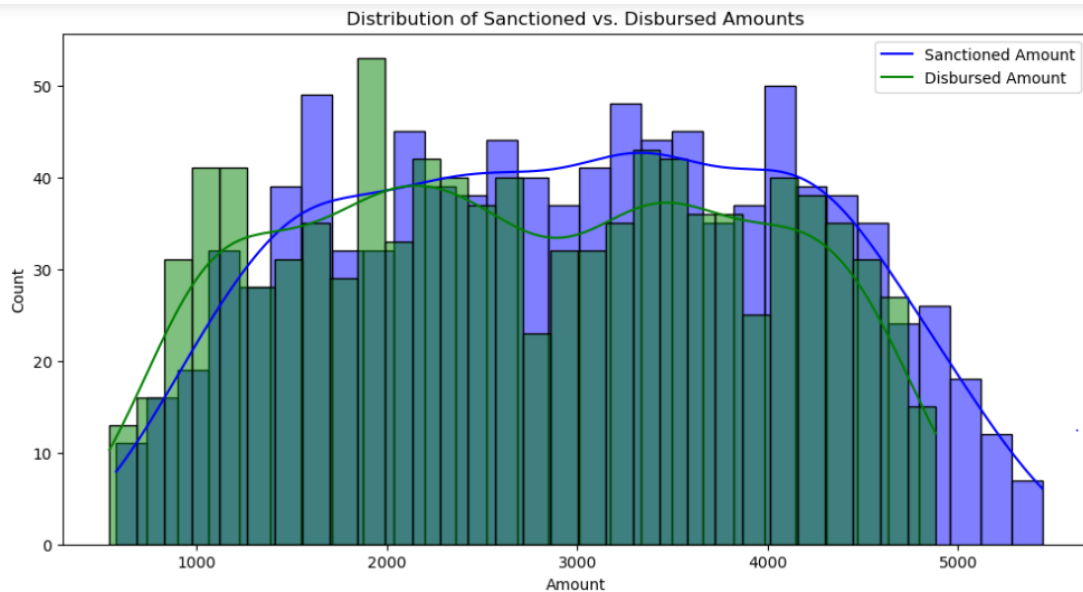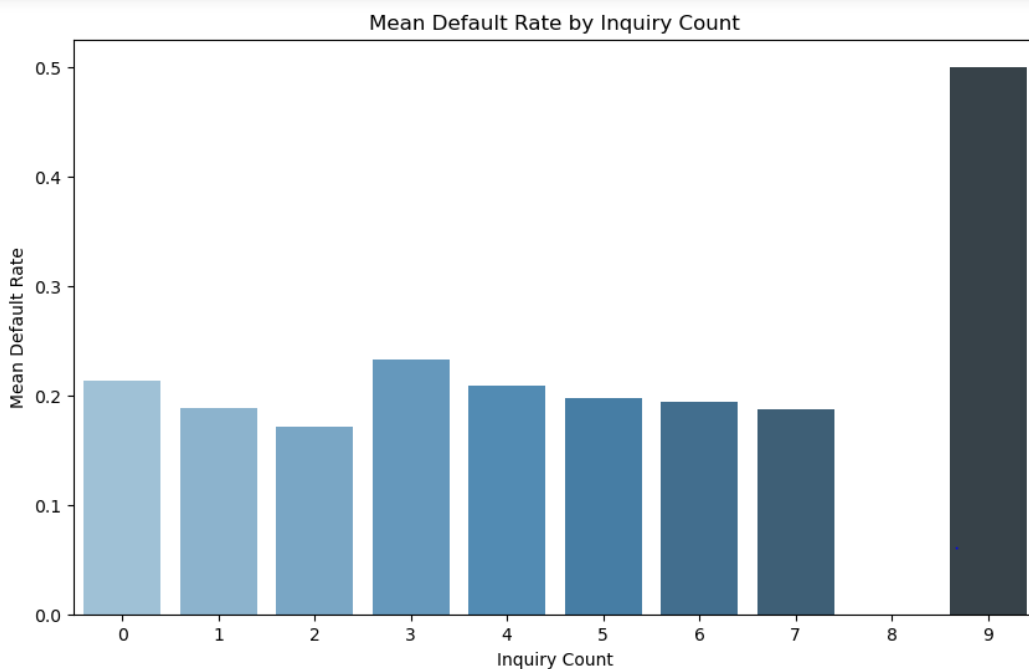
## Credit Bureau Score Distribution by Default Status



n=483689a683c800f74946f973704506c4019b5399330...

```python
In [64]:  1  # Boxplot to compare primary and secondary loan amounts by default status
          2  plt.figure(figsize=(12, 6))
          3  sns.boxplot(data=data, x='default_status', y='primary_loan_amount', hue='default_status', palette='Set2')
          4  plt.title('Primary Loan Amounts by Default Status')
          5  plt.xlabel('Default Status')
          6  plt.ylabel('Primary Loan Amount')
          7  plt.legend(['Non-Defaulter', 'Defaulter'])
          8  plt.show()
          9
         10  plt.figure(figsize=(12, 6))
         11  sns.boxplot(data=data, x='default_status', y='secondary_loan_amount', hue='default_status', palette='Set2')
         12  plt.title('Secondary Loan Amounts by Default Status')
         13  plt.xlabel('Default Status')
         14  plt.ylabel('Secondary Loan Amount')
         15  plt.legend(['Non-Defaulter', 'Defaulter'])
         16  plt.show()
         17
```

```python
1  # Comparing sanctioned and disbursed amounts
2  plt.figure(figsize=(12, 6))
3  sns.histplot(data=data, x='sanctioned_amount', bins=30, kde=True, color='blue', alpha=0.5)
4  sns.histplot(data=data, x='disbursed_amount', bins=30, kde=True, color='green', alpha=0.5)
5  plt.title('Distribution of Sanctioned vs. Disbursed Amounts')
6  plt.xlabel('Amount')
7  plt.ylabel('Count')
8  plt.legend(['Sanctioned Amount', 'Disbursed Amount'])
9  plt.show()
10
```



Distribution of Sanctioned vs. Disbursed Amounts

```python
1  # Relationship between inquiry count and default status
2  plt.figure(figsize=(10, 6))
3  sns.barplot(data=data, x='inquiry_count', y='default_status', estimator=np.mean, errorbar=None, palette='Blues_d')
4  plt.title('Mean Default Rate by Inquiry Count')
5  plt.xlabel('Inquiry Count')
6  plt.ylabel('Mean Default Rate')
7  plt.show()
8
```



Mean Default Rate by Inquiry Count

```
In [69]:    1  # Plotting credit history features by default status
            2  fig, axes = plt.subplots(1, 3, figsize=(18, 6), sharey=True)
            3
            4  sns.barplot(data=data, x='new_loans_last_six_months', y='default_status', ax=axes[0], estimator=np.mean,  errorbar=None, pal
            5  axes[0].set_title('Mean Default Rate by New Loans in Last 6 Months')
            6  axes[0].set_xlabel('New Loans in Last 6 Months')
            7  axes[0].set_ylabel('Mean Default Rate')
            8
            9  sns.barplot(data=data, x='loans_defaulted_last_six_months', y='default_status', ax=axes[1], estimator=np.mean,  errorbar=Non
           10  axes[1].set_title('Mean Default Rate by Loans Defaulted in Last 6 Months')
           11  axes[1].set_xlabel('Loans Defaulted in Last 6 Months')
           12  axes[1].set_ylabel('Mean Default Rate')
           13
           14  sns.barplot(data=data, x='time_since_first_loan', y='default_status', ax=axes[2], estimator=np.mean,  errorbar=None, palette
           15  axes[2].set_title('Mean Default Rate by Time Since First Loan')
           16  axes[2].set_xlabel('Time Since First Loan')
           17  axes[2].set_ylabel('Mean Default Rate')
           18
           19  plt.tight_layout()
           20  plt.show()
```