

Logistic Regression

Analysed by: Pavan Eleti

Introduction:

Loantap is a leading financial technology company based in India, specializing in providing flexible and innovative loan products to individuals and businesses. With a focus on customer-centric solutions, Loantap leverages technology to offer hassle-free borrowing experiences, including personal loans, salary advances, and flexible EMI options. Their commitment to transparency, speed, and convenience has established them as a trusted partner for borrowers seeking efficient financial solutions.

- LoanTap is at the forefront of offering tailored financial solutions to millennials.
- · Their innovative approach seeks to harness data science for refining their credit underwriting process.
- The focus here is the Personal Loan segment. A deep dive into the dataset can reveal patterns in borrower behavior and creditworthiness.
- Analyzing this dataset can provide crucial insights into the financial behaviors, spending habits, and potential risk associated with each borrower
- The insights gained can optimize loan disbursal, balancing customer outreach with risk management.

Our Task:

 As a data scientist at LoanTap, you are tasked with analyzing the dataset to determine the creditworthiness of potential borrowers. Your ultimate objective is to build a logistic regression model, evaluate its performance, and provide actionable insights for the underwriting process.

Features of the dataset:

· Column Profiling:

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
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
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

Exploratory Data Analysis

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import ttest_ind,chi2_contingency
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split, KFold, cross_val_score
from \ sklearn.preprocessing \ import \ MinMaxScaler, \ Label Encoder, \ Standard Scaler
from sklearn.metrics import (
    accuracy_score, confusion_matrix, classification_report,
    roc_auc_score, roc_curve, auc, precision_recall_curve, average_precision_score,
    ConfusionMatrixDisplay, RocCurveDisplay,f1_score,recall_score,precision_score
from statsmodels.stats.outliers_influence import variance_inflation_factor
from imblearn.over_sampling import SMOTE
import warnings
warnings.filterwarnings("ignore")
lt_data =pd.read_csv('loantap-data.csv')
df = lt_data.copy()
df.head()
```

0 10000.0 36 months 11.44 329.48 B B4 Marketing 10+ years RENT 117000.0 1 8000.0 36 months 11.99 265.68 B B5 Credit analyst 4 years MORTGAGE 65000.0 2 15600.0 36 months 10.49 506.97 B B3 Statistician <1 year		loan_amnt	term	int_rate	installment	grade	sub_grade	emp_title	emp_length	home_ownership	annual_inc	 open
1 8000.0 months 11.99 265.68 B BS analyst 4 years MORTGAGE 65000.0 months 2 15600.0 months 10.49 506.97 B B3 Statistician <1 year	(10000.0		11.44	329.48	В	В4	Marketing	10+ years	RENT	117000.0	
2 15600.0 months 10.49 506.97 B B3 Statistician <1 year RENT 43057.0 3 7200.0 36 6.49 220.65 A A2 Client Advocate 6 years RENT 54000.0 4 24375.0 60 months 17.27 609.33 C C5 Management Inc.	1	8000.0		11.99	265.68	В	B5		4 years	MORTGAGE	65000.0	
3 7200.0 months 6.49 220.65 A A2 Advocate 6 years HENT 54000.0 Destiny 4 24375.0 60 17.27 609.33 C C5 Management 9 years MORTGAGE 55000.0 Inc.	2	2 15600.0		10.49	506.97	В	В3	Statistician	< 1 year	RENT	43057.0	
4 24375.0 60 17.27 609.33 C C5 Management 9 years MORTGAGE 55000.0 Inc.	3	3 7200.0		6.49	220.65	Α	A2		6 years	RENT	54000.0	
5 rows × 27 columns	4	1 24375.0		17.27	609.33	С	C5	Management	9 years	MORTGAGE	55000.0	
	5	rows × 27 colun	nns									

pd.set_option('display.max_columns', None)

Exploration of data :

df.shape

→ (396030, 27)

df.info()

<<class 'pandas.core.frame.DataFrame'>
 RangeIndex: 396030 entries, 0 to 396029
 Data columns (total 27 columns):

Data	columns (total 27 c	olumns):	
#	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
                         396030 non-null object
13 purpose
                         394274 non-null object
14 title
                         396030 non-null
                                          float64
 15 dti
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
                         395754 non-null
                                          float64
 20
    revol_util
 21
    total_acc
                         396030 non-null
                                          float64
    initial_list_status
                         396030 non-null
 22
                                          object
   application_type
23
                         396030 non-null
                                          object
                         358235 non-null float64
 24
    mort_acc
25 pub_rec_bankruptcies
                         395495 non-null float64
26 address
                         396030 non-null object
dtypes: float64(12), object(15)
memory usage: 81.6+ MB
```

df.columns

Statistical Summary

df.describe().T

→		count	mean	std	min	25%	50%	75%	max
	loan_amnt	396030.0	14113.888089	8357.441341	500.00	8000.00	12000.00	20000.00	40000.00
	int_rate	396030.0	13.639400	4.472157	5.32	10.49	13.33	16.49	30.99
	installment	396030.0	431.849698	250.727790	16.08	250.33	375.43	567.30	1533.81
	annual_inc	396030.0	74203.175798	61637.621158	0.00	45000.00	64000.00	90000.00	8706582.00
	dti	396030.0	17.379514	18.019092	0.00	11.28	16.91	22.98	9999.00
	open_acc	396030.0	11.311153	5.137649	0.00	8.00	10.00	14.00	90.00
	pub_rec	396030.0	0.178191	0.530671	0.00	0.00	0.00	0.00	86.00
	revol_bal	396030.0	15844.539853	20591.836109	0.00	6025.00	11181.00	19620.00	1743266.00
	revol_util	395754.0	53.791749	24.452193	0.00	35.80	54.80	72.90	892.30
	total_acc	396030.0	25.414744	11.886991	2.00	17.00	24.00	32.00	151.00
	mort_acc	358235.0	1.813991	2.147930	0.00	0.00	1.00	3.00	34.00
	pub_rec_bankruptcies	395495.0	0.121648	0.356174	0.00	0.00	0.00	0.00	8.00

df.describe(include='object').T

₹		count	unique	top	freq
	term	396030	2	36 months	302005
	grade	396030	7	В	116018
	sub_grade	396030	35	В3	26655
	emp_title	373103	173105	Teacher	4389
	emp_length	377729	11	10+ years	126041
	home_ownership	396030	6	MORTGAGE	198348
	verification_status	396030	3	Verified	139563
	issue_d	396030	115	Oct-2014	14846
	loan_status	396030	2	Fully Paid	318357
	purpose	396030	14	debt_consolidation	234507
	title	394274	48816	Debt consolidation	152472
	earliest_cr_line	396030	684	Oct-2000	3017
	initial_list_status	396030	2	f	238066
	application_type	396030	3	INDIVIDUAL	395319

396030 393700 USCGC Smith\r\nFPO AE 70466

▼ ■■ Duplicate Detection

address

df[df.duplicated()]

loan_amnt term int_rate installment grade sub_grade emp_title emp_length home_ownership annual_inc verification

✓ ■ Insights

• The dataset does not contain any duplicates.

? Null Detection

df.isna().any()[df.isna().any()]

		_
\rightarrow	emp_title	True
_	emp_length	True
	title	True
	revol_util	True
	mort_acc	True
	<pre>pub_rec_bankruptcies</pre>	True
	dtype: bool	

df.isna().sum().sort_values(ascending=False)

\overline{z}	mort_acc	37795
_	emp_title	22927
	emp_length	18301
	title	1756
	<pre>pub_rec_bankruptcies</pre>	535
	revol_util	276
	loan_amnt	0
	dti	0
	application_type	0
	initial_list_status	0
	total_acc	0
	revol_bal	0
	pub_rec	0
	open_acc	0
	earliest_cr_line	0
	purpose	0
	term	0
	loan_status	0
	issue_d	0
	verification_status	0
	annual_inc	0
	home_ownership	0
	sub_grade	0
	grade	0
	installment	0
	int_rate	0

address dtype: int64

```
def missing_data(df):
    total_missing_df = df.isnull().sum().sort_values(ascending =False)
    percent_missing_df = (df.isnull().sum()/df.isna().count()*100).sort_values(ascending=False)
    missing_data_df = pd.concat([total_missing_df, percent_missing_df], axis=1, keys=['Total', 'Percent'])
    return missing_data_df
```

missing_pct = missing_data(df)
missing_pct[missing_pct['Total']>0]

₹		Total	Percent
	mort_acc	37795	9.543469
	emp_title	22927	5.789208
	emp_length	18301	4.621115
	title	1756	0.443401
	pub_rec_bankruptcies	535	0.135091
	revol_util	276	0.069692



Following columns has missing values

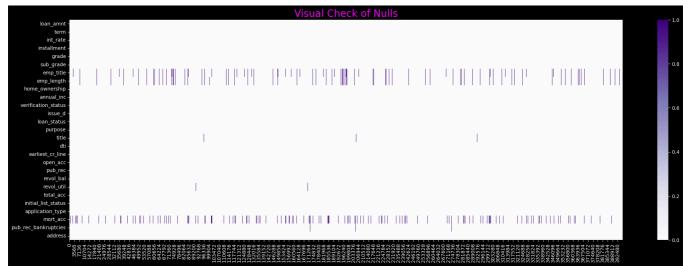
- 1. emp_title has 5.78% missing values
- 2. emp_length has 4.62% missing values
- 3. title has 0.44% missing values
- 4. revol_until has 0.06% missing values
- 5. mort_acc has 9.54% missing values
- 6. pub_rec_bankruptcies has 0.13% missing values

Action

• Since ML algorithm do not work on columns which has missing values so we need to impute these missing values.

```
plt.figure(figsize=(25,8))
plt.style.use('dark_background')
sns.heatmap(df.isnull().T,cmap='Purples')
plt.title('Visual Check of Nulls',fontsize=20,color='magenta')
plt.show()
```





```
df.isna().sum().sum()
# since there are 81590 rows are null , we cant drop na ...

    81590

#checking the unique values for columns
for _ in df.columns:
    print()
    print(f'Total Unique Values in {_} column are :- {df[_].nunique()}')
    print(f'Unique Values in {_} column are :-\n {df[_].unique()}')
    print(f'Value_counts of {_} column :-\n {df[_].value_counts()}')
    print()
    print('-'*120)
```

```
Total Unique Values in address column are :- 393700
Unique Values in address column are :-
 ['0174 Michelle Gateway\r\nMendozaberg, OK 22690'
 '1076 Carney Fort Apt. 347\r\nLoganmouth, SD 05113' 87025 Mark Dale Apt. 269\r\nNew Sabrina, WV 05113'
 '953 Matthew Points Suite 414\r\nReedfort, NY 70466'
 '7843 Blake Freeway Apt. 229\r\nNew Michael, FL 29597'
 '787 Michelle Causeway\r\nBriannaton, AR 48052']
Value_counts of address column :-
 address
USCGC Smith\r\nFPO AE 70466
                                                          8
USS Johnson\r\nFPO AE 48052
                                                          8
USNS Johnson\r\nFP0 AE 05113
                                                          8
USS Smith\r\nFPO AP 70466
                                                          8
USNS Johnson\r\nFP0 AP 48052
                                                          7
455 Tricia Cove\r\nAustinbury, FL 00813
                                                          1
7776 Flores Fall\r\nFernandezshire, UT 05113
                                                          1
6577 Mia Harbors Apt. 171\r\nRobertshire, OK 22690
8141 Cox Greens Suite 186\r\nMadisonstad, VT 05113
                                                          1
787 Michelle Causeway\r\nBriannaton, AR 48052
Name: count, Length: 393700, dtype: int64
```

✓ ■Null Treatment:

df.describe().T

```
df.loc[df['revol_util'].isna(),'revol_util'] = 0.0
df.loc[df['mort_acc'].isna(),'mort_acc'] = 0.0
df.loc[df['pub_rec_bankruptcies'].isna(),'pub_rec_bankruptcies'] = 0.0
df.loc[df['emp_title'].isna(),'emp_title'] = 'No Employee Title'
df.loc[df['title'].isna(),'title'] = 'Unavailable'
df['emp_length'] = df['emp_length'].fillna('< 1 year')</pre>
df.isna().sum()
                             0
→ loan_amnt
                             0
     term
     int_rate
     installment
                             0
     grade
     sub_grade
     emp_title
                             0
     emp_length
     home_ownership
    annual_inc
     verification_status
     issue_d
     loan_status
    purpose
                             0
     title
                             0
    dti
                             0
     earliest_cr_line
     open_acc
     pub_rec
     revol_bal
                             0
     revol_util
     total acc
                             0
     initial_list_status
     application_type
                             0
    mort acc
                             0
     pub_rec_bankruptcies
                             0
     address
                             0
     dtype: int64
```

		count	mean	std	min	25%	50%	75%	max
	loan_amnt	396030.0	14113.888089	8357.441341	500.00	8000.00	12000.00	20000.00	40000.00
	int_rate	396030.0	13.639400	4.472157	5.32	10.49	13.33	16.49	30.99
	installment	396030.0	431.849698	250.727790	16.08	250.33	375.43	567.30	1533.81
	annual_inc	396030.0	74203.175798	61637.621158	0.00	45000.00	64000.00	90000.00	8706582.00
	dti	396030.0	17.379514	18.019092	0.00	11.28	16.91	22.98	9999.00
	open_acc	396030.0	11.311153	5.137649	0.00	8.00	10.00	14.00	90.00
	pub_rec	396030.0	0.178191	0.530671	0.00	0.00	0.00	0.00	86.00
	revol_bal	396030.0	15844.539853	20591.836109	0.00	6025.00	11181.00	19620.00	1743266.00
	revol_util	396030.0	53.754260	24.484857	0.00	35.80	54.80	72.90	892.30
	total_acc	396030.0	25.414744	11.886991	2.00	17.00	24.00	32.00	151.00
	mort_acc	396030.0	1.640873	2.111249	0.00	0.00	1.00	3.00	34.00
	pub_rec_bankruptcies	396030.0	0.121483	0.355962	0.00	0.00	0.00	0.00	8.00

df.describe(include='object').T

<u> </u>					
		count	unique	top	freq
	term	396030	2	36 months	302005
	grade	396030	7	В	116018
	sub_grade	396030	35	B3	26655
	emp_title	396030	173106	No Employee Title	22927
	emp_length	396030	11	10+ years	126041
	home_ownership	396030	6	MORTGAGE	198348
	verification_status	396030	3	Verified	139563
	issue_d	396030	115	Oct-2014	14846
	loan_status	396030	2	Fully Paid	318357
	purpose	396030	14	debt_consolidation	234507
	title	396030	48817	Debt consolidation	152472
	earliest_cr_line	396030	684	Oct-2000	3017
	initial_list_status	396030	2	f	238066
	application_type	396030	3	INDIVIDUAL	395319
	address	396030	393700	USCGC Smith\r\nFPO AE 70466	8

→ Feature Engineering

```
df['pub_rec'] = [1 if i > 1 else 0 for i in df['pub_rec']]
df['mort_acc'] = [1 if i > 1 else 0 for i in df['mort_acc']]
df['pub_rec_bankruptcies'] = [1 if i > 1 else 0 for i in df['pub_rec_bankruptcies']]
```

df.sample()

df['address']

```
₹
                        term int_rate installment grade sub_grade
                                                                         emp_title emp_length home_ownership annual_inc veri
           loan_amnt
                                                                       Superintendent
                          36
     60136
                                                                   C1
                                                                                                                     95000.0
              35000.0
                                  12 29
                                             1167.36
                                                                                       10+ years
                                                                                                     MORTGAGE
                                                                                of
                                                                         Maintenance
```

```
#Split issue_date into month and year
df[['issue_month', 'issue_year']] = df['issue_d'].str.split('-', expand=True)
df.drop(['issue_d'], axis=1, inplace=True)

#Split er_cr_line date into month and year
df[['er_cr_line_m', 'er_cr_line_y']] = df['earliest_cr_line'].str.split('-', expand=True)
df.drop(['earliest_cr_line'], axis=1, inplace=True)
```

```
<del>_</del>
    0
                  0174 Michelle Gateway\r\nMendozaberg, OK 22690
     1
               1076 Carney Fort Apt. 347\r\nLoganmouth, SD 05113
     2
               87025 Mark Dale Apt. 269\r\nNew Sabrina, WV 05113
     3
                         823 Reid Ford\r\nDelacruzside, MA 00813
     4
                          679 Luna Roads\r\nGreggshire, VA 11650
     396025
                12951 Williams Crossing\r\nJohnnyville, DC 30723
     396026
               0114 Fowler Field Suite 028\r\nRachelborough, ...
               953 Matthew Points Suite 414\r\nReedfort, NY 7...
     396027
               7843 Blake Freeway Apt. 229\r\nNew Michael, FL..
     396028
     396029
                   787 Michelle Causeway\r\nBriannaton, AR 48052
    Name: address, Length: 396030, dtype: object
#Split address into State and Zip code
import re
df[['state','zipcode']] = df['address'].str.extract(r'([A-Z]{2}) (\d{5})')
df.drop(['address'], axis=1, inplace=True)
df['state'].nunique() , df['zipcode'].nunique()
→ (54, 10)
df['state'].isna().sum() , df['zipcode'].isna().sum()
→ (0, 0)
df['emp\_length\_yrs'] = df['emp\_length'].str.extract('(\d+)')
df.drop(['emp_length'], axis=1, inplace=True)
df['term'] = df['term'].str.split().str[0].astype('object')
df.sample()
\overline{\mathbf{T}}
            loan_amnt term int_rate installment grade sub_grade emp_title home_ownership annual_inc verification_status
                                                                        Insurance
     43629
               16425.0
                                 17.57
                                              590.27
                                                                                            RENT
                                                                                                       42000.0
                                                                                                                       Source Verified
                                                                        Consultant
```

df.shape

→ (396030, 30)

List of categorical columns
cat_cols = df.select_dtypes(include='object')

List of numerical columns
num_cols = df.select_dtypes(exclude='object')

cat_cols.sample(3)

₹		term	grade	sub_grade	emp_title	home_ownership	verification_status	loan_status	purpose	title	j
	297504	60	Е	E3	Admin Supervisor	RENT	Verified	Charged Off	debt_consolidation	Debt consolidation	
	348683	60	С	С3	BUSINESS CONSULTANT	MORTGAGE	Verified	Charged Off	credit_card	Credit card refinancing	
	361643	60	G	G2	Supply Chain Manager	RENT	Source Verified	Fully Paid	other	Other	

num_cols.sample(3)

 $\overline{\rightarrow}$ loan_amnt int_rate installment annual_inc dti open_acc pub_rec revol_bal revol_util total_acc mort_acc p 225076 14400.0 14.09 492.79 50000.0 34.70 13.0 0 5148.0 28.0 38.0 0 0 98.2 0 313487 16000.0 16.55 393.79 70000.0 27.93 6.0 14527.0 21.0 . . .

num_cols.skew()

loan_amnt 0.777285 int_rate 0.420669

```
0.983598
installment
                          41.042725
annual_inc
                         431.051225
dti
open_acc
                           1.213019
pub_rec
                           6.812303
revol_bal
                          11.727515
revol_util
                          -0.074238
total_acc
                           0.864328
                           0.412225
mort acc
pub_rec_bankruptcies
                          12.936099
dtype: float64
```

- ✓

 √ Insights
 - · Features are Right skewed

Action

· Need to apply log transformations in order to normalise them

```
df1.sample()

loan_amnt term int_rate installment grade sub_grade emp_title home_ownership annual_inc verification_statu

241830 8000.0 36 9.99 258.1 B B3 Employee RENT 36000.0 Verifie
```

Q1. What percentage of customers have fully paid their Loan Amount?

df['loan_status'].value_counts(normalize=True)*100

Fully Paid 80.387092
Charged Off 19.612908
Name: proportion, dtype: float64

- ✓

 √ Insights:
 - Target variable distribution is 80%-20%. Data is significantly imbalanced

Graphical Analysis:

```
uni / bi / multi variate Analysis
```

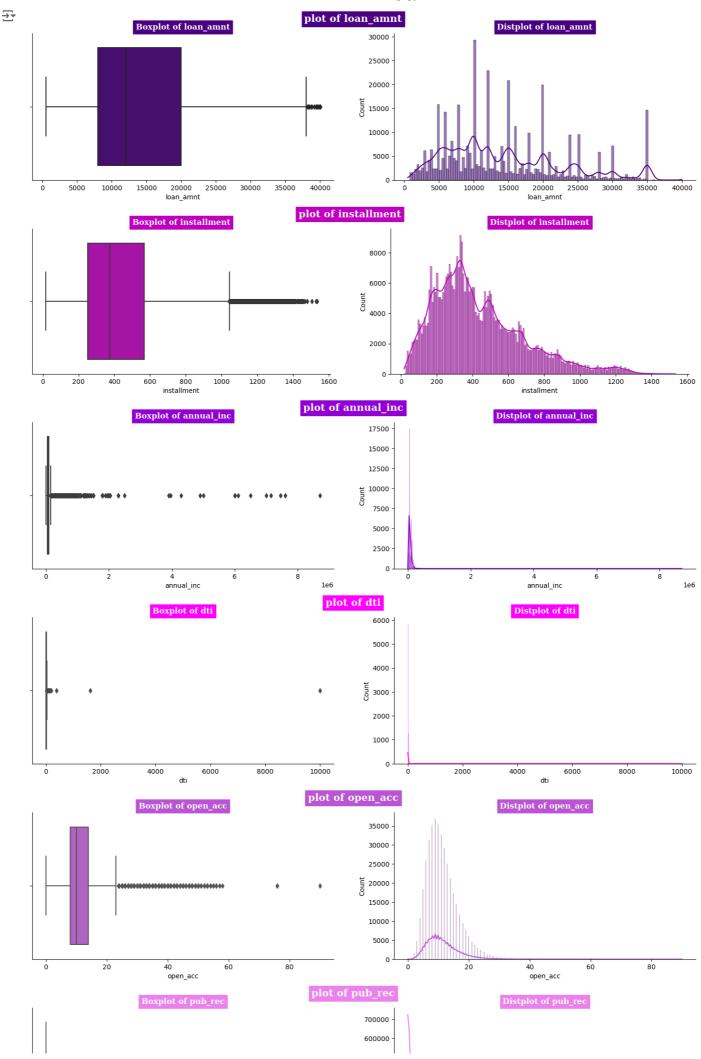
```
cp = ['indigo','m','darkviolet','magenta','mediumorchid','violet','purple','orchid','mediumpurple','deeppink','blueviolet','
```

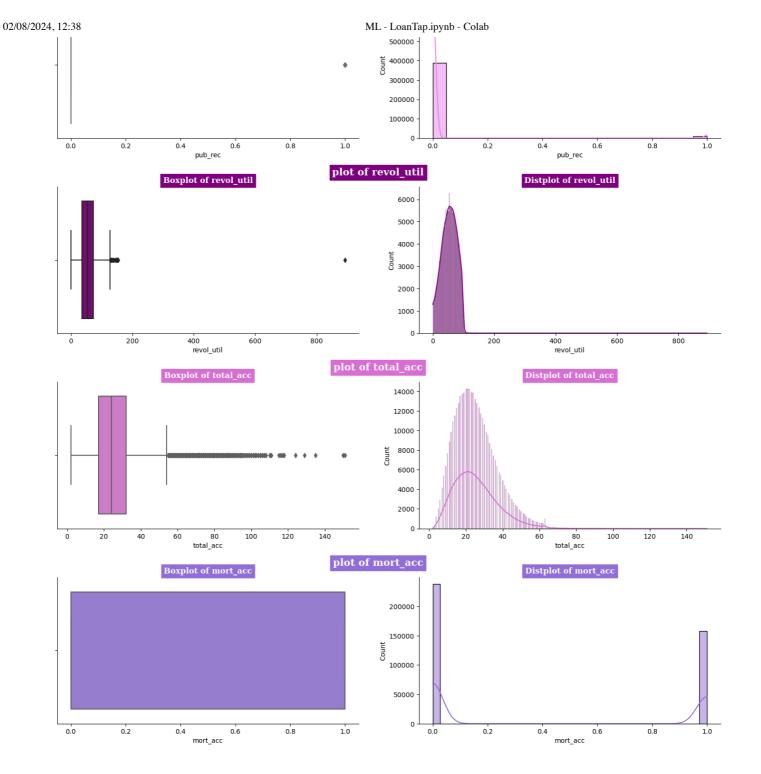
num_cols.iloc[:,[0,2,3,4,5,6,8,9,10]].sample()

```
        Source
        loan_amnt
        installment
        annual_inc
        dti
        open_acc
        pub_rec
        revol_util
        total_acc
        mort_acc

        52024
        10000.0
        317.54
        105000.0
        25.98
        16.0
        0
        55.9
        43.0
        1
```

```
plt.style.use('default')
plt.style.use('seaborn-bright')
outlier_graphical_cols = num_cols.iloc[:,[0,2,3,4,5,6,8,9,10]]
for _,col in enumerate(outlier_graphical_cols.columns):
    plt.figure(figsize=(18,4))
    plt.suptitle(f'plot of {col}',fontsize=15,fontfamily='serif',fontweight='bold',backgroundcolor=cp[_],color='w')
    plt.subplot(121)
    sns.boxplot(x=df[col],color=cp[_])
    plt.title(f'Boxplot of {col}',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[_],color='w')
    plt.subplot(122)
    sns.histplot(x=df[col], kde=True,color=cp[_])
    plt.title(f'Distplot of {col}',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[_],color='w')
    sns.despine()
    plt.show()
```





✓ √ Insights:

- 1. The analysis suggests a prevalence of outliers, prompting further investigation into outlier detection techniques.
- 2. Among the numerical features, Potential outliers may still be present.
- 3. Notably, features such as Pub_rec, Mort_acc, and Pub_rec_bankruptcies display a sparse distribution of unique values, indicating the potential benefit of generating binary features from these variables.

```
#Countplots of various categorical features w.r.t. to target variable loan_status
plt.figure(figsize=(16,17))
plt.suptitle('Countplots of various categorical features w.r.t. to target variable loan_status',
             fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor=cp[1], color='w')
plt.subplot(321)
sns.countplot(data=df, x='loan_status',palette=cp)
plt.title('Loan Status Counts',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[2],color='w')
plt.subplot(322)
sns.countplot(data=df, x='loan_status', hue='term',palette=cp)
plt.title('Term wise loan status count',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[3],color='w')
sns.countplot(data=df, x='home_ownership', hue='loan_status',palette=cp)
plt.title('Loan Status Vs Home Ownership',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[4],color='w')
plt.subplot(324)
\verb|sns.countplot(data=df, x='verification_status', hue='loan_status', palette=cp)|
plt.title('Loan Status Vs Verification Status',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[5],color=
plt.subplot(325)
\verb|sns.countplot(data=df, x='issue\_month', hue='loan\_status', palette=cp)|
plt.title('Loan Status Vs issue_month',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[6],color='w')
plt.subplot(326)
sns.countplot(data=df, x='zipcode', hue='loan_status',palette=cp)
plt.title('Loan Status Vs zipcode',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[7],color='w')
sns.despine()
plt.show()
```





```
zip_codes = ["11650", "86630", "93700"]
states = df[df['zipcode'].isin(zip_codes)]['state']
for zip_code, state in zip(zip_codes, states):
    print(f"Zip code: {zip_code}, State: {state}")

    Zip code: 11650, State: VA
    Zip code: 86630, State: MI
    Zip code: 93700, State: MD
```

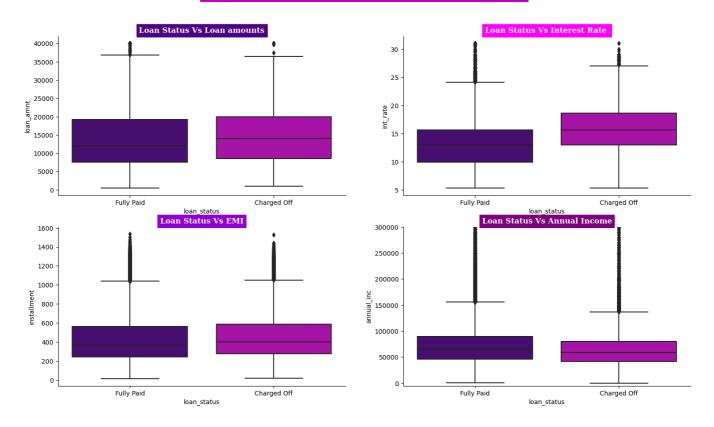
✓ Observations:

- It's been observed that loans haven't been completely repaid in zip codes 11650, 86630, and 93700.
- · Loans haven't been repaid by borrowers residing in 'VA', 'MI', and 'MD'.

```
#Boxplot of various cont. features w.r.t. target variable loan_status
plt.figure(figsize=(18,10))
plt.suptitle('Boxplot of various cont. features w.r.t. target variable loan_status',
             fontsize=14,fontfamily='serif',fontweight='bold',backgroundcolor=cp[1],color='w')
plt.subplot(221)
sns.boxplot(data=df, x='loan_status', y='loan_amnt',palette=cp)
plt.title('Loan Status Vs Loan amounts',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[0],color='w')
sns.boxplot(data=df, x='loan_status', y='int_rate',palette=cp)
plt.title('Loan Status Vs Interest Rate ',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[3],color='w')
plt.subplot(223)
sns.boxplot(data=df, x='loan_status', y='installment',palette=cp)
plt.title('Loan Status Vs EMI',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[2],color='w')
plt.subplot(224)
sns.boxplot(data=df, x='loan_status', y='annual_inc',palette=cp)
plt.ylim(bottom=-5000, top=300000)
plt.title('Loan Status Vs Annual Income',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[6],color='w')
sns.despine()
plt.show()
```



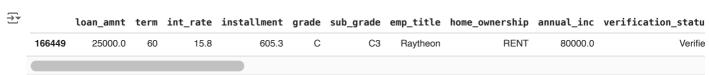
Boxplot of various cont. features w.r.t. target variable loan_status



Observations:

- Charged Off customers exhibit a notably higher median interest rate compared to Fully Paid customers.
- The median annual income of Charged Off customers is lower than that of Fully Paid customers.
- Charged Off customers tend to have a higher median EMI compared to Fully Paid customers.
- The median loan amount for Charged Off customers surpasses that of Fully Paid customers.

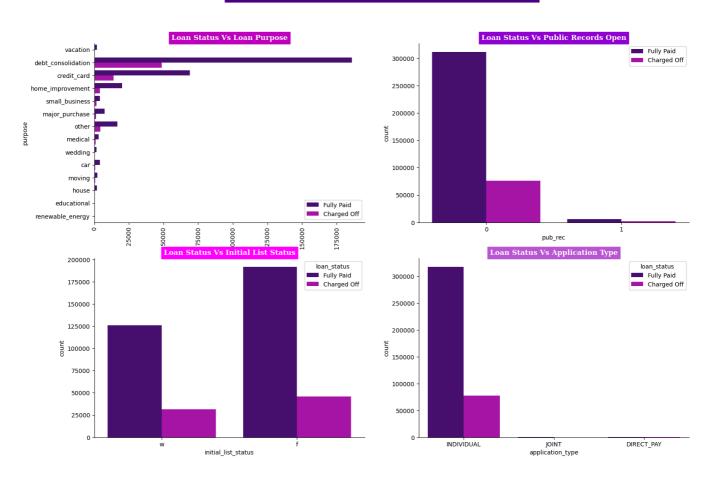
df.sample()



```
#Countplot of categorical variables w.r.t. target variable loan_status
plt.figure(figsize=(18,12))
plt.suptitle('Countplot of categorical variables w.r.t. target variable loan_status',
             fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor=cp[0], color='w')
plt.subplot(221)
sns.countplot(data=df, y='purpose', hue='loan_status',palette=cp)
plt.xticks(rotation=90)
plt.title('Loan Status Vs Loan Purpose',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[1],color='w')
plt.legend(loc=4)
plt.subplot(222)
sns.countplot(data=df, x='pub_rec',hue='loan_status',palette=cp)
plt.title('Loan Status Vs Public Records Open',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[2],color=
plt.legend(loc=1)
plt.subplot(223)
sns.countplot(data=df, x='initial_list_status', hue='loan_status',palette=cp)
plt.title('Loan Status Vs Initial List Status',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[3],color=
plt.subplot(224)
sns.countplot(data=df, x='application_type',hue='loan_status',palette=cp)
plt.title('Loan Status Vs Application Type',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[4],color='w'
sns.despine()
plt.show()
```

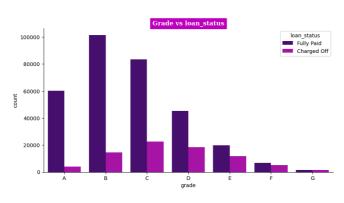
$\overline{\mathbf{T}}$

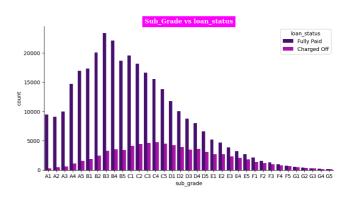
Countplot of categorical variables w.r.t. target variable loan_status





Countplot of categorical variables w.r.t. target variable loan_status





Observations:

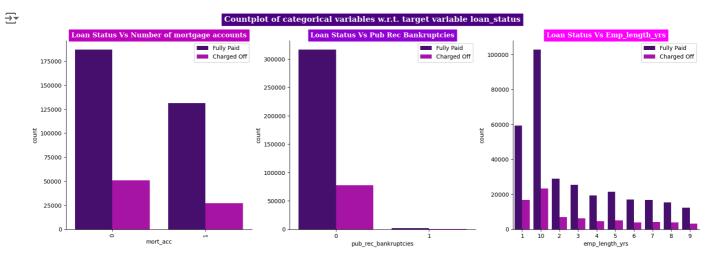
- Top 2 loan purpose categories are Debit Consolidation and Credit Card
- Topmost loan type application is INDIVIDUAL
- The distribution of open_acc appears to be relatively normal when visualized graphically.
- Charged Off and Fully Paid categories exhibit similar distributions.

df.sample()

₹		loan_amnt	term	int_rate	installment	grade	sub_grade	emp_title	home_ownership	annual_inc	verification_statu
	67006	28000.0	36	14.27	960.65	С	C2	ADP Total Souce (Specialty Manufacturing	MORTGAGE	125000.0	Verifi

```
#Countplot for various categorical features w.r.t. target variable loan_status
```

```
plt.figure(figsize=(20,6))
plt.suptitle('Countplot of categorical variables w.r.t. target variable loan_status',
             fontsize=14, fontfamily='serif', fontweight='bold', backgroundcolor=cp[0], color='w')
plt.subplot(131)
sns.countplot(data=df, x='mort_acc',hue='loan_status',palette=cp)
plt.xticks(rotation=90)
plt.title('Loan Status Vs Number of mortgage accounts',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[1
plt.legend(loc=1)
plt.subplot(132)
\verb|sns.countplot(data=df, x='pub\_rec\_bankruptcies', hue='loan\_status', palette=cp)|
plt.title('Loan Status Vs Pub Rec Bankruptcies',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[2],color
plt.legend(loc=1)
plt.subplot(133)
order = sorted(df.emp_length_yrs.unique().tolist())
sns.countplot(data=df, x='emp_length_yrs',hue='loan_status',order=order,palette=cp)
plt.title('Loan Status Vs Emp_length_yrs',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[3],color='w')
plt.legend(loc=1)
sns.despine()
plt.show()
```

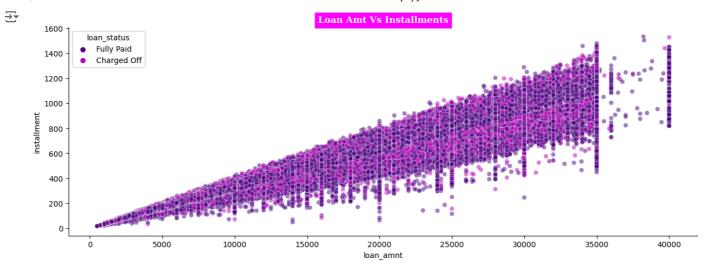


✓ Q2. Comment about the correlation between Loan Amount and Installment features.

df[['loan_amnt', 'installment']].corr()



```
plt.figure(figsize = (15,5))
sns.scatterplot(data = df, x = 'loan_amnt', y = 'installment', alpha = 0.5, hue = 'loan_status', palette = cp)
plt.title('Loan Amt Vs Installments',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[3],color='w')
sns.despine()
plt.show()
```



The correlation coefficient measures the strength and direction of the linear relationship between two variables. In this case, the correlation coefficient between 'loan_amnt' and 'installment' is quite high, approximately 0.95, indicating a strong positive linear relationship between these two variables.

- Loan Terms: Understanding the relationship between loan amount and installment payments is crucial for setting appropriate loan terms. Lenders can adjust loan terms such as interest rates and repayment periods based on the borrower's ability to handle installment payments associated with different loan amounts.
- Potential Multicollinearity: When building predictive models, it's essential to be cautious of multicollinearity between highly correlated
 predictor variables. Multicollinearity can lead to unstable estimates and difficulties in interpreting the model coefficients. Therefore, it
 might be necessary to address multicollinearity through techniques such as variable selection or regularization.

Q3. The majority of people have home ownership as ____.

 $(df['home_ownership'].value_counts(normalize=True)*100).to_frame()$

→ -	proportion
home_ownership	
MORTGAGE	50.084085
RENT	40.347953
OWN	9.531096
OTHER	0.028281
NONE	0.007828
ANY	0.000758

Insights:

- Mortgage holders comprise the majority with approximately 50.08%, indicating that a significant portion of individuals own homes through Mortgage agreements.
- Renters constitute a substantial portion, accounting for around 40.35% of home ownership types. This suggests a sizable demographic of individuals who opt for renting rather than owning a home.

pd.crosstab(df['grade'],df['loan_status'], normalize = 'index')

₹	loan_status	Charged Off	Fully Paid
	grade		
	Α	0.062879	0.937121
	В	0.125730	0.874270
	С	0.211809	0.788191
	D	0.288678	0.711322

0.373634

0.427880

0.478389

0.626366

0.572120

0.521611

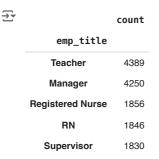
Ε

G

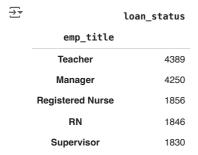
- True . Grade 'A' borrowers demonstrate a significantly high likelihood of fully repaying their loans, with approximately
 93.71% of loans being fully paid. This suggests that borrowers with the highest credit rating are more inclined to fulfill their loan obligations successfully.
- The proportion of charged-off loans for grade 'A' borrowers is relatively low, standing at approximately 6.29%. This indicates a low default rate among borrowers with the highest credit rating, emphasizing their creditworthiness and reliability in loan repayment.

Q5. Name the top 2 afforded job titles.

df[df['emp_title'] != 'No Employee Title']['emp_title'].value_counts().to_frame().head()



df.groupby('emp_title')['loan_status'].count().sort_values(ascending=False).to_frame()[1:6]



✓ √ Insights:

• The Most afforded job titles are Teachers & Managers.

```
plt.figure(figsize=(20,12))
sns.heatmap(num_cols.corr(), annot=True, cmap='Purples')
plt.title('Correlations - Heatmap',fontsize=12,fontfamily='serif',fontweight='bold',backgroundcolor=cp[2],color='w')
plt.show()
```



						Correlation	s - Heatma _l	9					1.0
loan_amnt -	1	0.17	0.95	0.34	0.017	0.2	-0.028	0.33	0.1	0.22	0.22	-0.028	
int_rate -	0.17	1	0.16	-0.057	0.079	0.012	0.037	-0.011	0.29	-0.036	-0.054	0.022	
installment -	0.95	0.16	1	0.33	0.016	0.19	-0.021	0.32	0.12	0.2	0.2	-0.025	- 0.8
annual_inc -	0.34	-0.057	0.33	1	-0.082	0.14	0.01	0.3	0.027	0.19	0.2	-0.0092	
dti -	0.017	0.079	0.016	-0.082	1	0.14	-0.013	0.064	0.088	0.1	0.0017	-0.0064	- 0.6
open_acc -	0.2	0.012	0.19	0.14	0.14	1	-0.0097	0.22	-0.13	0.68	0.13	-0.0076	
pub_rec -	-0.028	0.037	-0.021	0.01	-0.013	-0.0097	1	-0.045	-0.041	-0.0013	0.013	0.53	- 0.4
revol_bal -	0.33	-0.011	0.32	0.3	0.064	0.22	-0.045	1	0.23	0.19	0.18	-0.034	
revol_util -	0.1	0.29	0.12	0.027	0.088	-0.13	-0.041	0.23	1	-0.1	0.019	-0.035	- 0.2
total_acc -	0.22	-0.036	0.2	0.19	0.1	0.68	-0.0013	0.19	-0.1	1	0.33	0.016	
mort_acc -	0.22	-0.054	0.2	0.2	0.0017	0.13	0.013	0.18	0.019	0.33	1	0.017	- 0.0
pub_rec_bankruptcies -	-0.028	0.022	-0.025	-0.0092	-0.0064	-0.0076	0.53	-0.034	-0.035	0.016	0.017	1	
	loan_amnt -	int_rate -	installment -	annual_inc -	dti -	oben_acc -	- bnp_rec	revol_bal -	revol_util -	total_acc -	mort_acc -	pub_rec_bankruptcies -	

→ Observations:

- There exists a strong correlation between loan_amnt and installment, indicating that higher loan amounts correspond to larger installment payments.
- The variables total_acc and open_acc exhibit a significant correlation.
- There is a notable correlation between pub_rec_bankruptcies and pub_rec.

Qoutlier Treatment:

outlier treatment

```
def remove_outliers_zscore(df, threshold=2): #(considering 2 std.dev away from mean approx 95% of data)
    Remove outliers from a DataFrame using the Z-score method.
    Parameters:
        df (DataFrame): The input DataFrame.
        threshold (float): The Z-score threshold for identifying outliers.
                           Observations with a Z-score greater than this threshold
                           will be considered as outliers.
    Returns:
       DataFrame: The DataFrame with outliers removed.
    # Calculate Z-scores for numerical columns
    z_scores = (df[numerical_cols] - df[numerical_cols].mean()) / df[numerical_cols].std()
    # Identify outliers
    outliers = np.abs(z_scores) > threshold
    # Keep non-outliers for numerical columns
    df_cleaned = df[~outliers.any(axis=1)]
    return df_cleaned
cleaned_df = remove_outliers_zscore(df1)
print(cleaned_df.shape)

→ (311392, 30)
def clip_outliers_zscore(df, threshold=2):
    Clip outliers in a DataFrame using the Z-score method.
    Parameters:
        df (DataFrame): The input DataFrame.
        threshold (float): The Z-score threshold for identifying outliers.
                           Observations with a Z-score greater than this threshold
                           will be considered as outliers.
    Returns:
       DataFrame: The DataFrame with outliers clipped.
    # Calculate Z-scores for numerical columns
    z_scores = (df[numerical_cols] - df[numerical_cols].mean()) / df[numerical_cols].std()
    # Clip outliers
    {\tt clipped\_values = df[numerical\_cols].clip(df[numerical\_cols].mean() - threshold * df[numerical\_cols].std(),} \\
                                              df[numerical_cols].mean() + threshold * df[numerical_cols].std(),
    # Assign clipped values to original DataFrame
    df clipped = df.copv()
    df_clipped[numerical_cols] = clipped_values
    return df_clipped
clipped_df = clip_outliers_zscore(df1)
print(clipped_df.shape)

→ (396030, 30)

data = cleaned_df.copy()
cp_data = clipped_df.copy()
data.sample()
\overline{2}
             loan_amnt term int_rate installment grade sub_grade emp_title home_ownership annual_inc verification_statu
     110850
               14000.0
                          36
                                 11.67
                                              462.8
                                                        В
                                                                  B4
                                                                         Manager
                                                                                           RENT
                                                                                                     95000.0
                                                                                                                     Source Verifie
data['pub_rec_bankruptcies'].value_counts() , data['pub_rec'].value_counts()
    (pub_rec_bankruptcies
          311392
     Name: count, dtype: int64,
     pub_rec
          311392
     Name: count, dtype: int64)
```

cp_data['pub_rec_bankruptcies'].value_counts() , cp_data['pub_rec'].value_counts()

(pub_rec_bankruptcies 0.000000 393705 0.158662 2325 Name: count, dtype: int64, pub_rec 0.000000 388011 0.301947 8019 Name: count, dtype: int64)

data.shape

→ (311392, 30)

data.info()

<<class 'pandas.core.frame.DataFrame'>
 Index: 311392 entries, 0 to 396029
 Data columns (total 30 columns):

Data #	columns (total 30 column		.l Count	Dtype
0	loan amnt	311302	non-null	float64
1	term		non-null	object
	int rate	311392		float64
2	installment		non-null	float64
4	grade		non-null	object
5	sub grade		non-null	object
6	emp_title		non-null	object
7	home ownership		non-null	object
8	annual_inc	311392		float64
9	verification_status	311392	non-null	object
10	loan_status	311392	non-null	object
11	purpose	311392	non-null	object
12	title	311392	non-null	object
13	dti	311392	non-null	float64
14	open_acc		non-null	float64
15	pub_rec	311392	non-null	int64
16	revol_bal		non-null	float64
17	revol_util		non-null	float64
18	total_acc		non-null	float64
19	initial_list_status		non-null	object
20	application_type		non-null	object
21	mort_acc	311392		int64
22	<pre>pub_rec_bankruptcies</pre>	311392		int64
23	issue_month	311392		object
24	issue_year		non-null	object
25	er_cr_line_m er_cr_line_y		non-null	object
26			non-null	object
27	state		non-null	object
28	zipcode		non-null	object
29	emp_length_yrs		non-null	object
	es: float64(9), int64(s, obje	ect(18)	
memo	ry usage: 73.6+ MB			

Manual encoding:

data['loan_status']=data.loan_status.map({'Fully Paid':1, 'Charged Off':0})

data.head()

_ →		loan_amnt	term	int_rate	installment	grade	sub_grade	emp_title	home_ownership	annual_inc	verification_status	lı
	0	10000.0	36	11.44	329.48	В	B4	Marketing	RENT	117000.0	Not Verified	
	1	8000.0	36	11.99	265.68	В	B5	Credit analyst	MORTGAGE	65000.0	Not Verified	
	2	15600.0	36	10.49	506.97	В	ВЗ	Statistician	RENT	43057.0	Source Verified	
	3	7200.0	36	6.49	220.65	Α	A2	Client Advocate	RENT	54000.0	Not Verified	
	4	24375.0	60	17.27	609.33	С	C5	Destiny Management Inc.	MORTGAGE	55000.0	Verified	

[✓] Feature selection - done by hypothesis testing & VIF(multicolinearity)

Find VIF after modelling and remove features with high VIF (>5):

```
def calc vif(X):
    # Calculating the VIF
    vif=pd.DataFrame()
    vif['Feature']=X.columns
     vif['VIF']=[variance_inflation_factor(X.values,i) for i in range(X.shape[1])]
    vif['VIF']=round(vif['VIF'],2)
     vif=vif.sort_values(by='VIF',ascending=False)
     return vif
cat_cols = data.select_dtypes(include=['object']).columns.tolist()
for col in cat_cols:
    chi2, p, dof, expected = chi2_contingency(pd.crosstab(data[col], data['loan_status']))
    if p > 0.05:
        print('>>>>> Independent feature - Not Significant:',col,' >> p value:',p)
>>>>> Independent feature - Not Significant: emp_title >> p value: 0.5367121560200798
    >>>>> Independent feature - Not Significant: title >> p value: 1.0 >>>>> Independent feature - Not Significant: er_cr_line_m >> p value: 0.2722117086158036
    >>>>> Independent feature - Not Significant: state >> p value: 0.76047808977373
## dropping cols based on correlation(heatmap,hypothesis testing)
lt = data.drop(columns=['emp_title','title','sub_grade','er_cr_line_m','er_cr_line_y','initial_list_status',
                         'state','issue_month','issue_year','pub_rec','pub_rec_bankruptcies'],axis=1)
lt.shape
→ (311392, 19)
lt.sample()
\overline{2}
             loan_amnt term int_rate installment grade home_ownership annual_inc verification_status loan_status
                                                                 MODTOMOE
                                                                                 100000
#### Performing OneHotEncoding on feature having multiple variable
dummies=['zipcode', 'grade','purpose','home_ownership','verification_status','application_type']
ltd = pd.get_dummies(lt, columns=dummies, drop_first=True)*1
ltd.shape
→ (311392, 50)
ltd.dtypes
    loan_amnt
                                              float64
                                               object
     term
     int_rate
                                              float64
                                              float64
     installment
     annual inc
                                              float64
     loan_status
                                                int64
                                              float64
     dti
    open_acc
                                              float64
     revol_bal
                                              float64
     revol_util
                                              float64
                                              float64
     total_acc
                                                int64
     mort_acc
    emp_length_yrs
                                               object
    zipcode_05113
                                                int64
    zipcode_11650
                                                int64
     zipcode_22690
                                                int64
    zipcode_29597
                                                int64
     zipcode_30723
                                                int64
     zipcode_48052
                                                int64
     zipcode_70466
                                                int64
     zipcode_86630
                                                int64
     zipcode_93700
                                                int64
     grade_B
                                                int64
    grade_C
                                                int64
     grade_D
                                                int64
     grade_E
                                                int64
    grade_F
                                                int64
                                                int64
     grade_G
     purpose_credit_card
                                                int64
     purpose_debt_consolidation
                                                int64
```

int64

purpose_educational

purpose_home_improvement

```
int64
purpose_house
purpose_major_purchase
                                          int64
                                          int64
purpose_medical
purpose_moving
                                          int64
purpose_other
                                          int64
purpose_renewable_energy
                                          int64
purpose_small_business
                                          int64
purpose_vacation
                                          int64
purpose_wedding
                                          int64
home_ownership_MORTGAGE
                                          int64
home_ownership_OTHER
                                          int64
                                          int64
home_ownership_OWN
                                          int64
home_ownership_RENT
                                          int64
verification_status_Source Verified
                                          int64
verification_status_Verified
                                          int64
application_type_INDIVIDUAL
                                          int64
application_type_JOINT
                                          int64
dtype: object
```

ltd.sample(8)

_		$\overline{\mathbf{v}}$
	٠	

	loan_amnt	term	int_rate	installment	annual_inc	loan_status	dti	open_acc	revol_bal	revol_util	total_acc	n
118504	15000.0	36	10.99	491.01	85000.0	0	17.05	11.0	44706.0	88.0	17.0	
20036	26000.0	36	11.99	863.45	100000.0	1	13.22	13.0	19601.0	56.0	24.0	
388815	9175.0	36	13.35	310.70	40000.0	1	15.12	12.0	2447.0	40.1	27.0	
388094	13000.0	60	18.24	331.82	82000.0	1	18.29	12.0	20040.0	89.5	35.0	
254903	9600.0	36	15.31	334.25	65000.0	1	11.78	5.0	8346.0	76.6	11.0	
264585	22500.0	36	15.61	786.71	81000.0	1	18.42	21.0	21860.0	86.7	33.0	
368842	15000.0	36	8.90	476.30	120000.0	1	9.56	12.0	13292.0	41.9	32.0	
44417	24000.0	60	13.99	558.32	75000.0	1	16.26	8.0	12780.0	61.0	28.0	

✓ Model:

Minmax scaling the data

```
scaler = MinMaxScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
X_train = pd.DataFrame(X_train, columns=X.columns)
X_test = pd.DataFrame(X_test, columns=X.columns)
X_train.head()
```



	loan_amnt	term	int_rate	installment	annual_inc	dti	open_acc	revol_bal	revol_util	total_acc	mort_acc	emp_le
0	0.379538	0.0	0.339161	0.411590	0.207250	0.465341	0.368421	0.171897	0.419816	0.276596	0.0	
1	0.643564	1.0	0.680070	0.524221	0.367868	0.252652	0.473684	0.221905	0.590398	0.340426	0.0	
2	0.168317	0.0	0.208625	0.176198	0.134712	0.357576	0.368421	0.052236	0.304392	0.212766	0.0	
3	0.379538	1.0	0.680070	0.307444	0.367868	0.449242	0.315789	0.255109	0.767109	0.297872	1.0	
4	0.368812	0.0	0.543706	0.421460	0.246109	0.315530	0.263158	0.090649	0.614913	0.361702	0.0	

✓ ■Model-1

```
#Fit the Model on training data
logreg_model = LogisticRegression()
logreg_model.fit(X_train, y_train)

* LogisticRegression
    LogisticRegression()

#Predit the data on test dataset
y_train_pred = logreg_model.predict(X_train)
y_test_pred = logreg_model.predict(X_test)
```

logreg_model.score(X_test, y_test) , logreg_model.score(X_test, y_test_pred)

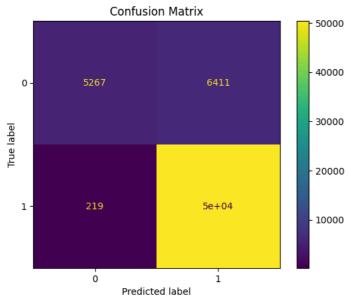
```
(0.8935435700637454, 1.0)
```

If logreg_model.score(X_test, y_test) consistently returns 1, it would imply that your model is predicting the test set perfectly, which could be a sign of overfitting, data leakage, or an issue with the evaluation process.

```
#Model Evaluation
print('Train Accuracy :', logreg_model.score(X_train, y_train).round(2))
print('Train F1 Score:',f1_score(y_train,y_train_pred).round(2))
print('Train Recall Score:',recall_score(y_train,y_train_pred).round(2))
print('Train Precision Score:',precision_score(y_train,y_train_pred).round(2))
print('\nTest Accuracy :',logreg_model.score(X_test,y_test).round(2))
print('Test F1 Score:',f1_score(y_test,y_test_pred).round(2))
print('Test Recall Score:',recall_score(y_test,y_test_pred).round(2))
print('Test Precision Score:',precision_score(y_test,y_test_pred).round(2))
# Confusion Matrix
cm = confusion_matrix(y_test, y_test_pred)
disp = ConfusionMatrixDisplay(cm)
disp.plot()
plt.title('Confusion Matrix')
plt.show()
```

Train Accuracy: 0.89
Train F1 Score: 0.94
Train Recall Score: 1.0
Train Precision Score: 0.89

Test Accuracy: 0.89 Test F1 Score: 0.94 Test Recall Score: 1.0 Test Precision Score: 0.89



print(classification_report(y_test,y_test_pred))

_	precision	recall	f1-score	support
0 1	0.96 0.89	0.45 1.00	0.61 0.94	11678 50601
accuracy macro avg weighted avg	0.92 0.90	0.72 0.89	0.89 0.78 0.88	62279 62279 62279

· Here the recall value for the 'charged off' is very low, Hence will build a better model

✓ ■ Model-2

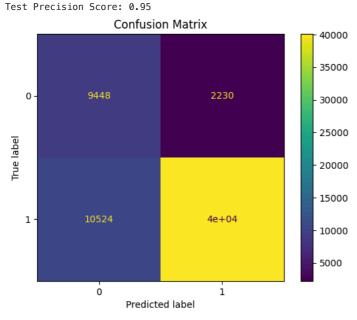
Oversampling to balance the target variable

```
sm=SMOTE(random_state=42)
X_train_res, y_train_res = sm.fit_resample(X_train,y_train.ravel())
print(f"Before OverSampling, count of label 1: {sum(y_train == 1)}")
print(f"Before OverSampling, count of label 0: {sum(y_train == 0)}")
print(f"After OverSampling, count of label 1: {sum(y_train_res == 1)}")
print(f"After OverSampling, count of label 0: {sum(y_train_res == 0)}")
```

Before OverSampling, count of label 1: 202401
Before OverSampling, count of label 0: 46712
After OverSampling, count of label 1: 202401
After OverSampling, count of label 0: 202401

```
model = LogisticRegression()
model.fit(X_train_res, y_train_res)
train_preds = model.predict(X_train)
test_preds = model.predict(X_test)
#Model Evaluation
print('Train Accuracy :', model.score(X_train, y_train).round(2))
print('Train F1 Score:',f1_score(y_train,train_preds).round(2))
print('Train Recall Score:',recall_score(y_train,train_preds).round(2))
print('Train Precision Score:',precision_score(y_train,train_preds).round(2))
print('\nTest Accuracy :',model.score(X_test,y_test).round(2))
print('Test F1 Score:',f1_score(y_test,test_preds).round(2))
print('Test Recall Score:',recall_score(y_test,test_preds).round(2))
print('Test Precision Score:',precision_score(y_test,test_preds).round(2))
# Confusion Matrix
cm = confusion_matrix(y_test, test_preds)
disp = ConfusionMatrixDisplay(cm)
disp.plot()
plt.title('Confusion Matrix')
plt.show()
    Train Accuracy: 0.79
     Train F1 Score: 0.86
    Train Recall Score: 0.79
    Train Precision Score: 0.95
```

Test Accuracy : 0.8 Test F1 Score: 0.86 Test Recall Score: 0.79



y_pred = test_preds
print(classification_report(y_test,y_pred))

→	precision	recall	f1-score	support
0 1	0.47 0.95	0.81 0.79	0.60 0.86	11678 50601
accuracy macro avg weighted avg	0.71 0.86	0.80 0.80	0.80 0.73 0.81	62279 62279 62279

Observations:

- The model demonstrates a high recall score, successfully identifying 80% of actual defaulters.
- · However, the precision for the positive class (defaulters) is low; only 47% of predicted defaulters are actually defaulters.
- This high recall and low precision indicate that while the model is effective at flagging most defaulters, it also results in many false positives. Consequently, many deserving customers may be denied loans.
- The low precision adversely affects the F1 score, reducing it to 60%, despite an overall accuracy of 80%. This highlights the trade-off between precision and recall in the model's performance.

Explanation:

- The model is good at catching most people who don't pay back their loans it catches 80% of them.
- But, when it says someone won't pay back, it's right only half of the time.47% So, there's a chance it's making mistakes and wrongly flagging people.
- Because of these mistakes, some people who deserve loans might not get them.
- Even though the model seems okay overall, its balance between being right and not making mistakes isn't great. It's like a seesaw; when one side goes up, the other goes down.

✓ ■Regularization Model

```
#Try with different regularization factor lamda and choose the best to build the model
```

```
lamb = np.arange(0.01, 10000, 10)

train_scores = []

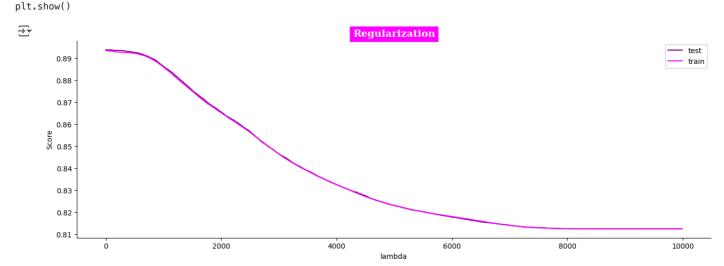
for lam in lamb:
    model = LogisticRegression(C = 1/lam)
    model.fit(X_train, y_train)

    tr_score = model.score(X_train, y_train)
    te_score = model.score(X_test, y_test)

    train_scores.append(tr_score)
    test_scores.append(te_score)

#Plot the train and test scores with respect lambda values i
```

```
#Plot the train and test scores with respect lambda values i.e. regularization factors
ran = np.arange(0.01, 10000, 10)
plt.figure(figsize=(16,5))
sns.lineplot(x=ran,y=test_scores,color='purple',label='test')
sns.lineplot(x=ran,y=train_scores,color='magenta',label='train')
plt.title('Regularization',fontsize=14,fontfamily='serif',fontweight='bold',backgroundcolor='magenta',color='w')
plt.xlabel("lambda")
plt.ylabel("Score")
sns.despine()
```



```
#Check the index of best test score and the check the best test score
print(np.argmax(test_scores))
print(test_scores[np.argmax(test_scores)])

$\frac{1}{2}$
0.8939289327060486
```

#Calculate the best lambda value based on the index of best test score

```
best_lamb = 0.01 + (10*2)
best_lamb
```

```
<del>→</del> 20.01
```

```
#Fit the model using best lambda
```

```
reg_model = LogisticRegression(C=1/best_lamb)
reg_model.fit(X_train, y_train)
```

```
LogisticRegression (C=0.04997501249375312)
```

#Predict the y_values and y_probability values

```
y_reg_pred = reg_model.predict(X_test)
y_reg_pred_proba = reg_model.predict_proba(X_test)
```

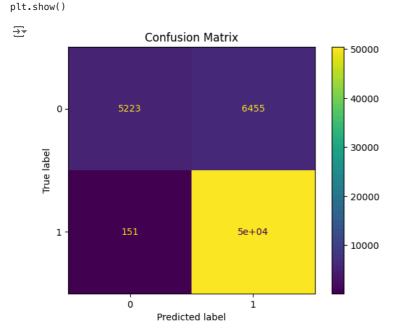
#Print model score

print(f'Logistic Regression Model Score with best lambda: ',end='')
print(round(model.score(X_test, y_test)*100,2),'%')

→ Logistic Regression Model Score with best lambda: 89.39 %

Confusion Matrix

```
cm = confusion_matrix(y_test, y_reg_pred)
disp = ConfusionMatrixDisplay(cm)
disp.plot()
plt.title('Confusion Matrix')
```



print(classification_report(y_test, y_reg_pred))

₹	precision	recall	f1-score	support
0 1	0.97 0.89	0.45 1.00	0.61 0.94	11678 50601
accuracy macro avg weighted avg	0.93 0.90	0.72 0.89	0.89 0.78 0.88	62279 62279 62279

→ ○ Observations from classification report:

Regularized model

Precision: 89%Recall: 100%F1-score: 94%Accuracy: 89%

◆ K-fold - Cross_validation

· cross validation accuracy has to be approx 89%

```
x=scaler.fit_transform(X)
kfold = KFold(n_splits=10)
accuracy = np.mean(cross_val_score(reg_model,x,y,cv=kfold,scoring='accuracy'))
print("Cross Validation accuracy : {:.3f}".format(accuracy))
Fr Cross Validation accuracy: 0.894
cm = confusion_matrix(y_test, y_reg_pred)
cm_df = pd.DataFrame(cm, index=['Defaulter','Fully paid'], columns=['Defaulter','Fully paid'])
cm_df
₹
               Defaulter Fully paid
                    5223
      Defaulter
                                6455
     Fully paid
                     151
                               50450
```

✓ √ Insights:

- TN = 5223 (True Negative: Correctly predicted Charged Off)
- TP = 50450 (True Positive: Correctly predicted Fully Paid)
- FP = 6455 (False Positive: Predicted Fully Paid but actually Charged Off)
- FN = 151 (False Negative: Predicted Charged Off but actually Fully Paid)
- Actual Negative (Charged Off) = 5223 + 6455 = 11678
- Actual Positive (Fully Paid) = 151 + 50450 = 50601
- Predicted Negative (Charged Off) = 5223 + 151 = 5374
- Predicted Positive (Fully Paid) = 6455 + 50450 = 56905

```
#Collect the model coefficients and print those in dataframe format
coeff_df = pd.DataFrame()
coeff_df['Features'] = X_train_res.columns
coeff_df['Weights'] = model.coef_[0]
coeff_df['ABS_Weights'] = abs(coeff_df['Weights'])
coeff_df = coeff_df.sort_values(['ABS_Weights'], ascending=False)
coeff_df
```

.,			
₹	Features	Weights	ABS_Weights
13	zipcode_11650	-7.658994	7.658994
20	zipcode_93700	-7.655336	7.655336
19	zipcode_86630	-7.631667	7.631667
17	zipcode_48052	-2.484366	2.484366
12	zipcode_05113	2.473869	2.473869
15	zipcode_29597	2.466530	2.466530
16	zipcode_30723	-2.442974	2.442974
18	zipcode_70466	-2.432947	2.432947
14	zipcode_22690	-2.425458	2.425458
4	annual_inc	1.159623	1.159623
5	dti	-1.147357	1.147357
24	grade_E	-1.134186	1.134186
23	grade_D	-0.968284	0.968284
22	grade_C	-0.764751	0.764751
25	grade_F	-0.746807	0.746807
37	purpose_small_business	-0.538707	0.538707
6	open_acc	-0.500688	0.500688
1	term	-0.492242	0.492242
2	int_rate	-0.436760	0.436760
9	total_acc	0.413223	0.413223
21		-0.374553	0.374553
39	purpose_wedding	0.342119	0.342119
8		-0.336643	0.336643
47	application_type_INDIVIDUAL		0.301003
3	installment		0.273351
7		0.260325	0.260325
26		-0.244293	0.244293
48	application_type_JOINT	0.239988	0.239988
	erification status Source Verified		0.206324
40			0.205282
30	home_ownership_MORTGAGE purpose home improvement		0.202956
<pre>imp_featur plt.figure sns.barplo plt.title(plt.xlabel plt.yticks</pre>	re = coeff_df.sort_values(re(figsize=(15,10)) re(y = imp_feature['Feature x = imp_feature['Weights "Feature Importance for M .("Weights") re(fontsize=8) .("Features") re()	by='Weight es'], '],color=	ts',ascending= 'm')