

Context

LoanTap is an online platform committed to delivering customized loan products to millennials. They innovate in an otherwise dull loan segment, to deliver instant, flexible loans on consumer friendly terms to salaried professionals and businessmen.

The data science team at LoanTap is building an underwriting layer to determine the creditworthiness of MSMEs as well as individuals.

LoanTap deploys formal credit to salaried individuals and businesses 4 main financial instruments:

→ Personal Loan → EMI Free Loan → Personal Overdraft → Advance Salary Loan

This project will focus on the underwriting process behind Personal Loan only.

→ 1. Define Problem Statement and perform Exploratory Data Analysis

Problem Statement :

Objective: Determine the creditworthiness of applicants to decide if a personal loan should be approved and define optimal repayment terms.

Extended Problem Definition :

LoanTap aims to minimize the risk of non-performing assets (NPAs) by accurately identifying potential loan defaulters while maximizing the loan issuance to creditworthy customers. This approach should balance maximizing profit (interest earned on loans) with minimizing losses (defaults) through an effective underwriting process.

b) Observations on shape of data, data types of all the attributes, conversion of categorical attributes to 'category' (If required), missing value detection, statistical summary.

```
In [1]: ┏━━━ 1 import numpy as np  
      2 import pandas as pd  
      3 import matplotlib.pyplot as plt  
      4 import seaborn as sns
```

```
In [2]: ┏━━━ 1 df = pd.read_csv('C:/SCALER\Case_Studies/7_LoanTrap/logistic_regressi
```

In [3]: 1 df.sample(5)

Out[3]:

	loan_amnt	term	int_rate	installment	grade	sub_grade	emp_title	emp_length
184742	4000.0	36 months	8.90	127.02	A	A5	Lead Pastor	7 year
180934	7200.0	36 months	7.62	224.37	A	A3	Mountain Park Health Center	4 year
20219	12175.0	36 months	16.29	429.79	C	C4	MPP Group Of Companies	9 year
1938	17500.0	60 months	8.18	356.35	B	B1	Operations Coordinator	9 year
8903	20000.0	60 months	11.99	444.79	C	C1	NaN	NaN

5 rows × 27 columns

In [4]: 1 df.shape

Out[4]: (396030, 27)

Feature description :-

loan_amnt : The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value.

term : The number of payments on the loan. Values are in months and can be either 36 or 60.

int_rate : Interest Rate on the loan

installment : The monthly payment owed by the borrower if the loan originates.

grade : LoanTap assigned loan grade

sub_grade : LoanTap assigned loan subgrade

emp_title :The job title supplied by the Borrower when applying for the loan.

emp_length : Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years.

home_ownership : The home ownership status provided by the borrower during registration or obtained from the credit report.

annual_inc : The self-reported annual income provided by the borrower during registration.

verification_status : Indicates if income was verified by LoanTap, not verified, or if the income source was verified

issue_d : The month which the loan was funded

loan_status : Current status of the loan - Target Variable

purpose : A category provided by the borrower for the loan request.

title : The loan title provided by the borrower

dti : A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LoanTap loan, divided by the borrower's self-reported monthly income.

earliest_cr_line : The month the borrower's earliest reported credit line was opened open_acc : The number of open credit lines in the borrower's credit file.

pub_rec : Number of derogatory public records

revol_bal : Total credit revolving balance

revol_util : Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.

total_acc : The total number of credit lines currently in the borrower's credit file

initial_list_status : The initial listing status of the loan. Possible values are – W, F

application_type : Indicates whether the loan is an individual application or a joint application with two co-borrowers mort_acc : Number of mortgage accounts.

pub_rec_bankruptcies : Number of public record bankruptcies Address: Address of the individual

In [5]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 396030 entries, 0 to 396029
Data columns (total 27 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   loan_amnt        396030 non-null  float64
 1   term              396030 non-null  object 
 2   int_rate          396030 non-null  float64
 3   installment       396030 non-null  float64
 4   grade             396030 non-null  object 
 5   sub_grade          396030 non-null  object 
 6   emp_title          373103 non-null  object 
 7   emp_length         377729 non-null  object 
 8   home_ownership     396030 non-null  object 
 9   annual_inc         396030 non-null  float64
 10  verification_status 396030 non-null  object 
 11  issue_d            396030 non-null  object 
 12  loan_status         396030 non-null  object 
 13  purpose             396030 non-null  object 
 14  title               394274 non-null  object 
 15  dti                 396030 non-null  float64
 16  earliest_cr_line    396030 non-null  object 
 17  open_acc            396030 non-null  float64
 18  pub_rec              396030 non-null  float64
 19  revol_bal            396030 non-null  float64
 20  revol_util           395754 non-null  float64
 21  total_acc            396030 non-null  float64
 22  initial_list_status 396030 non-null  object 
 23  application_type      396030 non-null  object 
 24  mort_acc              358235 non-null  float64
 25  pub_rec_bankruptcies 395495 non-null  float64
 26  address              396030 non-null  object 

dtypes: float64(12), object(15)
memory usage: 81.6+ MB
```

In [6]: 1 df.isnull().sum().sort_values(ascending=False).head(10)

```
Out[6]: mort_acc      37795
emp_title      22927
emp_length      18301
title          1756
pub_rec_bankruptcies  535
revol_util      276
loan_amnt          0
dti              0
application_type  0
initial_list_status 0
dtype: int64
```

In [7]: 1 df.describe()

Out[7]:

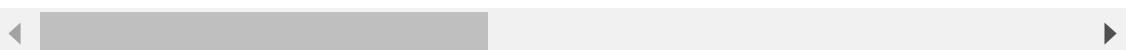
	loan_amnt	int_rate	installment	annual_inc	dti	op
count	396030.000000	396030.000000	396030.000000	3.960300e+05	396030.000000	396030.
mean	14113.888089	13.639400	431.849698	7.420318e+04	17.379514	11.
std	8357.441341	4.472157	250.727790	6.163762e+04	18.019092	5.
min	500.000000	5.320000	16.080000	0.000000e+00	0.000000	0.
25%	8000.000000	10.490000	250.330000	4.500000e+04	11.280000	8.
50%	12000.000000	13.330000	375.430000	6.400000e+04	16.910000	10.
75%	20000.000000	16.490000	567.300000	9.000000e+04	22.980000	14.
max	40000.000000	30.990000	1533.810000	8.706582e+06	9999.000000	90.



In [8]: 1 df.describe(include='object')

Out[8]:

	term	grade	sub_grade	emp_title	emp_length	home_ownership	verification_s
count	396030	396030	396030	373103	377729	396030	396030
unique	2	7	35	173105	11	6	396030
top	36 months	B	B3	Teacher	10+ years	MORTGAGE	VERIFIED
freq	302005	116018	26655	4389	126041	198348	130000



c) Univariate Analysis

In [9]: 1 object_features = []
2 numerical_features = []
3
4 for col in df.columns:
5 if df[col].dtype==object:
6 object_features.append(col)
7 else:
8 numerical_features.append(col)

In [10]: 1 print(object_features)

```
['term', 'grade', 'sub_grade', 'emp_title', 'emp_length', 'home_ownership', 'verification_status', 'issue_d', 'loan_status', 'purpose', 'title', 'earliest_cr_line', 'initial_list_status', 'application_type', 'addreses']
```

```
In [11]: 1 print(numerical_features)
```

```
['loan_amnt', 'int_rate', 'installment', 'annual_inc', 'dti', 'open_ac  
c', 'pub_rec', 'revol_bal', 'revol_util', 'total_acc', 'mort_acc', 'pub_  
rec_bankruptcies']
```

```
In [12]: 1 # List of numerical features (columns)  
2 numerical_features = ['loan_amnt', 'int_rate', 'installment', 'annual  
3           'pub_rec', 'revol_bal', 'revol_util', 'total_ac  
4  
5 # Check skewness for each numerical feature  
6 skewness_dict = {}  
7 for feature in numerical_features:  
8     skewness_dict[feature] = df[feature].skew()  
9  
10 # Display skewness for each numerical column  
11 for feature, skew in skewness_dict.items():  
12     print(f"Skewness of {feature}: {skew}")
```

```
Skewness of loan_amnt: 0.777285467097746  
Skewness of int_rate: 0.4206694719582165  
Skewness of installment: 0.983598160890438  
Skewness of annual_inc: 41.042724746560665  
Skewness of dti: 431.05122535490415  
Skewness of open_acc: 1.2130188444585455  
Skewness of pub_rec: 16.576564199464457  
Skewness of revol_bal: 11.727515124126919  
Skewness of revol_util: -0.07177802032853098  
Skewness of total_acc: 0.8643276369429  
Skewness of mort_acc: 1.6001324380874855  
Skewness of pub_rec_bankruptcies: 3.4234403681961583
```

In [13]:

```
1 from scipy import stats
2
3 def apply_log_transformation(df, features):
4     for feature in features:
5         df[feature] = np.log1p(df[feature]) # Log1p to handle zero values
6     return df
7
8 def apply_sqrt_transformation(df, features):
9     for feature in features:
10        df[feature] = np.sqrt(df[feature])
11    return df
12
13 log_transformed_features = ['loan_amnt', 'annual_inc', 'dti', 'pub_rec']
14 sqrt_transformed_features = ['installment', 'open_acc', 'total_acc']
15
16 df = apply_log_transformation(df, log_transformed_features)
17
18 df = apply_sqrt_transformation(df, sqrt_transformed_features)
19
20 print("\nSkewness after transformation:")
21 for feature in df.columns:
22     if feature in numerical_features:
23         print(f"Skewness of {feature}: {df[feature].skew()}")
```

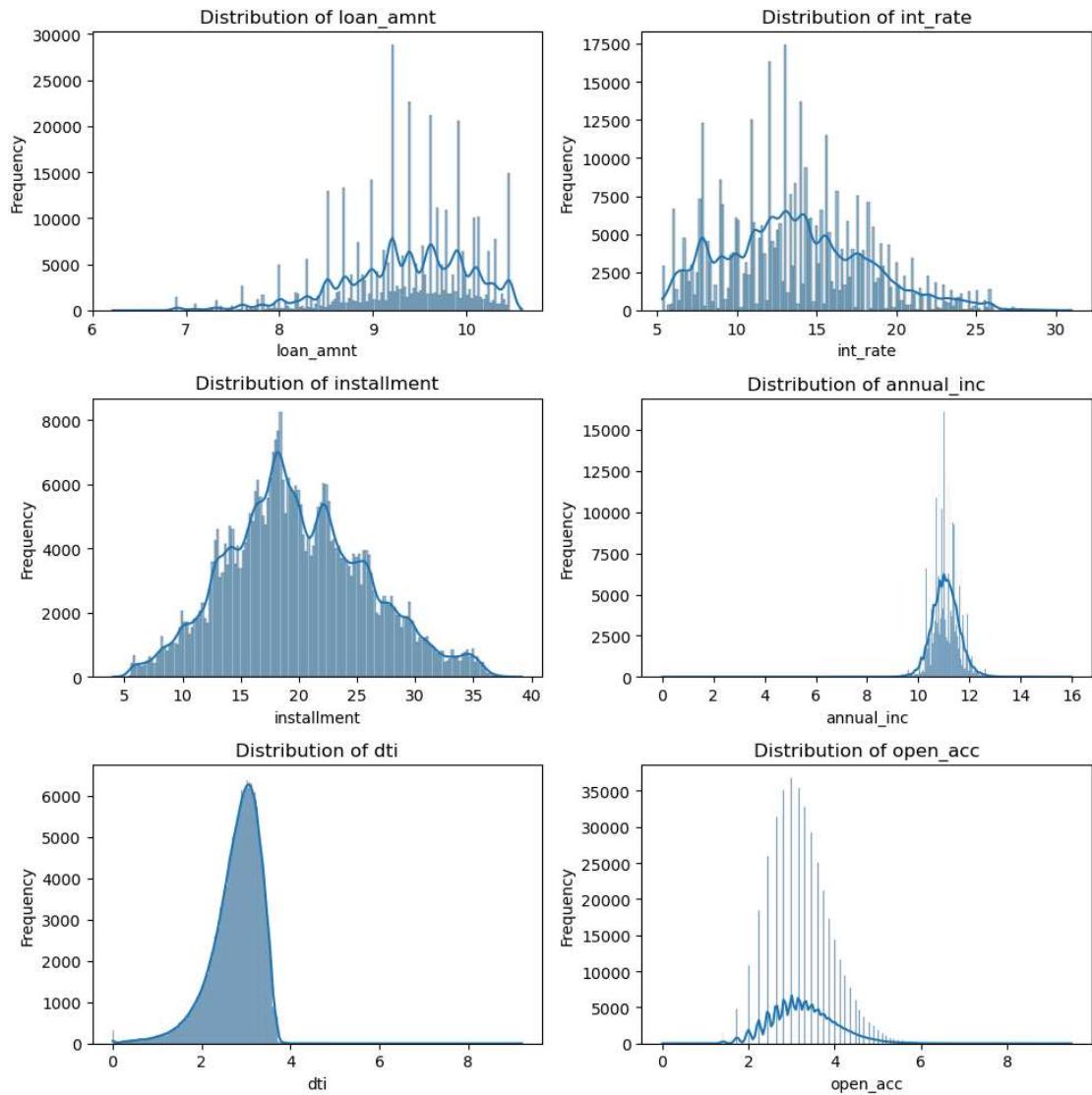
Skewness after transformation:
Skewness of loan_amnt: -0.6680611188267545
Skewness of int_rate: 0.4206694719582165
Skewness of installment: 0.26130316409228166
Skewness of annual_inc: 0.16622501307433268
Skewness of dti: -1.230713217981101
Skewness of open_acc: 0.4269654089750745
Skewness of pub_rec: 2.5354727777923403
Skewness of revol_bal: -2.8518972748189184
Skewness of revol_util: -0.07177802032853098
Skewness of total_acc: 0.17769891689572875
Skewness of mort_acc: 0.29647893480908233
Skewness of pub_rec_bankruptcies: 2.6217121173045226

In [14]: ►

```

1 numerical_features1 = ['loan_amnt', 'int_rate', 'installment', 'annual_inc']
2
3 plt.figure(figsize=(10,10))
4
5 for i,var in enumerate(numerical_features1,1):
6     plt.subplot(3,2,i)
7     sns.histplot(df[var],kde=True)
8     plt.title(f'Distribution of {var}')
9     plt.xlabel(var)
10    plt.ylabel('Frequency')
11 plt.tight_layout()
12 plt.show()

```



Based on these distribution plots, we can make the following observations:

- Loan Amount (loan_amnt):

The distribution appears multimodal, with peaks at certain loan amounts. This suggests that specific loan amounts are more common, possibly due to LoanTap's product design or

► Interest Rate (int_rate):

The interest rate distribution appears right-skewed, with a higher concentration of loans in the 10-20% range. Few loans have interest rates above 25%, indicating that higher rates are less common.

► Installment (installment):

This distribution is also right-skewed, with most installments clustered between 200 and 500. This likely reflects typical monthly payments for the majority of borrowers and suggests that smaller installments are more popular.

► Annual Income (annual_inc):

The annual income distribution is highly right-skewed, with a long tail reaching into the higher income levels. Most borrowers have an income below a certain range (likely under 200,000), but there are a few individuals with very high incomes, possibly skewing the mean.

► Debt-to-Income Ratio (dti):

The dti distribution is also right-skewed, with a large concentration at the lower end (close to zero). This suggests that many borrowers have manageable debt relative to their income. A few outliers with very high dti ratios indicate borrowers who may be over-leveraged, although these cases seem rare.

► Open Accounts (open_acc):

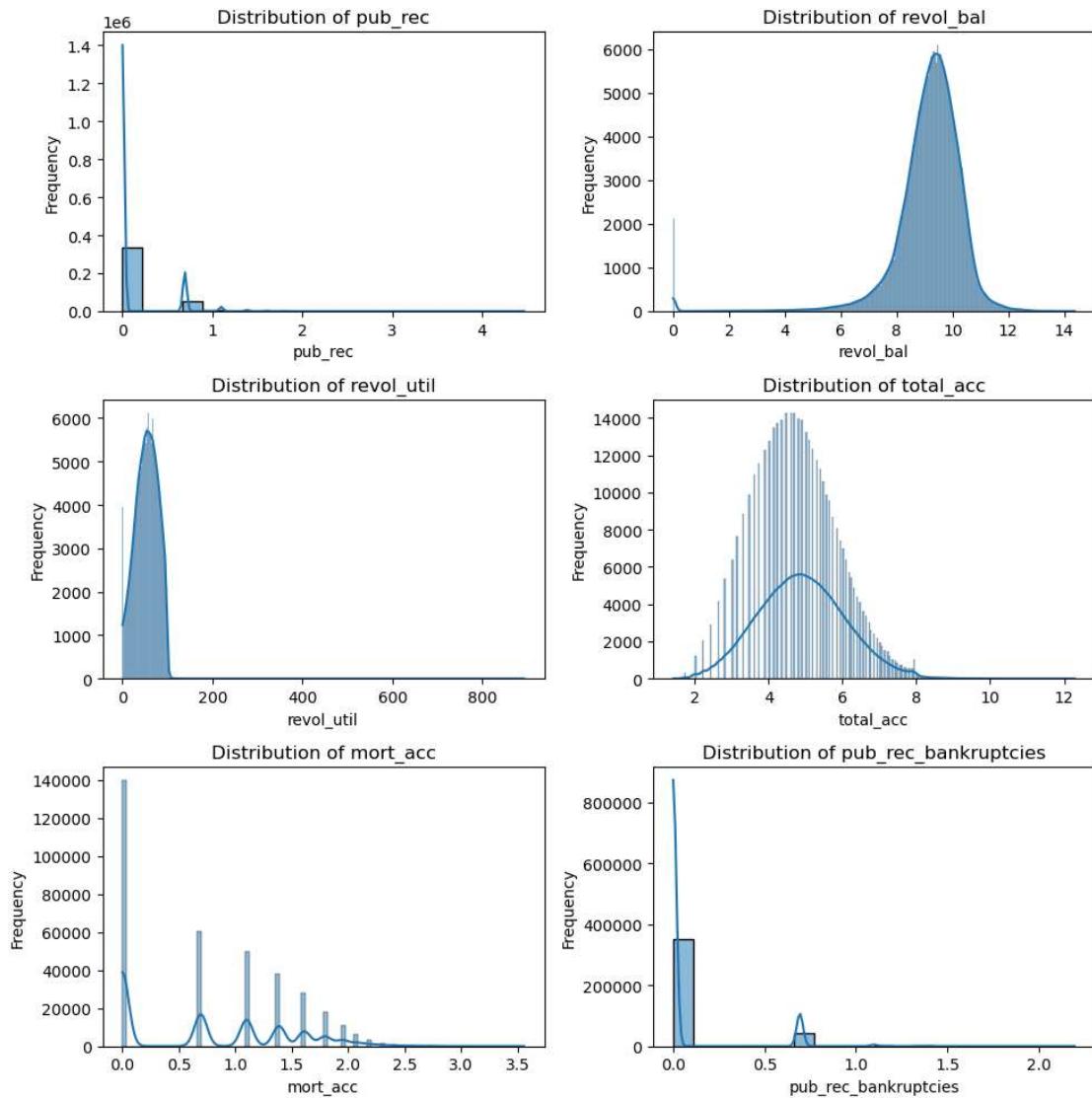
The distribution of open accounts is concentrated below 20, indicating that most borrowers have a relatively small number of open credit lines. There are fewer borrowers with a higher number of open accounts, suggesting that such cases are less common and possibly reflect higher credit activity.

In [15]: ►

```

1 numerical_features2 = ['pub_rec', 'revol_bal', 'revol_util', 'total_a
2
3 plt.figure(figsize=(10,10))
4
5 for i,var in enumerate(numerical_features2,1):
6     plt.subplot(3,2,i)
7     sns.histplot(df[var],kde=True)
8     plt.title(f'Distribution of {var}')
9     plt.xlabel(var)
10    plt.ylabel('Frequency')
11 plt.tight_layout()
12 plt.show()

```



Based on these distribution plots, we can make the following observations:

► pub_rec (Public Records) :

The high positive skew in the pub_rec distribution indicates that most individuals have a clean credit history with few or no public records. However, a small segment of the population has a significantly higher number of public records, which could be due to factors like bankruptcies, foreclosures, or tax liens.

► revol_bal (Revolving Balance) :

The right-skewed distribution of revol_bal suggests that most individuals have low revolving balances (credit card debt), while a few individuals carry substantial credit card debt. This could be attributed to factors like income levels, spending habits, and access to credit.

► revol_util (Revolving Utilization Rate) :

The negative skew in the revol_util distribution indicates that most individuals maintain low credit utilization rates, which is generally considered a positive factor for creditworthiness. However, a smaller group of individuals have high utilization rates, potentially putting them at risk of higher interest rates and reduced credit scores.

► total_acc (Total Number of Accounts) :

The approximately normal distribution of total_acc suggests that the number of total credit accounts is relatively evenly distributed across the population. This indicates a diverse range of credit histories, with some individuals having many accounts and others having fewer.

► mort_acc (Number of Mortgage Accounts) :

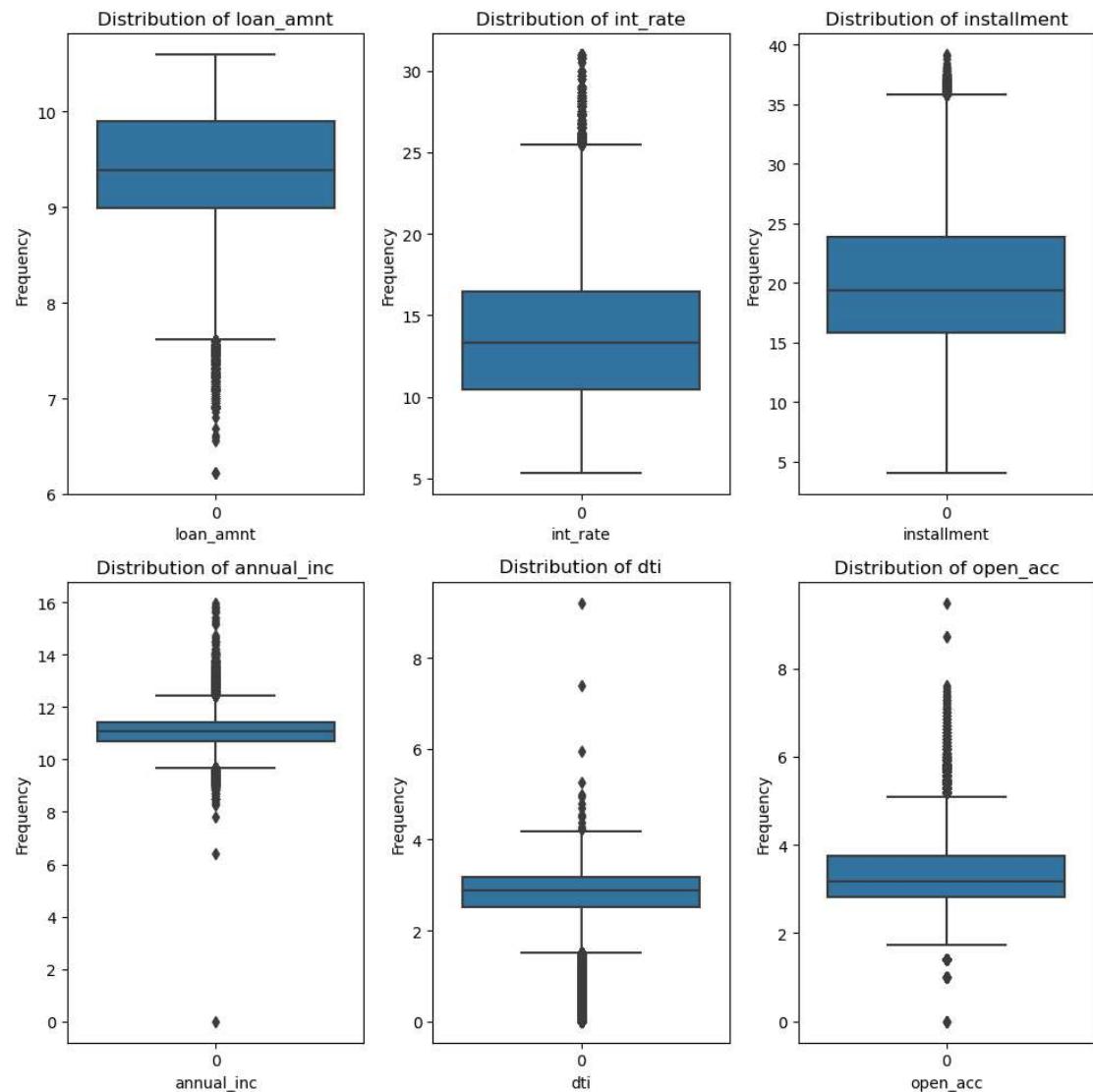
The high positive skew in the mort_acc distribution indicates that most individuals have few or no mortgage accounts. However, a smaller segment of the population has multiple mortgage accounts, possibly due to factors like homeownership, real estate investments, or refinancing.

► pub_rec_bankruptcies (Public Record Bankruptcies) :

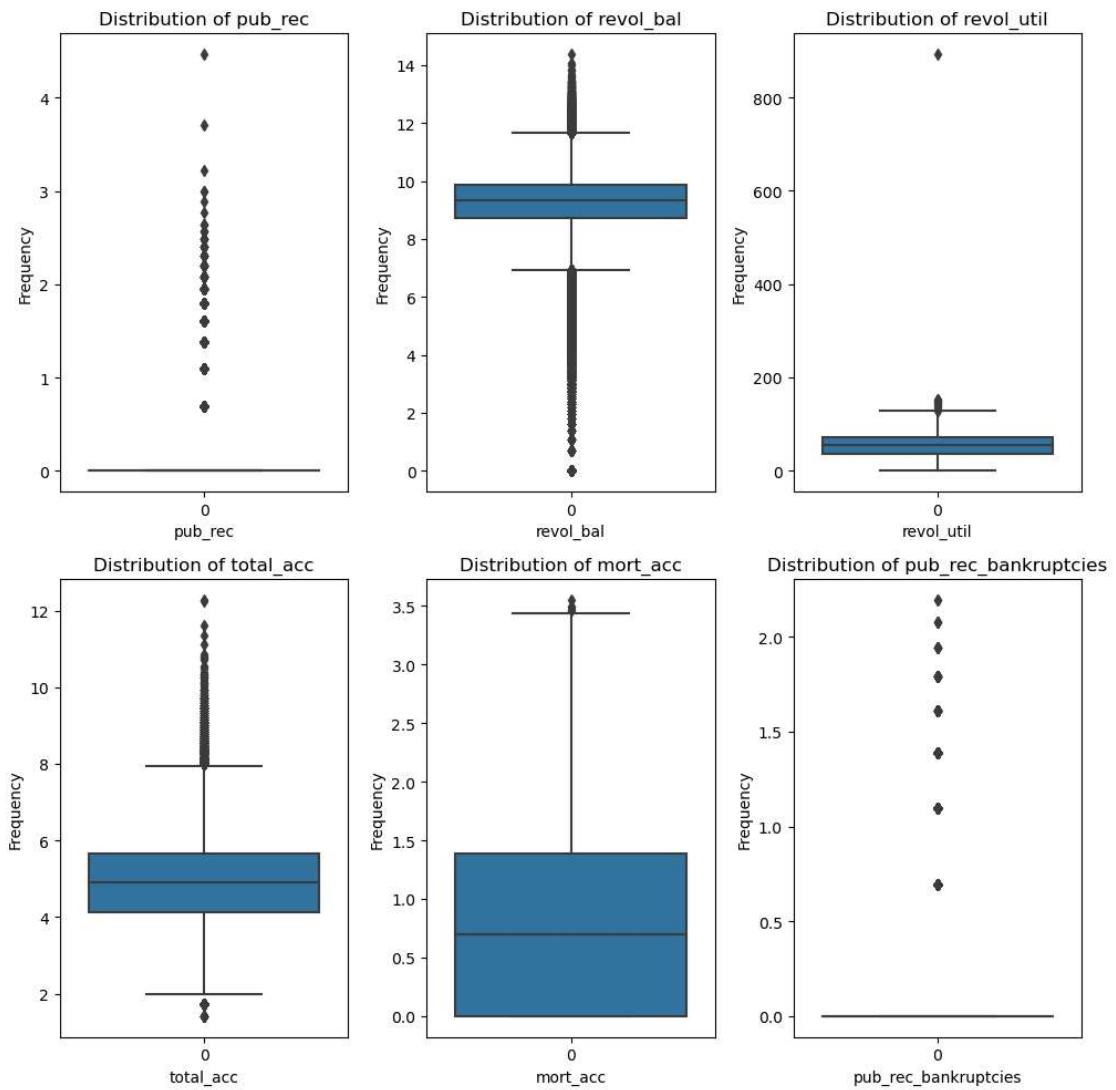
The high positive skew in the pub_rec_bankruptcies distribution indicates that most individuals

In [16]:

```
1 numerical_features1 = ['loan_amnt', 'int_rate', 'installment', 'annual_inc', 'dti', 'open_acc']
2
3 plt.figure(figsize=(10,10))
4
5 for i,var in enumerate(numerical_features1,1):
6     plt.subplot(2,3,i)
7     sns.boxplot(df[var])
8     plt.title(f'Distribution of {var}')
9     plt.xlabel(var)
10    plt.ylabel('Frequency')
11 plt.tight_layout()
12 plt.show()
```



```
In [17]: ┏━
1 numerical_features2 = ['pub_rec', 'revol_bal', 'revol_util', 'total_a
2
3 plt.figure(figsize=(10,10))
4
5 for i,var in enumerate(numerical_features2,1):
6     plt.subplot(2,3,i)
7     sns.boxplot(df[var])
8     plt.title(f'Distribution of {var}')
9     plt.xlabel(var)
10    plt.ylabel('Frequency')
11 plt.tight_layout()
12 plt.show()
```



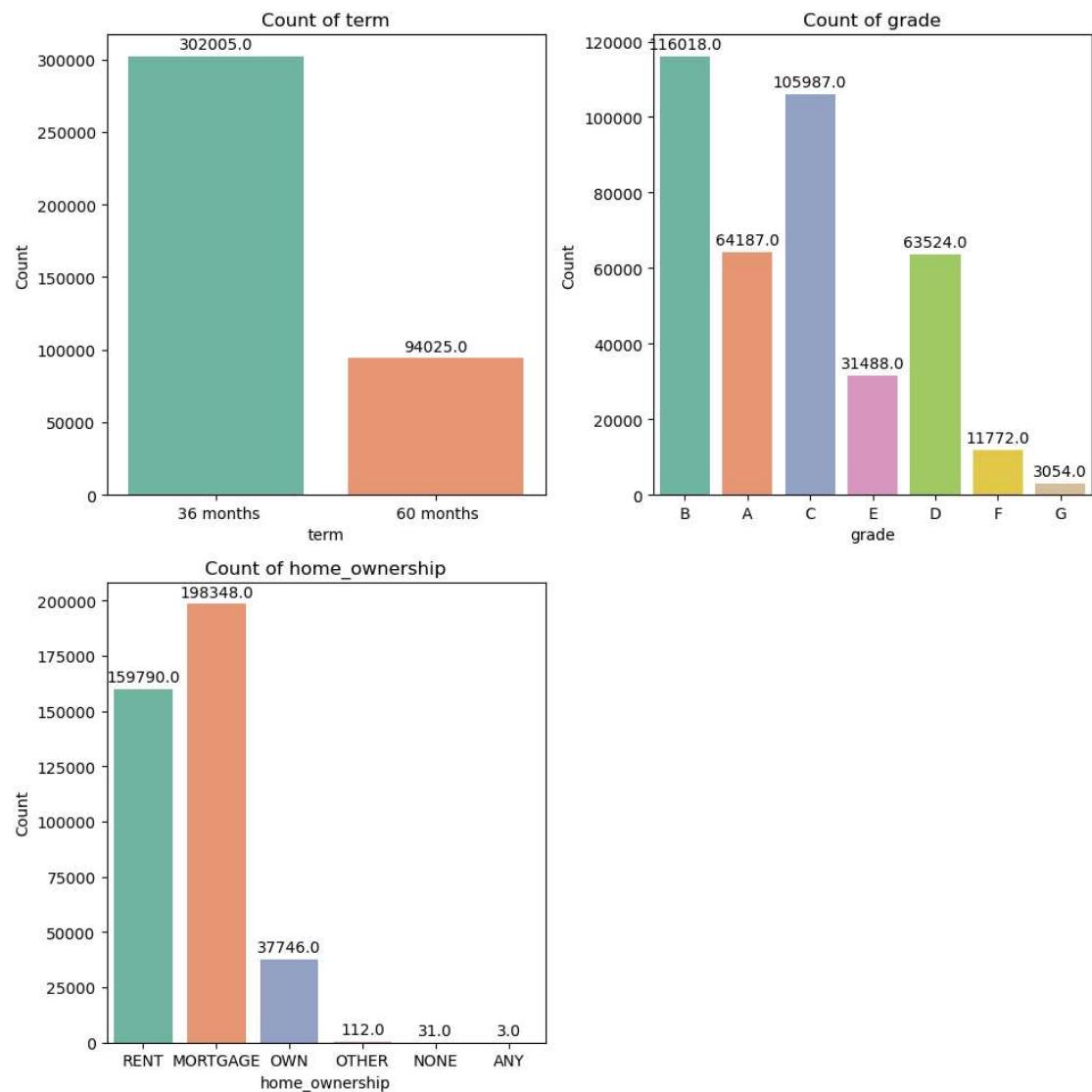
```
In [18]: ┏━
1 categorical_vars = ['term', 'grade', 'emp_length', 'home_ownership',
2                     'loan_status', 'initial_list_status', 'applicati
```

In [19]: ┌

```

1 categorical_vars1 = ['term', 'grade', 'home_ownership']
2
3 plt.figure(figsize=(10, 10))
4 for i, var in enumerate(categorical_vars1, 1):
5     plt.subplot(2, 2, i)
6     ax = sns.countplot(x=var, data=df, palette='Set2')
7     plt.title(f'Count of {var}')
8     plt.xlabel(var)
9     plt.ylabel('Count')
10
11    for p in ax.patches:
12        ax.annotate(f'{p.get_height()}', 
13                    (p.get_x() + p.get_width() / 2., p.get_height()), 
14                    ha='center', va='baseline', fontsize=10, color='black', 
15                    textcoords='offset points')
16 plt.tight_layout()
17 plt.show()

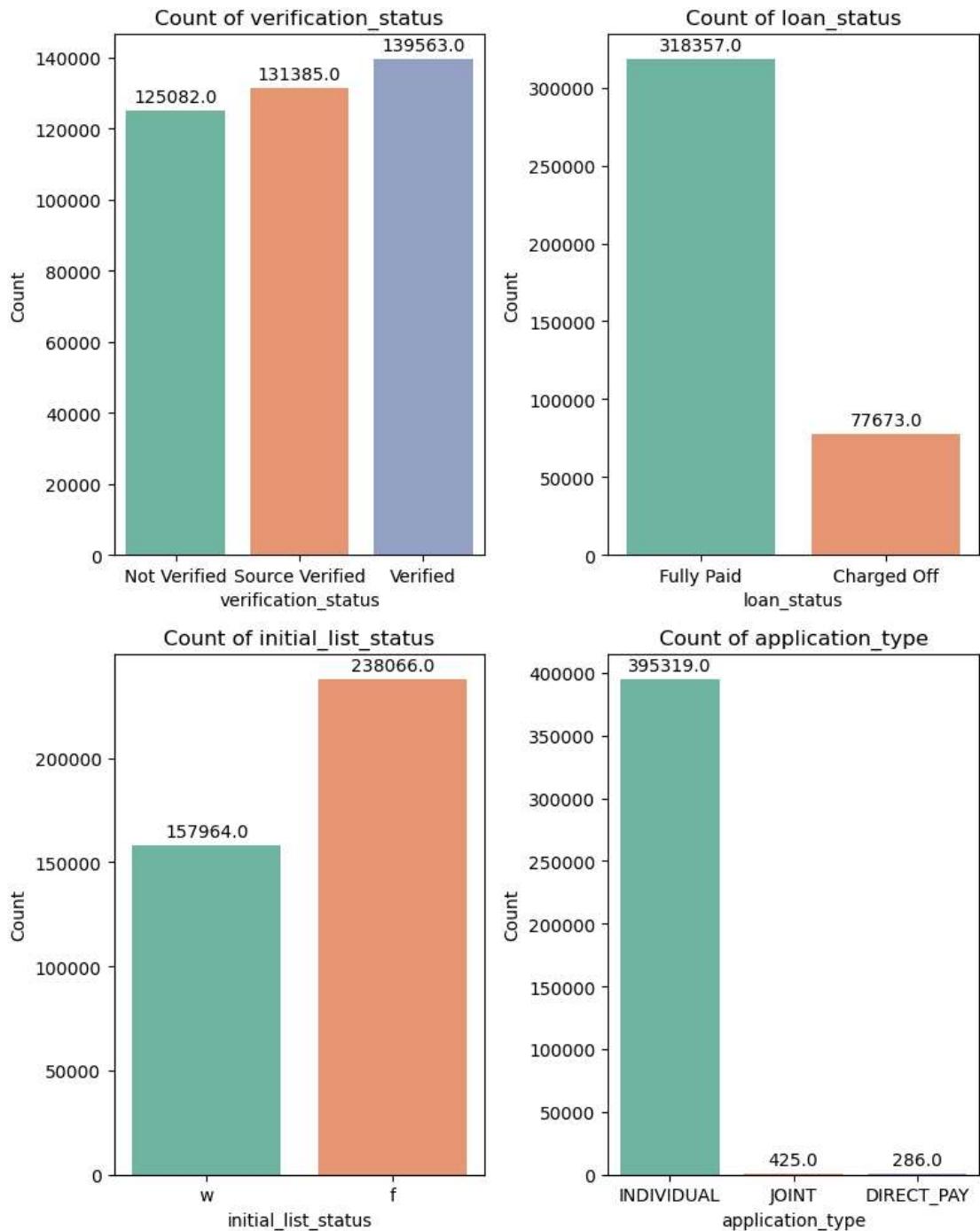
```



Observations :-

- term : The data shows a highly imbalanced distribution between 36 months and 60 months , with 36 months having significantly more observations.The countplot shows majority of people took loan for 36 months rest took for 60 months.
- grade : Grade B has the highest count, while Category G is the least frequent in the dataset.
- home_ownership : The 90 % of the data is distributed among rent, mortage, and own.

```
In [20]: ┏ 1 categorical_vars2 = ['verification_status', 'loan_status', 'initial_l  
2  
3 plt.figure(figsize=(8, 10))  
4 for i, var in enumerate(categorical_vars2, 1):  
5     plt.subplot(2, 2, i)  
6     ax = sns.countplot(x=var, data=df, palette='Set2')  
7     plt.title(f'Count of {var}')  
8     plt.xlabel(var)  
9     plt.ylabel('Count')  
10  
11    for p in ax.patches:  
12        ax.annotate(f'{p.get_height()}',  
13                     (p.get_x() + p.get_width() / 2., p.get_height()),  
14                     ha='center', va='baseline', fontsize=10, color='b  
15                     textcoords='offset points')  
16  
17 plt.tight_layout()  
18 plt.show()
```

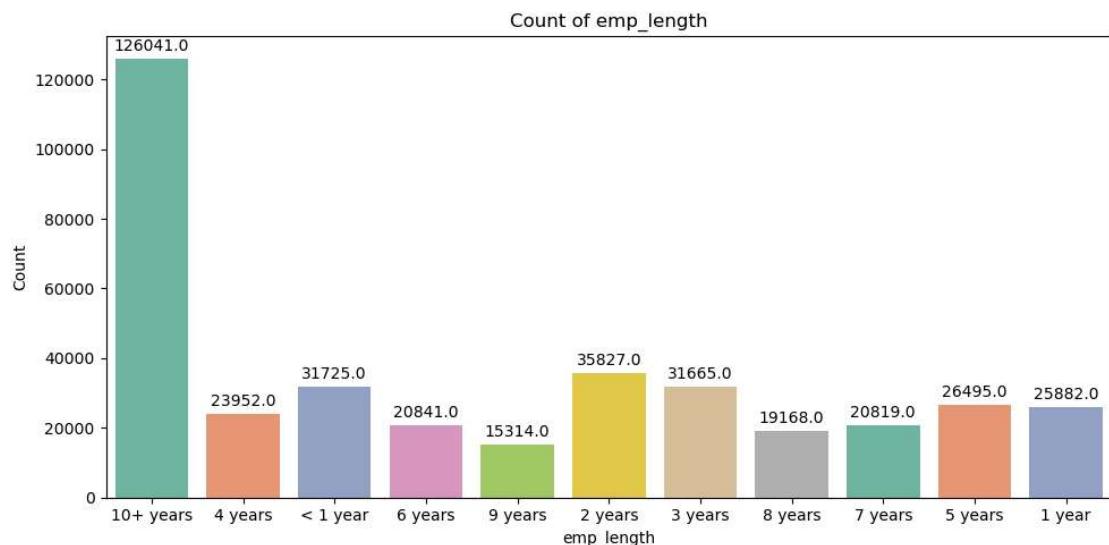


Observations :-

- **verification_status :** The countplot displays an almost uniform distribution across the categories, indicating balanced data. Number of verified people > source verified people > not verified people.
- **loan_status :** The data shows a highly imbalanced distribution between Fully paid and charged off , with most of the people has fully paid their loan.
- **initial_list_status :** Majority of people having only part (F) of the loan is available for funding or investment at first where other people is listed as a whole loan (W), and it is available for full funding or investment.

► application_type : Almost 99 % of people has application type as individual. Where the number of people with Joint and direct pay has very few data points, which might indicate an outlier or an uncommon event

```
In [21]: ┌─ 1 plt.figure(figsize=(10, 5))
2
3 ax = sns.countplot(x=df['emp_length'], data=df, palette='Set2')
4 plt.title('Count of emp_length')
5 plt.xlabel('emp_length')
6 plt.ylabel('Count')
7
8 for p in ax.patches:
9     ax.annotate(f'{p.get_height()}', (p.get_x() + p.get_width() / 2., p.get_height()),
10                 ha='center', va='baseline', fontsize=10, color='black',
11                 textcoords='offset points')
12
13 plt.tight_layout()
14 plt.show()
```



Observation :-

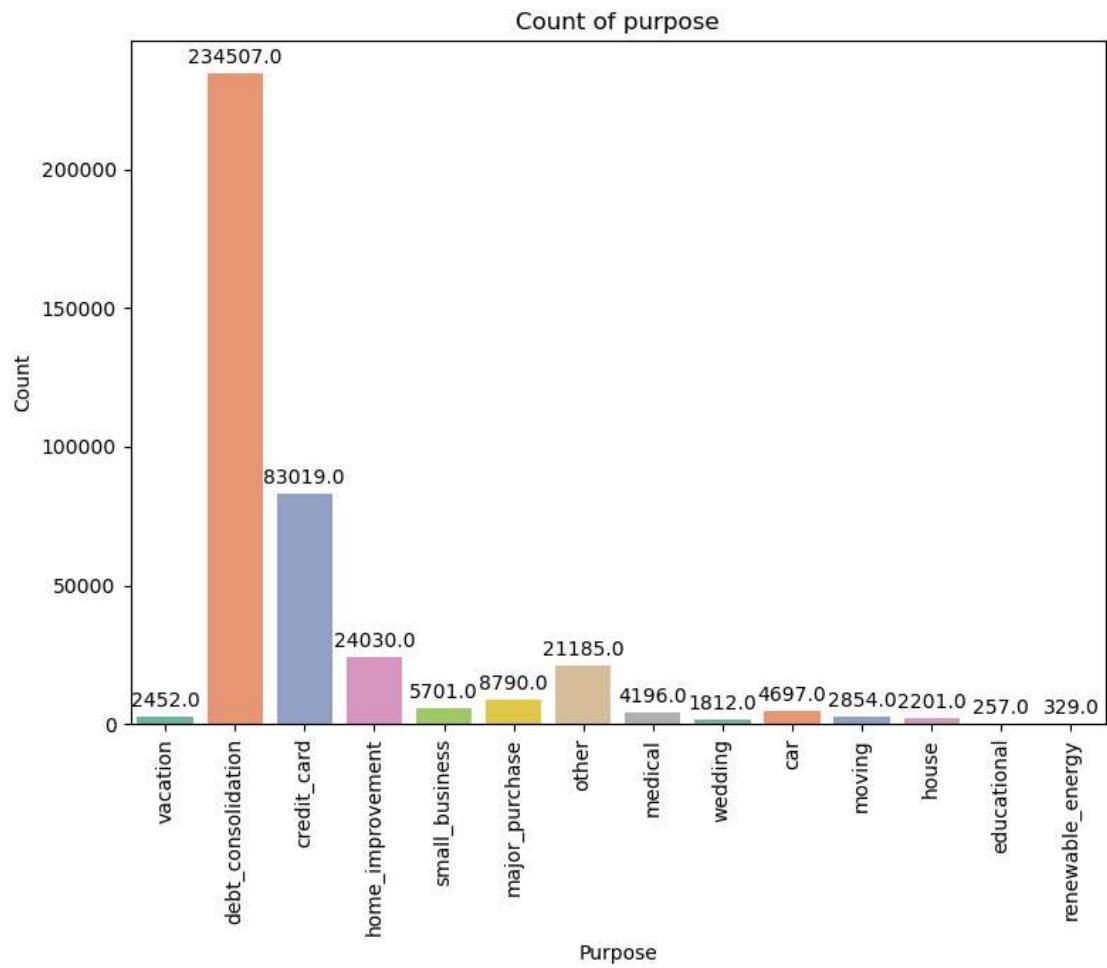
People having 10+ years experience has the highest count, while People having 9 years is the least frequent in the dataset. The countplot shows a right-skewed distribution with a large number of instances in 10+ years, moderate between <1 years, 2 years, 3 years, 5 years, 1 year, 4 years, 6 years and fewer in 8 years, 9 years suggesting that the model may have a bias towards People having 10+ years experience.

In [22]:

```

1 plt.figure(figsize=(8, 7))
2
3 ax = sns.countplot(x=df['purpose'], data=df, palette='Set2')
4 plt.title('Count of purpose')
5 plt.xlabel('Purpose')
6 plt.ylabel('Count')
7 plt.xticks(rotation=90)
8
9 for p in ax.patches:
10     ax.annotate(f'{p.get_height():.0f}', 
11                 (p.get_x() + p.get_width() / 2., p.get_height()),
12                 ha='center', va='baseline', fontsize=10, color='black',
13                 textcoords='offset points')
14 plt.tight_layout()
15 plt.show()

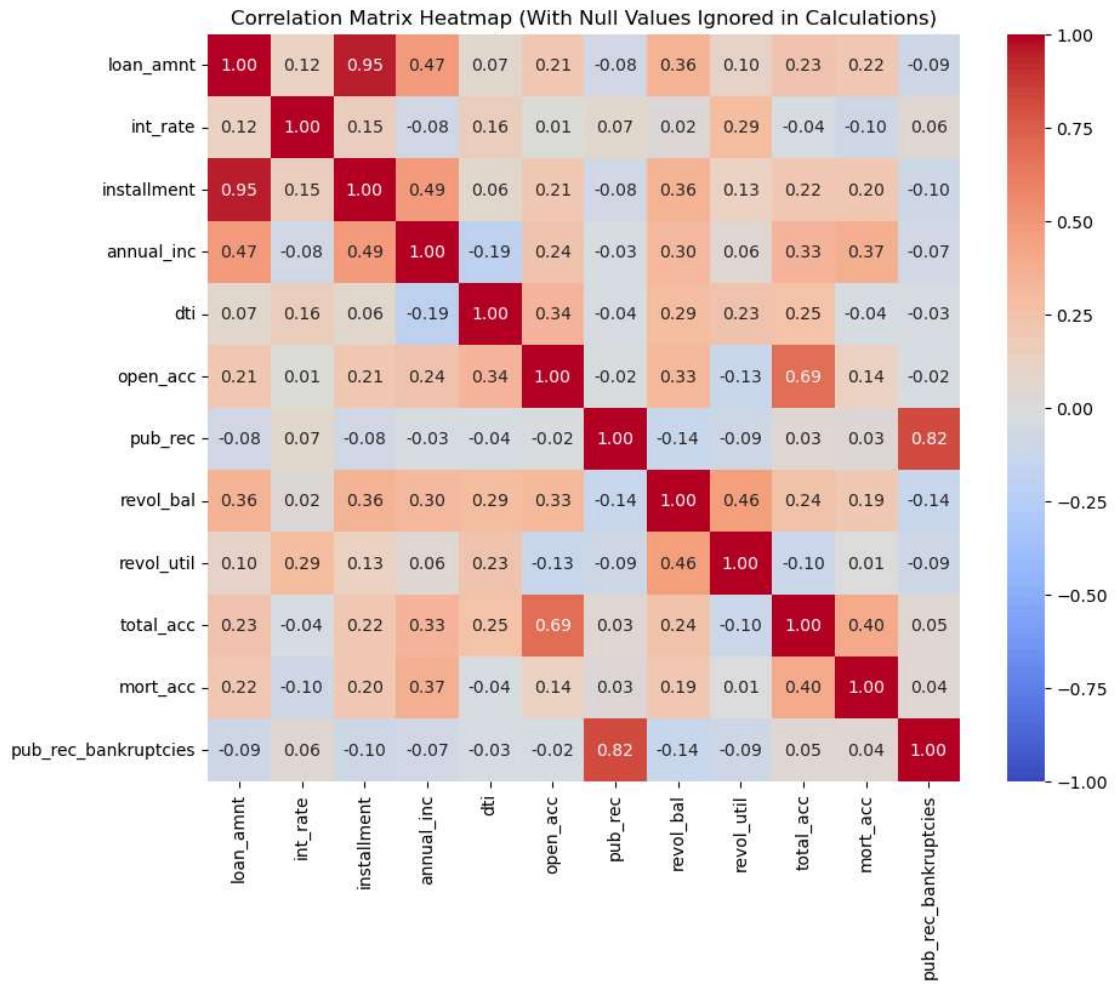
```



Observation :-

debt consolidation has the highest count, while educational is the least frequent in the dataset. small business, mdeical , car , moving ,house has very few data points, vacation ,wedding, educational, renewable energy might indicate an outlier.

```
In [30]: 1 correlation_matrix = df[numerical_features].corr()
2
3 plt.figure(figsize=(10, 8))
4 sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f")
5 plt.title("Correlation Matrix Heatmap (With Null Values Ignored in Ca")
6 plt.show()
```



→ 2. Data Preprocessing

a) Duplicate value check

```
In [31]: 1 df.duplicated().sum()
```

Out[31]: 0

b) Missing value treatment

```
In [32]: 1 null_percentage = df.isnull().sum()
2 null_percentage.sort_values(ascending = False).head(8)
```

```
Out[32]: mort_acc           37795
emp_title          22927
emp_length         18301
title              1756
pub_rec_bankruptcies   535
revol_util          276
loan_amnt            0
dti                  0
dtype: int64
```

```
In [33]: 1 df.shape
```

```
Out[33]: (396030, 27)
```

```
In [34]: 1 df['revol_util'] = df.groupby('grade')['revol_util'].transform(lambda
```

```
In [35]: 1 df['pub_rec_bankruptcies'] = df.groupby('grade')['pub_rec_bankruptcie
```

```
In [36]: 1 df['title'] = df.groupby('grade')['title'].transform(lambda x: x.fill
```

```
In [37]: 1 df['emp_length'] = df.groupby('grade')['emp_length'].transform(lambda
```

```
In [38]: 1 df['emp_title'] = df.groupby(['grade', 'emp_length'])['emp_title'].tr
```

```
In [39]: 1 df['mort_acc'] = df.groupby('grade')['mort_acc'].transform(lambda x:
```

```
In [40]: 1 null_percentage = df.isnull().sum()
2 null_percentage.sort_values(ascending = False).head(8)
```

```
Out[40]: loan_amnt          0
title                0
pub_rec_bankruptcies 0
mort_acc              0
application_type      0
initial_list_status    0
total_acc              0
revol_util              0
dtype: int64
```

c) Outlier treatment

```
In [41]: 1 numerical_columns = []
2 for i in df.columns:
3     if (df[i].dtypes==float):
4         numerical_columns.append(i)
5
6 print(numerical_columns)
```

```
['loan_amnt', 'int_rate', 'installment', 'annual_inc', 'dti', 'open_ac-
c', 'pub_rec', 'revol_bal', 'revol_util', 'total_acc', 'mort_acc', 'pub-
rec_bankruptcies']
```

```
In [42]: 1 from scipy import stats
2
3 z_scores = stats.zscore(df[numerical_columns])
4 z_scores_df = pd.DataFrame(z_scores, columns=numerical_columns)
5 threshold = 3 # Typically values greater than 3 are considered outliers
6 outliers = (z_scores_df.abs() > threshold)
7 outlier_rows = df[outliers.any(axis=1)]
8
9 outlier_rows.shape
```

Out[42]: (25473, 27)

```
In [43]: 1 df_cleaned = df[~outliers.any(axis=1)]
2 df_cleaned.shape
```

Out[43]: (370557, 27)

```
In [44]: 1 df_cleaned.columns
```

```
Out[44]: Index(['loan_amnt', 'term', 'int_rate', 'installment', 'grade', 'sub_grade',
       'emp_title', 'emp_length', 'home_ownership', 'annual_inc',
       'verification_status', 'issue_d', 'loan_status', 'purpose', 'titl-
e',
       'dti', 'earliest_cr_line', 'open_acc', 'pub_rec', 'revol_bal',
       'revol_util', 'total_acc', 'initial_list_status', 'application_ty-
pe',
       'mort_acc', 'pub_rec_bankruptcies', 'address'],
      dtype='object')
```

In [45]: 1 df_cleaned

Out[45]:

	loan_amnt	term	int_rate	installment	grade	sub_grade	emp_title	emp_len
0	9.210440	36 months	11.44	18.151584	B	B4	Marketing	10+ years
1	8.987322	36 months	11.99	16.299693	B	B5	Credit analyst	4 years
2	9.655090	36 months	10.49	22.515994	B	B3	Statistician	< 1 year
3	8.881975	36 months	6.49	14.854292	A	A2	Client Advocate	6 years
4	10.101354	60 months	17.27	24.684611	C	C5	Destiny Management Inc.	9 years
...
396025	9.210440	60 months	10.99	14.743812	B	B4	licensed bankere	2 years
396026	9.952325	36 months	12.29	26.465449	C	C1	Agent	5 years
396027	8.517393	36 months	9.99	12.701181	B	B1	City Carrier	10+ years
396028	9.952325	60 months	15.31	22.428107	C	C2	Gracon Services, Inc	10+ years
396029	7.601402	36 months	13.61	8.244998	C	C2	Internal Revenue Service	10+ years

370557 rows × 27 columns

d) Feature engineering

In [46]: 1 # Convert to datetime format
2 df['issue_d'] = pd.to_datetime(df['issue_d'], format='%b-%Y')
3 df['earliest_cr_line'] = pd.to_datetime(df['earliest_cr_line'], forma

In [47]: 1 # Extract month and year
2 df['loan_issue_month'] = df['issue_d'].dt.month
3 df['loan_issue_year'] = df['issue_d'].dt.year
4 df['earliest_cr_line_month'] = df['earliest_cr_line'].dt.month
5 df['earliest_cr_line_year'] = df['earliest_cr_line'].dt.year

```
In [48]: ┌ 1 df['credit_history_length_years'] = (df['issue_d'] - df['earliest_cr_
```

```
In [49]: ┌ 1 df['credit_history_length_years'] = df['credit_history_length_years']
```

```
In [50]: ┌ 1 df_cleaned.drop(['issue_d', 'earliest_cr_line'], axis=1, inplace=True
```

C:\Users\Prajwal\AppData\Local\Temp\ipykernel_2948\3795973535.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df_cleaned.drop(['issue_d', 'earliest_cr_line'], axis=1, inplace=True)
```

e) Data Preparation for Modeling.

```
In [51]: ┌ 1 categorical_columns = df.select_dtypes(include=['object', 'category'])
 2
 3 for col in categorical_columns:
 4     unique_count = df[col].nunique()
 5     if unique_count < 8:
 6         print(f"{col}: {unique_count}")
```

term: 2
 grade: 7
 home_ownership: 6
 verification_status: 3
 loan_status: 2
 initial_list_status: 2
 application_type: 3

```
In [52]: ┌ 1 # Columns to one-hot encode
 2 columns_to_encode = ['grade', 'home_ownership', 'verification_status']
 3 df_cleaned = pd.get_dummies(df_cleaned, columns=columns_to_encode, dr
```

```
In [53]: ┌ 1 # One-hot encoding for the 'term' column
 2 df_cleaned = pd.get_dummies(df_cleaned, columns=['term'], prefix='ter
```

In [54]:

```

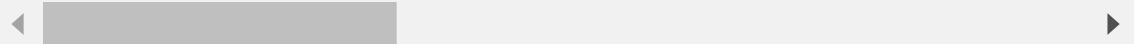
1 # Apply One-Hot Encoding to 'purpose' and 'emp_length'
2 df_cleaned = pd.get_dummies(df_cleaned, columns=['purpose', 'emp_length'])
3 df_cleaned.head(5)

```

Out[54]:

	loan_amnt	int_rate	installment	sub_grade	emp_title	annual_inc	title	dti	open_
0	9.210440	11.44	18.151584	B4	Marketing	11.669938	Vacation	3.304686	4.000000
1	8.987322	11.99	16.299693	B5	Credit analyst	11.082158	Debt consolidation	3.137666	4.123944
2	9.655090	10.49	22.515994	B3	Statistician	10.670303	Credit card refinancing	2.623944	3.605136
3	8.881975	6.49	14.854292	A2	Client Advocate	10.896758	Credit card refinancing	1.280934	2.449444
4	10.101354	17.27	24.684611	C5	Destiny Management Inc.	10.915107	Credit Card Refinance	3.553918	3.605136

5 rows × 57 columns



In [55]:

```

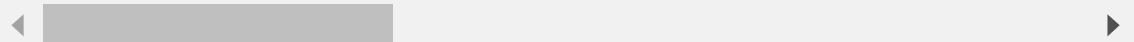
1 from sklearn.preprocessing import LabelEncoder
2 label_encoder = LabelEncoder()
3 df_cleaned['sub_grade_encoded'] = label_encoder.fit_transform(df_cleaned['sub_grade'])
4 df_cleaned.drop(['sub_grade'], axis=1, inplace=True)
5 df_cleaned.head()

```

Out[55]:

	loan_amnt	int_rate	installment	emp_title	annual_inc	title	dti	open_
0	9.210440	11.44	18.151584	Marketing	11.669938	Vacation	3.304686	4.000000
1	8.987322	11.99	16.299693	Credit analyst	11.082158	Debt consolidation	3.137666	4.123944
2	9.655090	10.49	22.515994	Statistician	10.670303	Credit card refinancing	2.623944	3.605136
3	8.881975	6.49	14.854292	Client Advocate	10.896758	Credit card refinancing	1.280934	2.449444
4	10.101354	17.27	24.684611	Destiny Management Inc.	10.915107	Credit Card Refinance	3.553918	3.605136

5 rows × 57 columns



```
In [56]: 1 categorical_columns = df_cleaned.select_dtypes(include=['object', 'ca
2
3 for col in categorical_columns:
4     unique_count = df_cleaned[col].nunique()
5     print(f"{col}: {unique_count}")
```

emp_title: 164695
title: 45763
address: 368494

```
In [57]: 1 # Frequency encoding for emp_title
2 emp_title_counts = df_cleaned['emp_title'].value_counts()
3 df_cleaned['emp_title_encoded'] = df_cleaned['emp_title'].map(emp_tit
4
5 # Frequency encoding for title
6 title_counts = df_cleaned['title'].value_counts()
7 df_cleaned['title_encoded'] = df_cleaned['title'].map(title_counts)
```

```
In [58]: 1 # Drop original 'emp_title' and 'title' columns
2 df_cleaned.drop(['emp_title', 'title', 'address'], axis=1, inplace=True)
```

```
In [59]: 1 df_cleaned.sample(5)
```

Out[59]:

	loan_amnt	int_rate	installment	annual_inc	dti	open_acc	pub_rec	revol_
54685	9.680406	18.49	20.262774	10.778977	2.520113	3.000000	0.693147	8.4648
79362	10.463132	13.98	28.531211	11.407576	3.154444	3.000000	0.000000	11.4780
128629	9.974691	21.18	24.148499	10.933125	3.180551	3.464102	0.000000	8.7544
2342	10.239996	11.14	30.307590	11.451061	1.803359	3.162278	0.000000	10.2786
230457	9.952325	11.14	26.247095	11.849405	1.791759	2.449490	0.000000	8.9316

5 rows × 56 columns



```
In [60]: 1 categorical_columns = df_cleaned.select_dtypes(include=['object', 'ca
2
3 for col in categorical_columns:
4     unique_count = df_cleaned[col].nunique()
5     print(f"{col}: {unique_count}")
```

```
In [61]: 1 # Select boolean columns and convert them to 0 and 1
2 bool_columns = df_cleaned.select_dtypes(include=['bool']).columns
3 df_cleaned[bool_columns] = df_cleaned[bool_columns].astype(int)
```

In [62]: 1 df_cleaned.dtypes

```
Out[62]: loan_amnt                      float64
          int_rate                       float64
          installment                     float64
          annual_inc                      float64
          dti                            float64
          open_acc                        float64
          pub_rec                         float64
          revol_bal                       float64
          revol_util                      float64
          total_acc                        float64
          mort_acc                        float64
          pub_rec_bankruptcies            float64
          grade_B                          int32
          grade_C                          int32
          grade_D                          int32
          grade_E                          int32
          grade_F                          int32
          grade_G                          int32
          home_ownership_MORTGAGE         int32
          home_ownership_NONE             int32
          home_ownership_OTHER             int32
          home_ownership_OWN              int32
          home_ownership_RENT             int32
          verification_status_Source     Verified int32
          verification_status_Verified    int32
          loan_status_Fully_Paid          int32
          initial_list_status_w           int32
          application_type_INDIVIDUAL    int32
          application_type_JOINT          int32
          term_ 60_months                 int32
          purpose_credit_card              int32
          purpose_debt_consolidation      int32
          purpose_educational              int32
          purpose_home_improvement        int32
          purpose_house                     int32
          purpose_major_purchase           int32
          purpose_medical                  int32
          purpose_moving                   int32
          purpose_other                     int32
          purpose_renewable_energy         int32
          purpose_small_business            int32
          purpose_vacation                  int32
          purpose_wedding                  int32
          emp_length_10+ years             int32
          emp_length_2 years               int32
          emp_length_3 years               int32
          emp_length_4 years               int32
          emp_length_5 years               int32
          emp_length_6 years               int32
          emp_length_7 years               int32
          emp_length_8 years               int32
          emp_length_9 years               int32
          emp_length_< 1 year              int32
          sub_grade_encoded                int32
          emp_title_encoded                int64
```

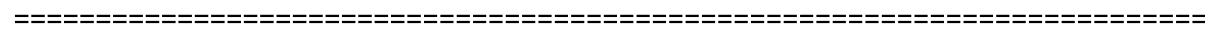
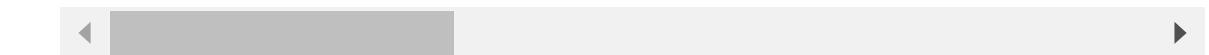
```
title_encoded                               int64
dtype: object
```

In [63]: 1 df_cleaned.sample(5)

Out[63]:

	loan_amnt	int_rate	installment	annual_inc	dti	open_acc	pub_rec	revol_
84162	7.824446	10.49	9.013878	11.407576	2.836150	4.242641	0.000000	10.7156
284167	9.903538	7.49	24.940730	11.925042	3.108614	3.605551	0.000000	11.4176
347691	10.257694	24.99	28.919716	11.225257	2.776332	3.162278	0.000000	9.1426
56394	8.987322	13.33	16.456913	10.373522	3.648057	3.162278	0.693147	9.8410
142721	9.210440	17.77	18.983677	11.722384	2.901422	3.162278	0.000000	9.7806

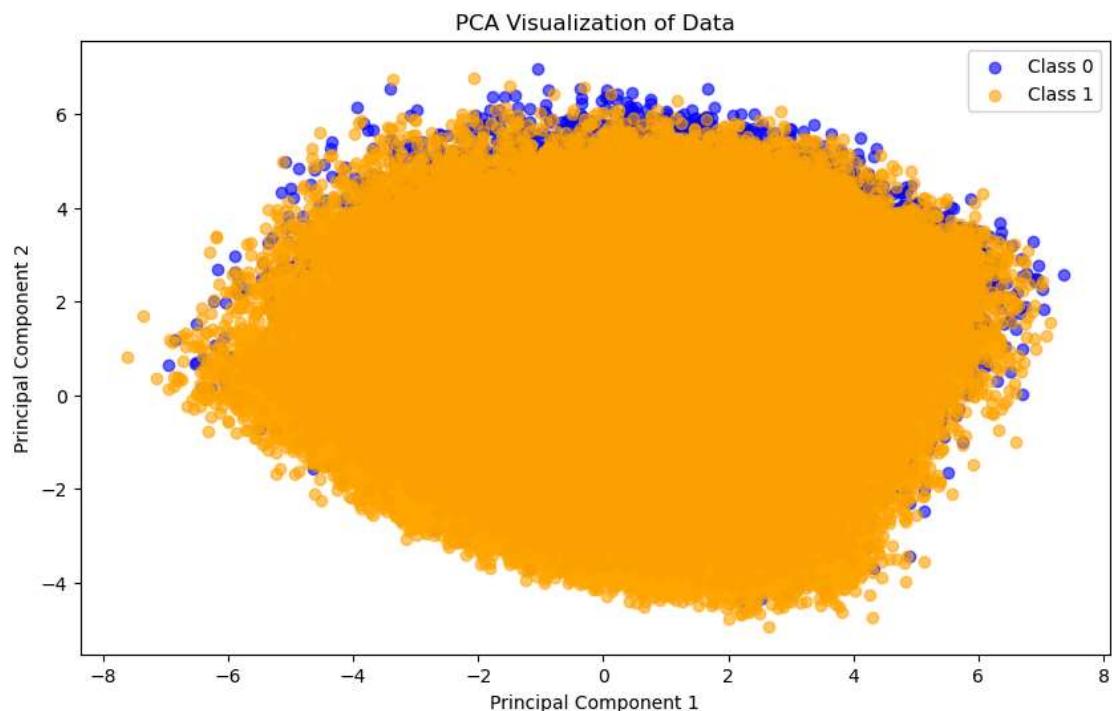
5 rows × 56 columns



→ 3. Model building

In [64]:

```
1 from sklearn.decomposition import PCA
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import StandardScaler
4
5 # Drop the target column and scale the features
6 X = df_cleaned.drop(columns=['loan_status_Fully Paid'])
7 y = df_cleaned['loan_status_Fully Paid']
8
9 scaler = StandardScaler()
10 X_scaled = scaler.fit_transform(X)
11
12 pca = PCA(n_components=2)
13 X_pca = pca.fit_transform(X_scaled)
14
15 plt.figure(figsize=(10, 6))
16 plt.scatter(X_pca[y == 0, 0], X_pca[y == 0, 1], alpha=0.6, label='Class 0')
17 plt.scatter(X_pca[y == 1, 0], X_pca[y == 1, 1], alpha=0.6, label='Class 1')
18
19 plt.title('PCA Visualization of Data')
20 plt.xlabel('Principal Component 1')
21 plt.ylabel('Principal Component 2')
22 plt.legend()
23 plt.show()
```

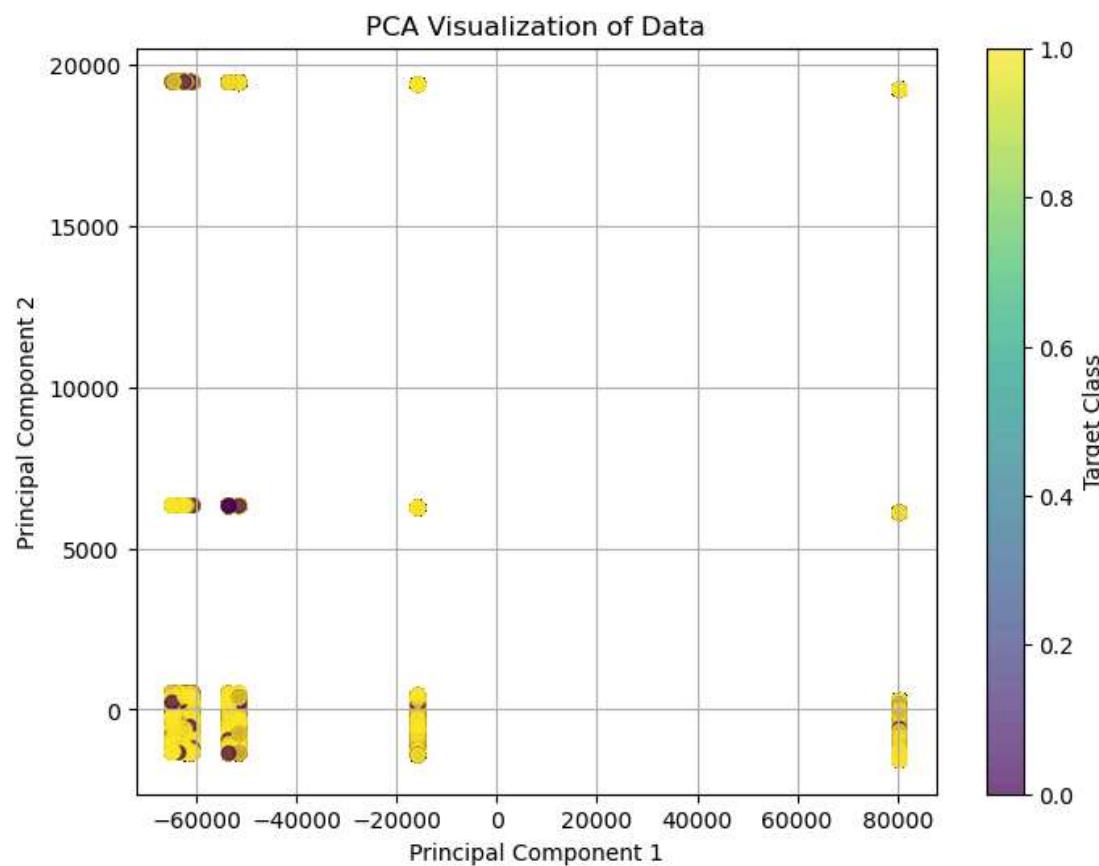


In [65]:

```

1 import matplotlib.pyplot as plt
2 from sklearn.decomposition import PCA
3
4 pca = PCA(n_components=2)
5 visualize = pca.fit_transform(X)
6
7 def scatter_visualize(visualize, target):
8     plt.figure(figsize=(8, 6))
9     scatter = plt.scatter(visualize[:, 0], visualize[:, 1], c=target,
10                           plt.colorbar(scatter, label='Target Class')
11                           plt.title('PCA Visualization of Data')
12                           plt.xlabel('Principal Component 1')
13                           plt.ylabel('Principal Component 2')
14                           plt.grid(True)
15                           plt.show()
16
17 scatter_visualize(visualize, y)

```



In [66]:

```
1 df_cleaned['loan_status_Fully Paid'].value_counts()
```

Out[66]:

loan_status_Fully Paid	count
1	297728
0	72829

Name: count, dtype: int64

In [67]:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0
5
6 scaler = StandardScaler()
7 X_train_scaled = scaler.fit_transform(X_train)
8 X_test_scaled = scaler.transform(X_test)

```

In [68]:

```

1 from imblearn.over_sampling import SMOTE
2
3 smote = SMOTE(random_state=42)
4 X_train_resampled, y_train_resampled = smote.fit_resample(X_train_sca
5 from collections import Counter
6 print("Before SMOTE:", Counter(y_train))
7 print("After SMOTE:", Counter(y_train_resampled))

```

Before SMOTE: Counter({1: 238182, 0: 58263})

After SMOTE: Counter({1: 238182, 0: 238182})

In [69]:

```

1 from sklearn.feature_selection import RFE
2 from sklearn.linear_model import LogisticRegression
3
4 model = LogisticRegression()
5
6 rfe = RFE(estimator=model, n_features_to_select=20) # Select the top
7 X_train_selected = rfe.fit_transform(X_train_resampled, y_train_resam
8
9 selected_features = rfe.support_
10 print("Selected Features: ", selected_features)
11
12 X_test_selected = rfe.transform(X_test_scaled)

```

Selected Features: [True True True True True True False True Tr
ue True False False
True True True False False False False False False True False False
False False False False True False False False True False False False
False True False True False False False False False False False False
False False False False True True True]

In [70]:

```
1 model.fit(X_train_selected, y_train_resampled)
```

Out[70]:

▼ LogisticRegression ⓘ ⓘ
https://scikit-learn.org/1.5/modules/generated/sklearn.linear_model.LogisticRegression.html

```
In [71]: 1 y_pred = model.predict(X_test_selected)
2
3 from sklearn.metrics import classification_report
4 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.31	0.66	0.43	14566
1	0.89	0.65	0.75	59546
accuracy			0.65	74112
macro avg	0.60	0.65	0.59	74112
weighted avg	0.77	0.65	0.68	74112

```
In [72]: 1 feature_names = X.columns
2 selected_feature_names = feature_names[rfe.support_]
3
4 coefficients = model.coef_[0]
5
6 coefficients_df = pd.DataFrame({
7     'Feature': selected_feature_names,
8     'Coefficient': coefficients
9 })
10
11 coefficients_df = coefficients_df.sort_values(by='Coefficient', key=abs)
12
13 print(coefficients_df)
```

	Feature	Coefficient
17	sub_grade_encoded	-0.768521
0	loan_amnt	-0.501306
2	installment	0.297446
3	annual_inc	0.274388
1	int_rate	0.221586
5	open_acc	-0.210705
4	dti	-0.203489
8	total_acc	0.158643
10	grade_C	-0.149826
7	revol_util	-0.147910
9	grade_B	-0.125599
12	home_ownership_RENT	-0.116720
19	title_encoded	-0.101696
11	grade_D	-0.089126
13	term_ 60 months	-0.081052
16	purpose_small_business	-0.077797
18	emp_title_encoded	-0.076978
6	revol_bal	0.076745
14	purpose_home_improvement	-0.070124
15	purpose_other	-0.054591

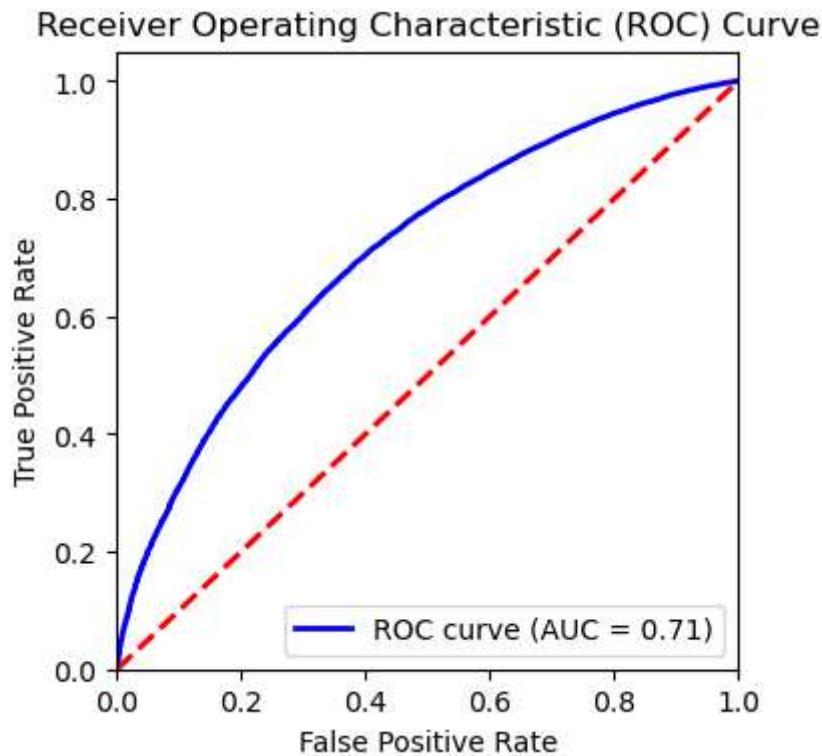
=====



→ 4. Results Evaluation

1. ROC AUC Curve & comments

```
In [77]: X_test_selected = rfe.transform(X_test_scaled)
y_prob = model.predict_proba(X_test_selected)[:, 1]
fpr, tpr, thresholds = roc_curve(y_test, y_prob)
roc_auc = roc_auc_score(y_test, y_prob)
plt.figure(figsize=(4, 4))
plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc:.2f})')
plt.plot([0, 1], [0, 1], color='red', linestyle='--', lw=2) # Diagonal line
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend(loc='lower right')
plt.show()
```



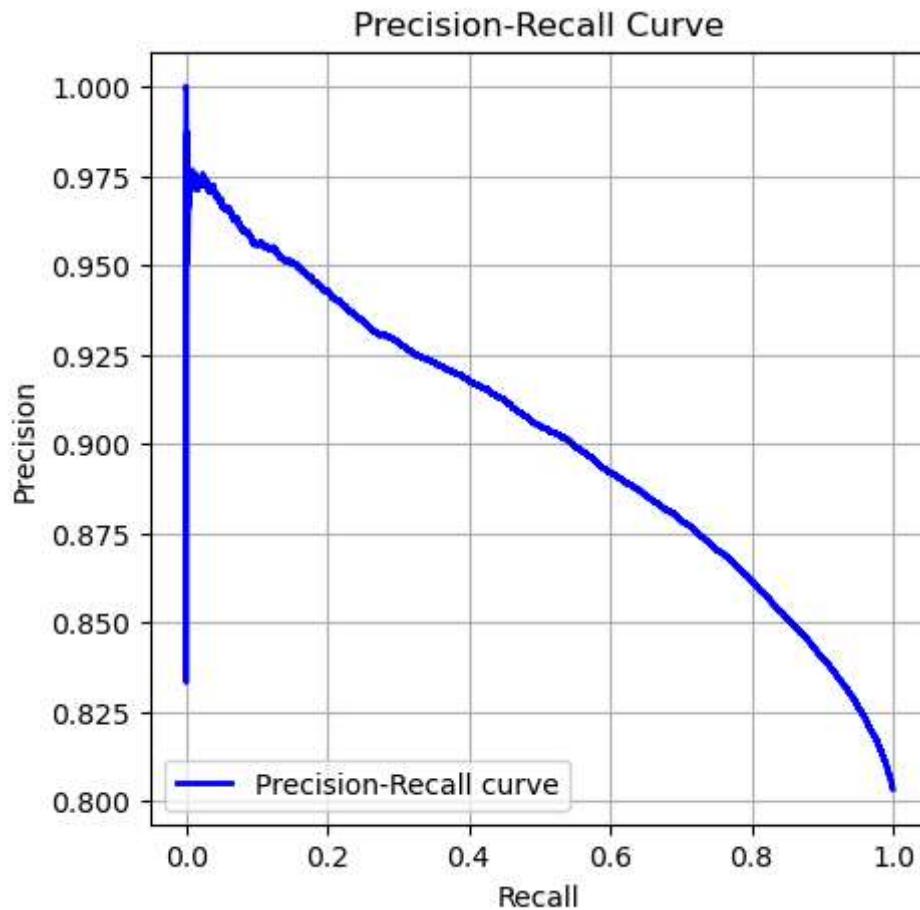
Comment on the ROC Curve :-

→ A model with an AUC of 0.65-0.7 is generally considered to have moderate performance. If your AUC is closer to 1, it would indicate a better ability of the model to distinguish between the two classes.

- The ROC curve should generally bow upwards, indicating a model that performs better than random chance. If the curve is near the diagonal line (red dashed line), the model has poor discriminative ability.
- As we can see our AUC = 0.71 and the ROC curve is also bow upwards which means our model has moderate performance not good .

2. Precision Recall Curve & comments

```
In [79]: ┌─ 1 from sklearn.metrics import precision_recall_curve, average_precision
  2 import matplotlib.pyplot as plt
  3
  4 X_test_selected = rfe.transform(X_test_scaled)
  5
  6 y_prob = model.predict_proba(X_test_selected)[:, 1]
  7
  8 precision, recall, thresholds = precision_recall_curve(y_test, y_prob)
  9
 10 plt.figure(figsize=(5, 5))
 11 plt.plot(recall, precision, color='blue', lw=2, label='Precision-Recall curve')
 12 plt.xlabel('Recall')
 13 plt.ylabel('Precision')
 14 plt.title('Precision-Recall Curve')
 15 plt.legend(loc='lower left')
 16 plt.grid(True)
 17 plt.show()
 18
 19 pr_auc = average_precision_score(y_test, y_prob)
 20 print(f'Average Precision (PR AUC): {pr_auc:.2f}')
```



Average Precision (PR AUC): 0.90

Comments on the Results :-

→ Precision-Recall Curve :

Based on your model's results, if the curve is near the top-right corner, it means your model is able to achieve both high precision and high recall. This is good for imbalanced datasets where the positive class (class 1) is more important.

→ PR AUC :

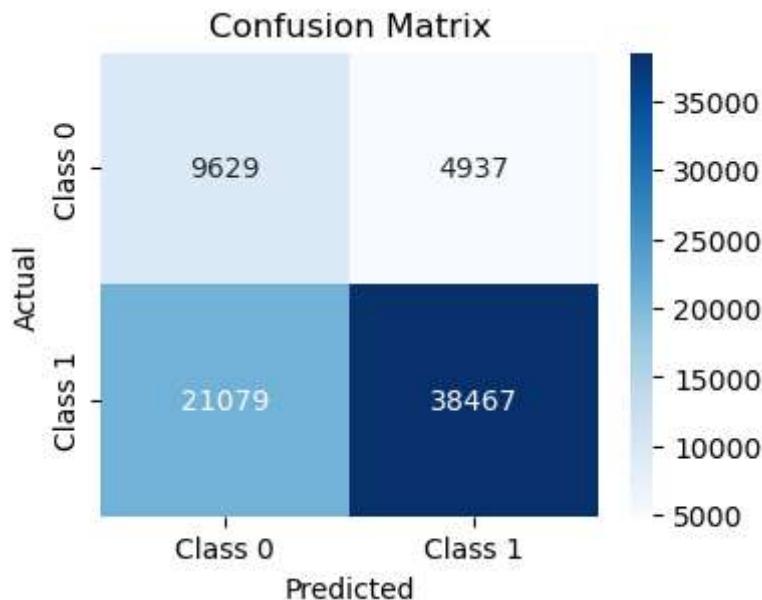
A PR AUC score of 0.7 or above is considered good. The closer the score is to 1, the better the model is at distinguishing between positive and negative classes with respect to both precision and recall. In our case (PR AUC) is 0.90 which means good.

3. Classification Report (Confusion Matrix etc)

```
In [80]: 1 from sklearn.metrics import classification_report, confusion_matrix
2
3 report = classification_report(y_test, y_pred)
4 print("Classification Report:\n", report)
5
6 cm = confusion_matrix(y_test, y_pred)
7
8 plt.figure(figsize=(4, 3))
9 sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", xticklabels=["Class 0", "Class 1"], yticklabels=["Actual 0", "Actual 1"])
10 plt.title('Confusion Matrix')
11 plt.xlabel('Predicted')
12 plt.ylabel('Actual')
13 plt.show()
```

Classification Report:

	precision	recall	f1-score	support
0	0.31	0.66	0.43	14566
1	0.89	0.65	0.75	59546
accuracy			0.65	74112
macro avg	0.60	0.65	0.59	74112
weighted avg	0.77	0.65	0.68	74112



4. Tradeoff Questions

- 1 a) How can we make sure that our model can detect real defaulters and there are less false positives? This is important as we can lose out on an opportunity to finance more individuals and earn interest on it,

This can be achieved through :-

- Adjusting decision threshold : Lowering the threshold for classifying a defaulter may increase recall but may decrease precision. You can use techniques like ROC curves or Precision-Recall curves to find the best threshold.
- Optimizing the F1-score : The current F1-score for class 1 is 0.75, which is decent but can be improved by increasing recall.
- Additional resampling techniques like undersampling the majority class or applying SMOTE-EN (SMOTE with edited nearest neighbors) could further help reduce false positives by making the model better at distinguishing between the classes.
- The current model may benefit from hyperparameter tuning using cross-validation techniques like GridSearchCV or RandomizedSearchCV. Consider using algorithms like Random Forest or XGBoost, which are better suited for handling class imbalances and may offer improved recall.
- Incorporating Domain Knowledge. Feature engineering and selection can help refine the model to focus on the most predictive features further reducing false positives

1 b) Since NPA (non-performing asset) is a real problem in this industry, it's important we play safe and shouldn't disburse loans to anyone

- Increase Recall for Defaulters (Class 1): Ensure more real defaulters are detected to minimize NPAs.
- Adjust Decision Threshold: Lower the threshold to classify defaulters more conservatively.
- Focus on F1-Score Optimization: Balance precision and recall to ensure safe decisions.
- Use Resampling Techniques: Improve model balance with advanced resampling like SMOTE-EN.
- Explore Robust Models: Use algorithms like Random Forest or XGBoost for better NPA detection.
- Incorporate Domain Knowledge: Prioritize features like debt-to-income ratio and credit history.
- Regular Model Monitoring: Continuously update the model with new NPA-related data.

=====



Actionable Insights :-

1. Imbalanced Data Handling: The model's recall for defaulters can be improved by using advanced techniques like SMOTE or SMOTE-ENN for better class balance.
2. Threshold Adjustment: Lowering the classification threshold can reduce the risk of misclassifying defaulters, ensuring fewer NPAs.

3. Key Feature Focus: Features like dti, loan_amount, and credit history should be monitored and weighted more heavily in decision-making.
 4. Improved Data Collection: Gather more granular data about borrower behavior, repayment history, and external economic factors.
 5. Use Ensemble Models: Algorithms like Random Forest and Gradient Boosting can capture complex patterns and improve prediction.
-
- 

Recommendations :-

1. Policy Adaptation: Implement stricter loan approval policies for high-risk profiles as flagged by the model.
 2. Regular Model Training: Continuously train the model with new data to adapt to changing borrower behavior.
 3. Business Rules Integration: Supplement the model with domain-specific rules (e.g., exclude borrowers with a high debt-to-income ratio).
 4. Pilot Testing: Run the updated model on historical data to validate its effectiveness in identifying defaulters.
 5. Performance Monitoring: Track key metrics like recall and F1-score for defaulters over time and adjust as needed.
-
- 