Importing necessary libraries -

```
1  import pandas as pd
2  import numpy as np
3  import seaborn as sns
4  from scipy import stats
5  import matplotlib.pyplot as plt
6  from sklearn.linear_model import LogisticRegression
7  from sklearn import metrics
8  from sklearn.metrics import confusion_matrix
9  from sklearn.metrics import classification_report
10 from sklearn.metrics import roc_auc_score
11 from sklearn.metrics import roc_curve
12 from sklearn.metrics import precision_recall_curve
13 from sklearn.model_selection import train_test_split, KFold, cross_val_score
14 from sklearn.preprocessing import MinMaxScaler
15 from sklearn.metrics import (
16     accuracy_score, confusion_matrix, classification_report,
17     roc_auc_score, roc_curve, auc,
18     plot_confusion_matrix, plot_roc_curve
19 )
20 from statsmodels.stats.outliers_influence import variance_inflation_factor
21 from imblearn.over_sampling import SMOTE
```

```
---------------------------------------------------------------------------
ImportError                               Traceback (most recent call last)
<ipython-input-1-e16e9118419d> in <cell line: 15>()
     13 from sklearn.model_selection import train_test_split, KFold, cross_val_score
     14 from sklearn.preprocessing import MinMaxScaler
---> 15 from sklearn.metrics import (
     16     accuracy_score, confusion_matrix, classification_report,
     17     roc_auc_score, roc_curve, auc,

ImportError: cannot import name 'plot_confusion_matrix' from 'sklearn.metrics' (/usr/local/lib/python3.10/dist-
packages/sklearn/metrics/__init__.py)

---------------------------------------------------------------------------
NOTE: If your import is failing due to a missing package, you can
manually install dependencies using either !pip or !apt.

To view examples of installing some common dependencies, click the
"Open Examples" button below.
---------------------------------------------------------------------------
```

OPEN EXAMPLES

Here is the information on this particular data set:

0. 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.

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

2. int_rate : Interest Rate on the loan

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

4. grade LC : assigned loan grade

5. sub_grade LC : assigned loan subgrade

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

7. 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.

8. home_ownership : The home ownership status provided by the borrower during registration or obtained from the credit report. Our values are: RENT, OWN, MORTGAGE, OTHER

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

10. verification_status : Indicates if income was verified by LC, not verified, or if the income source was verified

11. issue_d : The month which the loan was funded

12. loan_status : Current status of the loan

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

14. title : The loan title provided by the borrower

15. zip_code : The first 3 numbers of the zip code provided by the borrower in the loan application.

16. addr_state : The state provided by the borrower in the loan application

17. dti : A ratio calculated using the borrower's total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower's self-reported monthly income.
18. earliest_cr_line : The month the borrower's earliest reported credit line was opened
19. open_acc : The number of open credit lines in the borrower's credit file.
20. pub_rec : Number of derogatory public records
21. revol_bal : Total credit revolving balance
22. revol_util : Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.
23. total_acc : The total number of credit lines currently in the borrower's credit file
24. initial_list_status : The initial listing status of the loan. Possible values are – W, F
25. application_type : Indicates whether the loan is an individual application or a joint application with two co-borrowers
26. mort_acc : Number of mortgage accounts.
27. pub_rec_bankruptcies : Number of public record bankruptcies

Reading the data file -

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

    Mounted at /content/drive

```
1 !ls
```

    drive   logistic_regression.csv   sample_data

```
1 data = pd.read_csv('logistic_regression.csv')
2 data.head()
```

| annual_inc | ... | open_acc | pub_rec | revol_bal | revol_util | total_acc | initial_list_status |
|---|---|---|---|---|---|---|---|
| 117000.0 | ... | 16.0 | 0.0 | 36369.0 | 41.8 | 25.0 | w |
| 65000.0 | ... | 17.0 | 0.0 | 20131.0 | 53.3 | 27.0 | f |
| 43057.0 | ... | 13.0 | 0.0 | 11987.0 | 92.2 | 26.0 | f |
| 54000.0 | ... | 6.0 | 0.0 | 5472.0 | 21.5 | 13.0 | f |
| 55000.0 | ... | 13.0 | 0.0 | 24584.0 | 69.8 | 43.0 | f |

```
1 # Shape of the dataset -
2 print("No. of rows: ", data.shape[0])
3 print("No. of columns: ", data.shape[1])
```

    No. of rows:  231746
    No. of columns:  27

```
1 # Checking the distribution of outcome labels -
2 data.loan_status.value_counts(normalize=True)*100
```

    Fully Paid     80.423395
    Charged Off    19.576605
    Name: loan_status, dtype: float64

```
1 # Statistical summary of the dataset -
2 data.describe(include='all')
```

| | loan_amnt | term | int_rate | installment | grade | sub_grade | emp_title |
|---|---|---|---|---|---|---|---|
| count | 231746.000000 | 231746 | 231746.000000 | 231746.000000 | 231746 | 231746 | 218320 |
| unique | NaN | 2 | NaN | NaN | 7 | 35 | 109777 |
| top | NaN | 36 months | NaN | NaN | B | B3 | Teacher |
| freq | NaN | 176621 | NaN | NaN | 67866 | 15546 | 2557 |
| mean | 14107.757955 | NaN | 13.640928 | 431.586377 | NaN | NaN | NaN |
| std | 8353.939311 | NaN | 4.466482 | 250.592582 | NaN | NaN | NaN |
| min | 500.000000 | NaN | 5.320000 | 16.250000 | NaN | NaN | NaN |
| 25% | 8000.000000 | NaN | 10.490000 | 250.340000 | NaN | NaN | NaN |
| 50% | 12000.000000 | NaN | 13.330000 | 375.430000 | NaN | NaN | NaN |
| 75% | 20000.000000 | NaN | 16.490000 | 567.040000 | NaN | NaN | NaN |
| max | 40000.000000 | NaN | 30.990000 | 1533.810000 | NaN | NaN | NaN |

11 rows × 27 columns

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 231746 entries, 0 to 231745
Data columns (total 27 columns):
 #   Column              Non-Null Count   Dtype
---  ------              --------------   -----
 0   loan_amnt           231746 non-null  float64
 1   term                231746 non-null  object
 2   int_rate            231746 non-null  float64
 3   installment         231746 non-null  float64
 4   grade               231746 non-null  object
 5   sub_grade           231746 non-null  object
 6   emp_title           218320 non-null  object
 7   emp_length          221005 non-null  object
 8   home_ownership      231746 non-null  object
 9   annual_inc          231746 non-null  float64
 10  verification_status 231746 non-null  object
 11  issue_d             231746 non-null  object
 12  loan_status         231746 non-null  object
 13  purpose             231746 non-null  object
 14  title               230723 non-null  object
 15  dti                 231746 non-null  float64
 16  earliest_cr_line    231746 non-null  object
 17  open_acc            231746 non-null  float64
 18  pub_rec             231746 non-null  float64
 19  revol_bal           231746 non-null  float64
 20  revol_util          231576 non-null  float64
 21  total_acc           231746 non-null  float64
 22  initial_list_status 231746 non-null  object
 23  application_type    231746 non-null  object
 24  mort_acc            209558 non-null  float64
 25  pub_rec_bankruptcies 231426 non-null  float64
 26  address             231745 non-null  object
dtypes: float64(12), object(15)
memory usage: 47.7+ MB
```

## Correlation Heatmap -

A correlation heatmap is a heatmap that shows a 2D correlation matrix between two discrete dimensions, using colored cells to represent data from usually a monochromatic scale. The values of the first dimension appear as the rows of the table while of the second dimension as a column. The color of the cell is proportional to the number of measurements that match the dimensional value. This makes correlation heatmaps ideal for data analysis since it makes patterns easily readable and highlights the differences and variation in the same data. A correlation heatmap, like a regular heatmap, is assisted by a colorbar making data easily readable and comprehensible.

```
1 plt.figure(figsize=(12, 8))
2 sns.heatmap(data.corr(method='spearman'), annot=True, cmap='viridis')
3 plt.show()
```
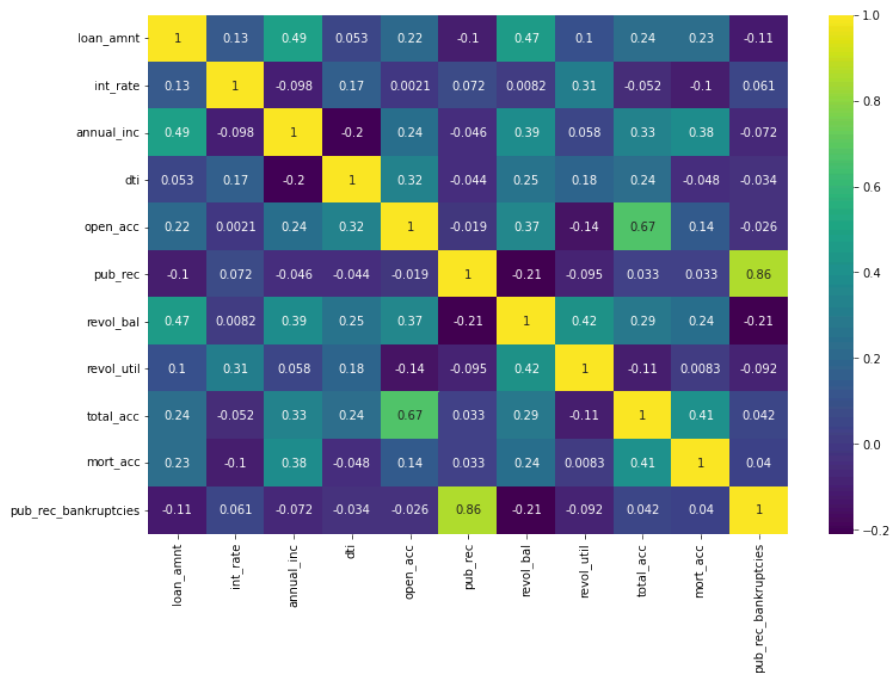
We noticed almost perfect correlation between "loan_amnt" the "installment" feature.

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

So, we can drop either one of those columns.

```
1 data.drop(columns=['installment'], axis=1, inplace=True)
```

```
1 plt.figure(figsize=(12, 8))
2 sns.heatmap(data.corr(method='spearman'), annot=True, cmap='viridis')
3 plt.show()
```

## Data Exploration -

1. The no of people those who have fully paid are 318357 and that of Charged Off are 77673.

```
1 data.groupby(by='loan_status')['loan_amnt'].describe()
```

| loan_status | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| Charged Off | 45368.0 | 15135.654536 | 8484.848789 | 1000.0 | 8700.0 | 14000.0 | 20000.0 | 40000.0 |
| Fully Paid | 186378.0 | 13857.548101 | 8302.548155 | 500.0 | 7500.0 | 12000.0 | 19200.0 | 40000.0 |

2. The majority of people have home ownership as Mortgage and Rent.

```
1 data['home_ownership'].value_counts()
```

```
MORTGAGE    116104
RENT         93551
OWN          22010
OTHER           59
NONE            19
ANY              3
Name: home_ownership, dtype: int64
```

3. Combininging the minority classes as 'OTHER'.

```
1 data.loc[(data.home_ownership == 'ANY') | (data.home_ownership == 'NONE'), 'home_ownership'] = 'OTHER'
2 data.home_ownership.value_counts()
```

```
MORTGAGE    116104
RENT         93551
OWN          22010
OTHER           81
Name: home_ownership, dtype: int64
```

```
1 data['home_ownership'].value_counts()
```

```
     MORTGAGE     116104
     RENT          93551
     OWN           22010
     OTHER            81
     Name: home_ownership, dtype: int64
```

```
1 # Checking the distribution of 'Other' -
2 data.loc[data['home_ownership']=='OTHER', 'loan_status'].value_counts()
```

```
     Fully Paid     64
     Charged Off    17
     Name: loan_status, dtype: int64
```

4. Coverting string to date-time format.

```
1 data['issue_d'] = pd.to_datetime(data['issue_d'])
2 data['earliest_cr_line'] = pd.to_datetime(data['earliest_cr_line'])
```

5. Saw some issues in title (Looks like it was filled manually and needs some fixing).

```
1 data['title'].value_counts()[:20]
```

```
     Debt consolidation          89183
     Credit card refinancing     30086
     Home improvement             8932
     Other                        7597
     Debt Consolidation          6854
     Major purchase              2887
     Consolidation               2305
     debt consolidation          2055
     Business                    1701
     Debt Consolidation Loan     1673
     Medical expenses            1605
     Car financing               1242
     Credit Card Consolidation   1049
     Vacation                    1025
     Moving and relocation        986
     consolidation                971
     Personal Loan                913
     Home Improvement             746
     Consolidation Loan           739
     Home buying                  674
     Name: title, dtype: int64
```

```
1 data['title'] = data.title.str.lower()
```

```
1 data.title.value_counts()[:10]
```

```
     debt consolidation          98370
     credit card refinancing     30256
     home improvement            10011
     other                        7637
     consolidation                3344
     major purchase               3024
     debt consolidation loan      2059
     business                     1749
     medical expenses             1652
     credit card consolidation    1567
     Name: title, dtype: int64
```
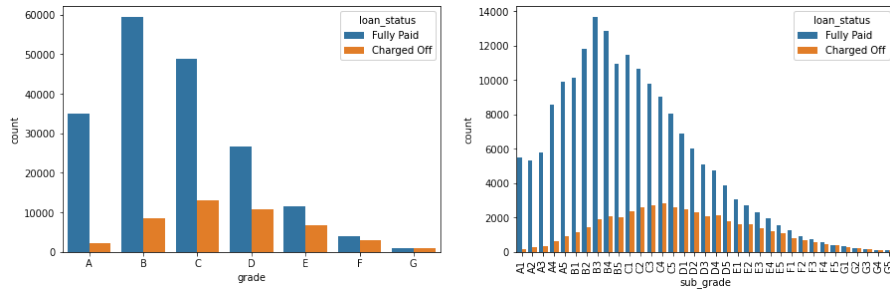
## ⌄ Visualization -

The grade of majority of people those who have fully paid the loan is 'B' and have subgrade 'B3'.

So from where we can infer that people with grade 'B' and subgrade 'B3' are more likely to fully pay the loan.

```
1 plt.figure(figsize=(15, 10))
2
3 plt.subplot(2, 2, 1)
4 grade = sorted(data.grade.unique().tolist())
5 sns.countplot(x='grade', data=data, hue='loan_status', order=grade)
6
7 plt.subplot(2, 2, 2)
8 sub_grade = sorted(data.sub_grade.unique().tolist())
9 g = sns.countplot(x='sub_grade', data=data, hue='loan_status', order=sub_grade)
10 g.set_xticklabels(g.get_xticklabels(), rotation=90);
```
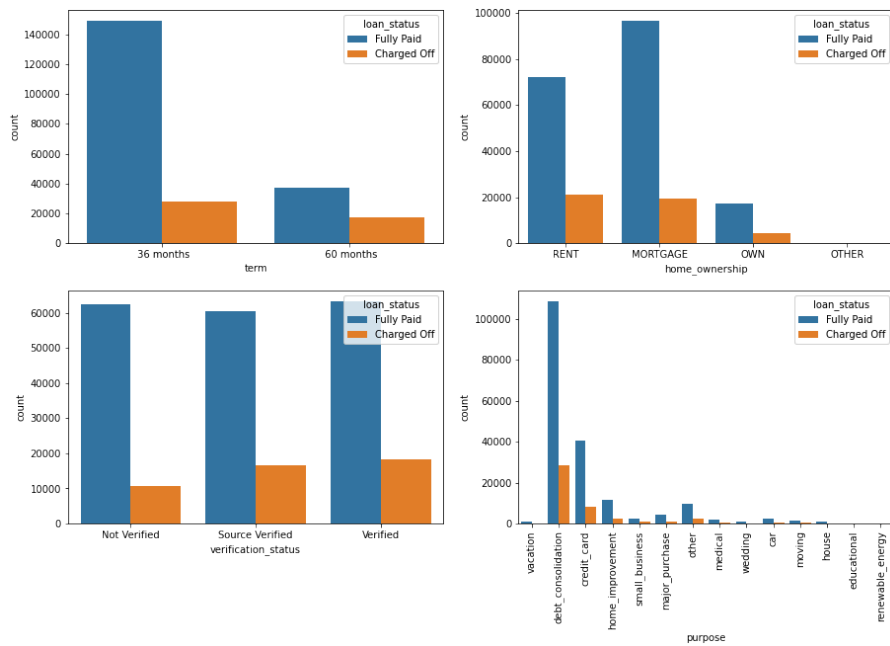


```
1 grade
```

```
['A', 'B', 'C', 'D', 'E', 'F', 'G']
```

```
1 plt.figure(figsize=(15, 20))
2
3 plt.subplot(4, 2, 1)
4 sns.countplot(x='term', data=data, hue='loan_status')
5
6 plt.subplot(4, 2, 2)
7 sns.countplot(x='home_ownership', data=data, hue='loan_status')
8
9 plt.subplot(4, 2, 3)
10 sns.countplot(x='verification_status', data=data, hue='loan_status')
11
12 plt.subplot(4, 2, 4)
13 g = sns.countplot(x='purpose', data=data, hue='loan_status')
14 g.set_xticklabels(g.get_xticklabels(), rotation=90);
```

1 Start coding or generate with AI.

Manager and Teacher are the most afforded loan job titles.

```
 1 plt.figure(figsize=(15, 12))
 2
 3 plt.subplot(2, 2, 1)
 4 order = ['< 1 year', '1 year', '2 years', '3 years', '4 years', '5 years',
 5          '6 years', '7 years', '8 years', '9 years', '10+ years',]
 6 g = sns.countplot(x='emp_length', data=data, hue='loan_status', order=order)
 7 g.set_xticklabels(g.get_xticklabels(), rotation=90);
 8
 9 plt.subplot(2, 2, 2)
10 plt.barh(data.emp_title.value_counts()[:30].index, data.emp_title.value_counts()[:30])
11 plt.title("The most 30 jobs title afforded a loan")
12 plt.tight_layout()
```

## Feature Engineering -

```python
1 def pub_rec(number):
2     if number == 0.0:
3         return 0
4     else:
5         return 1
6
7 def mort_acc(number):
8     if number == 0.0:
9         return 0
10     else:
11         return 1
12
13
14 def pub_rec_bankruptcies(number):
15     if number == 0.0:
16         return 0
17     else:
18         return 1
```

```python
1 data['pub_rec'] = data.pub_rec.apply(pub_rec)
2 data['mort_acc'] = data.mort_acc.apply(mort_acc)
3 data['pub_rec_bankruptcies'] = data.pub_rec_bankruptcies.apply(pub_rec_bankruptcies)
```

```python
1 plt.figure(figsize=(12, 30))
2
3 plt.subplot(6, 2, 1)
4 sns.countplot(x='pub_rec', data=data, hue='loan_status')
5
6 plt.subplot(6, 2, 2)
7 sns.countplot(x='initial_list_status', data=data, hue='loan_status')
8
9 plt.subplot(6, 2, 3)
10 sns.countplot(x='application_type', data=data, hue='loan_status')
11
12 plt.subplot(6, 2, 4)
13 sns.countplot(x='mort_acc', data=data, hue='loan_status')
14
15 plt.subplot(6, 2, 5)
16 sns.countplot(x='pub_rec_bankruptcies', data=data, hue='loan_status')
17
18 plt.show()
```

```
1 # Mapping of target variable -
2 data['loan_status'] = data.loan_status.map({'Fully Paid':0, 'Charged Off':1})
```

```
1 data.isnull().sum()/len(data)*100
```

```
loan_amnt              0.000000
term                   0.000000
int_rate               0.000000
grade                  0.000000
sub_grade              0.000000
emp_title              5.789208
emp_length             4.621115
home_ownership         0.000000
annual_inc             0.000000
verification_status    0.000000
issue_d                0.000000
loan_status            0.000000
purpose                0.000000
title                  0.443148
dti                    0.000000
earliest_cr_line       0.000000
open_acc               0.000000
pub_rec                0.000000
revol_bal              0.000000
revol_util             0.069692
total_acc              0.000000
initial_list_status    0.000000
application_type       0.000000
mort_acc               9.543469
pub_rec_bankruptcies   0.135091
address                0.000000
dtype: float64
```

## ⌄ Very Important: Mean Target Imputation

```
1 data.groupby(by='total_acc').mean()
```

| total_acc | loan_amnt | int_rate | annual_inc | loan_status | dti | open_acc | pub_ |
|---|---|---|---|---|---|---|---|
| 2.0 | 6672.222222 | 15.801111 | 64277.777778 | 0.222222 | 2.279444 | 1.611111 | 0.00( |
| 3.0 | 6042.966361 | 15.615566 | 41270.753884 | 0.220183 | 6.502813 | 2.611621 | 0.03: |
| 4.0 | 7587.399031 | 15.069491 | 42426.565969 | 0.214055 | 8.411963 | 3.324717 | 0.03 |
| 5.0 | 7845.734714 | 14.917564 | 44394.098003 | 0.203156 | 10.118328 | 3.921598 | 0.05! |
| 6.0 | 8529.019843 | 14.651752 | 48470.001156 | 0.215874 | 11.222542 | 4.511119 | 0.07( |
| ... | ... | ... | ... | ... | ... | ... | |
| 124.0 | 23200.000000 | 17.860000 | 66000.000000 | 1.000000 | 14.040000 | 43.000000 | 0.00( |
| 129.0 | 25000.000000 | 7.890000 | 200000.000000 | 0.000000 | 8.900000 | 48.000000 | 0.00( |
| 135.0 | 24000.000000 | 15.410000 | 82000.000000 | 0.000000 | 33.850000 | 57.000000 | 0.00( |
| 150.0 | 35000.000000 | 8.670000 | 189000.000000 | 0.000000 | 6.630000 | 40.000000 | 0.00( |
| 151.0 | 35000.000000 | 13.990000 | 160000.000000 | 1.000000 | 12.650000 | 26.000000 | 0.00( |

118 rows × 11 columns

```
1 total_acc_avg = data.groupby(by='total_acc').mean().mort_acc
2 # Saving mean of mort_acc according to total_acc_avg (you can pick any variable for your understanding)
```

```
1 def fill_mort_acc(total_acc, mort_acc):
2     if np.isnan(mort_acc):
3         return total_acc_avg[total_acc].round()
4     else:
5         return mort_acc
```

```
1 data['mort_acc'] = data.apply(lambda x: fill_mort_acc(x['total_acc'], x['mort_acc']), axis=1)
```

```
1 data.isnull().sum()/len(data)*100
```

```
loan_amnt             0.000000
term                  0.000000
int_rate              0.000000
grade                 0.000000
sub_grade             0.000000
emp_title             5.789208
emp_length            4.621115
home_ownership        0.000000
annual_inc            0.000000
verification_status   0.000000
issue_d               0.000000
loan_status           0.000000
purpose               0.000000
title                 0.443148
dti                   0.000000
earliest_cr_line      0.000000
open_acc              0.000000
pub_rec               0.000000
revol_bal             0.000000
revol_util            0.069692
total_acc             0.000000
initial_list_status   0.000000
application_type      0.000000
mort_acc              0.000000
pub_rec_bankruptcies  0.135091
address               0.000000
dtype: float64
```

```
1 # Current no. of rows -
2 data.shape
```

```
(396030, 26)
```

```
1 # Dropping rows with null values -
2 data.dropna(inplace=True)
```

```
1 # Remaining no. of rows -
2 data.shape
```
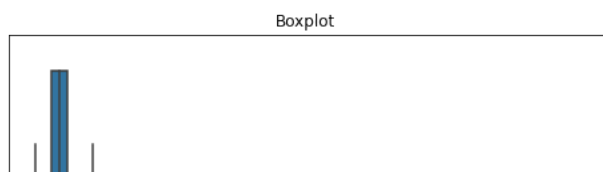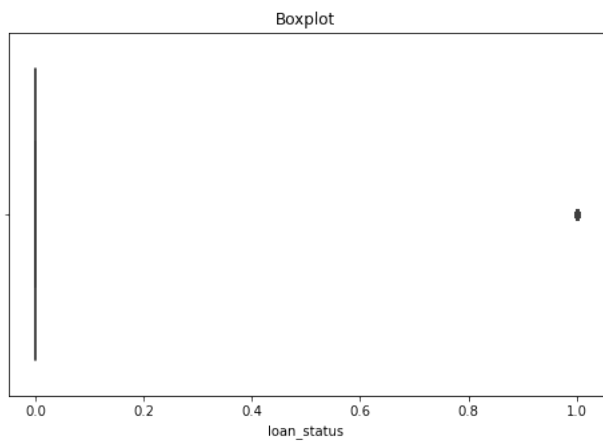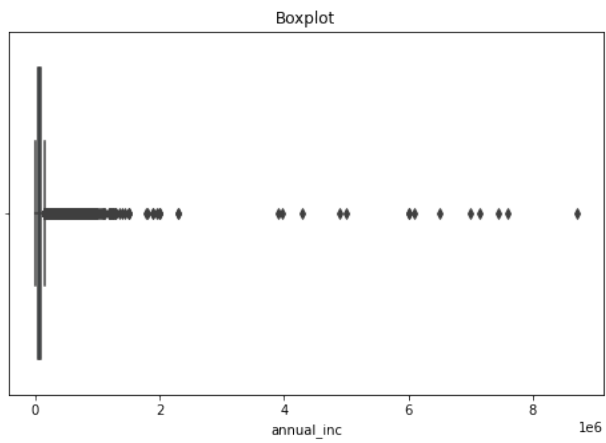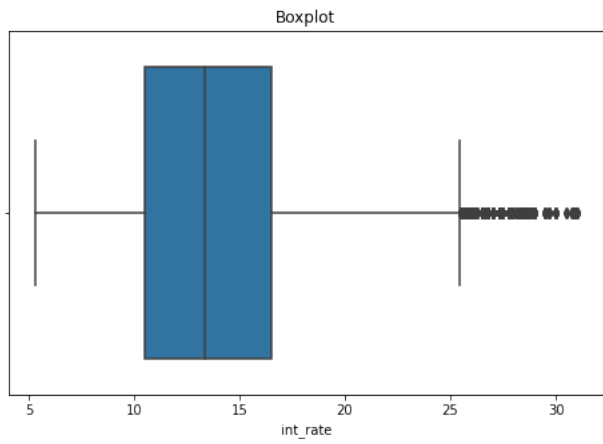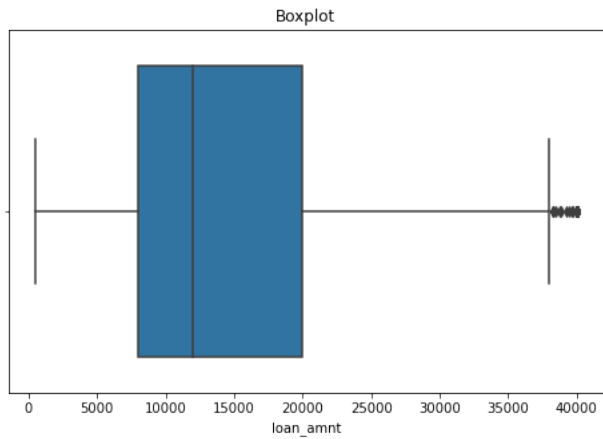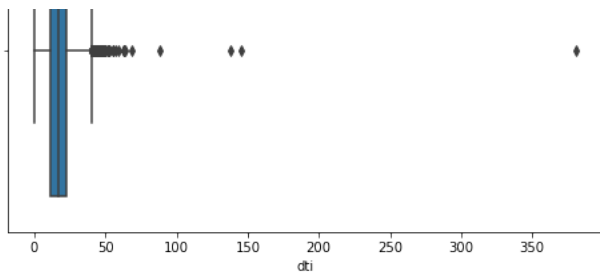
```
(370622, 26)
```

## Outlier Detection & Treatment -

```
1 numerical_data = data.select_dtypes(include='number')
2 num_cols = numerical_data.columns
3 len(num_cols)
```
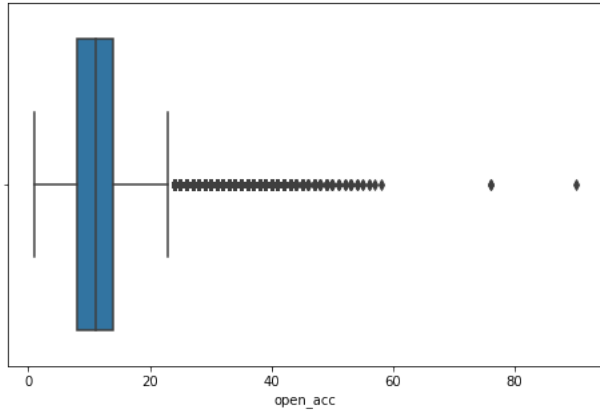
```
12
```

```
1 def box_plot(col):
2     plt.figure(figsize=(8, 5))
3     sns.boxplot(x=data[col])
4     plt.title('Boxplot')
5     plt.show()
6
7 for col in num_cols:
8     box_plot(col)
```
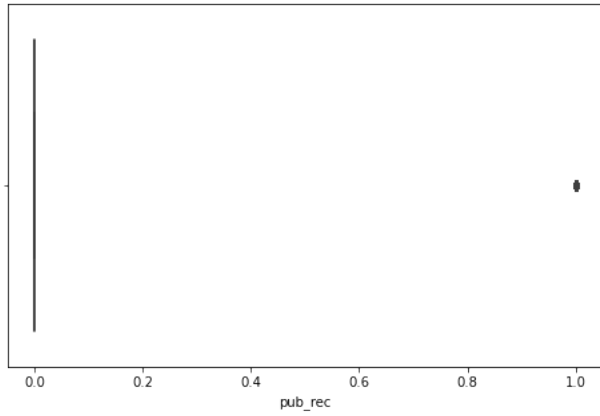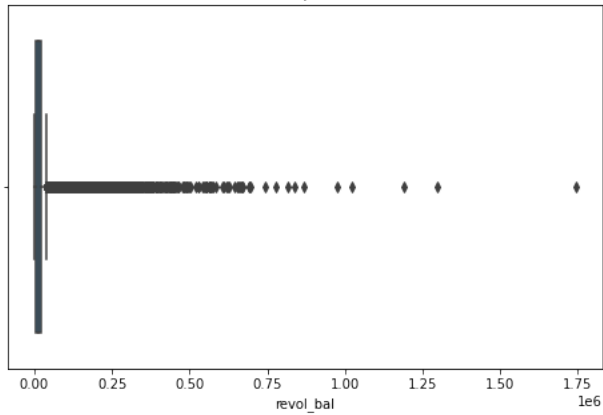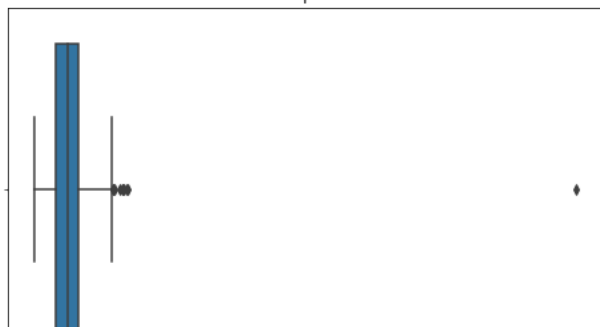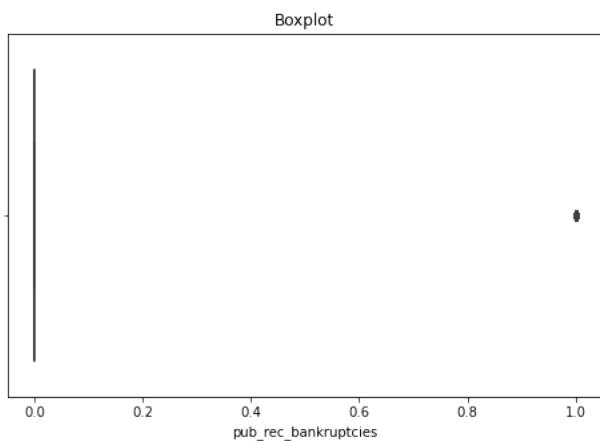
Boxplot



loan_amnt

Boxplot



int_rate

Boxplot



annual_inc

Boxplot



loan_status

Boxplot

Boxplot



Boxplot



Boxplot



Boxplot

Boxplot



total_acc
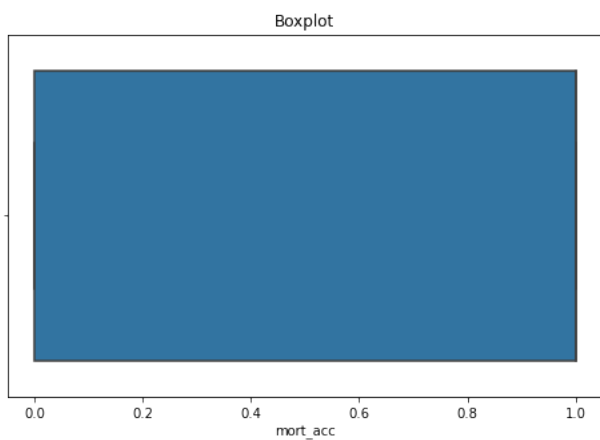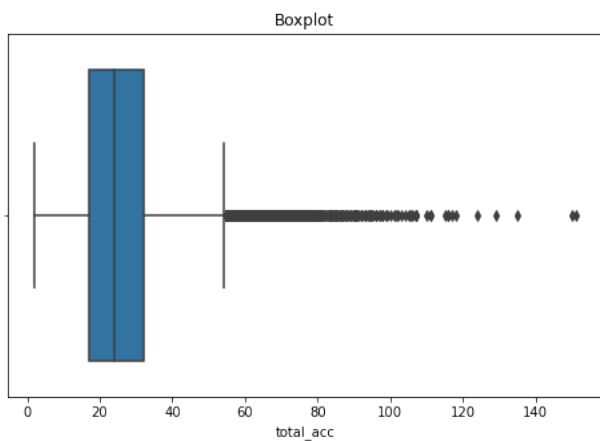
Boxplot



mort_acc

Boxplot



pub_rec_bankruptcies
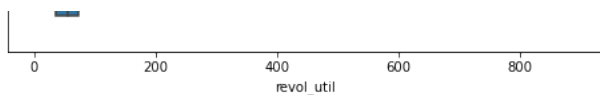
```
1 for col in num_cols:
2     mean = data[col].mean()
3     std = data[col].std()
4
5     upper_limit = mean+3*std
6     lower_limit = mean-3*std
7
8     data = data[(data[col]<upper_limit) & (data[col]>lower_limit)]
9
10 data.shape
```

```
(354519, 26)
```

## Data Preprocessing -

```
1 # Term -
2 data.term.unique()
```

```
    array([' 36 months', ' 60 months'], dtype=object)
```

```
1 Start coding or generate with AI.
```

```
1 term_values = {' 36 months': 36, ' 60 months': 60}
2 data['term'] = data.term.map(term_values)
```

```
1 # Initial List Status -
2 data['initial_list_status'].unique()
```

```
    array(['w', 'f'], dtype=object)
```

```
1 list_status = {'w': 0, 'f': 1}
2 data['initial_list_status'] = data.initial_list_status.map(list_status)
```

```
1 # Let's fetch ZIP from address and then drop the remaining details -
2 data['zip_code'] = data.address.apply(lambda x: x[-5:])
```

```
1 data['zip_code'].value_counts(normalize=True)*100
```

```
    70466    14.382022
    30723    14.277373
    22690    14.268347
    48052    14.127028
    00813    11.610097
    29597    11.537322
    05113    11.516731
    93700     2.774746
    11650     2.772771
    86630     2.733563
    Name: zip_code, dtype: float64
```

```
1 # Dropping some variables which IMO we can let go for now -
2 data.drop(columns=['issue_d', 'emp_title', 'title', 'sub_grade',
3                    'address', 'earliest_cr_line', 'emp_length'],
4                    axis=1, inplace=True)
```

## One-hot Encoding -

```
1 dummies = ['purpose', 'zip_code', 'grade', 'verification_status', 'application_type', 'home_ownership']
2 data = pd.get_dummies(data, columns=dummies, drop_first=True)
```

```
1 pd.set_option('display.max_columns', None)
2 pd.set_option('display.max_rows', None)
3
4 data.head()
```

|   | loan_amnt | term | int_rate | annual_inc | loan_status | dti | open_acc | pub_rec | revol_ba |
|---|-----------|------|----------|------------|-------------|------|----------|---------|----------|
| 0 | 10000.0 | 36 | 11.44 | 117000.0 | 0 | 26.24 | 16.0 | 0 | 36369 |
| 1 | 8000.0 | 36 | 11.99 | 65000.0 | 0 | 22.05 | 17.0 | 0 | 20131 |
| 2 | 15600.0 | 36 | 10.49 | 43057.0 | 0 | 12.79 | 13.0 | 0 | 11987 |
| 3 | 7200.0 | 36 | 6.49 | 54000.0 | 0 | 2.60 | 6.0 | 0 | 5472 |
| 4 | 24375.0 | 60 | 17.27 | 55000.0 | 1 | 33.95 | 13.0 | 0 | 24584 |

```
1 data.shape
```

```
    (354519, 49)
```

## Data Preparation for Modeling -

```
1 X = data.drop('loan_status', axis=1)
2 y = data['loan_status']
```

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30,
2                                                      stratify=y, random_state=42)
```

```
1 print(X_train.shape)
2 print(X_test.shape)
```

```
    (248163, 48)
    (106356, 48)
```

## MinMaxScaler -

For each value in a feature, MinMaxScaler subtracts the minimum value in the feature and then divides by the range. The range is the difference between the original maximum and original minimum.

MinMaxScaler preserves the shape of the original distribution. It doesn't meaningfully change the information embedded in the original data.

```
1 scaler = MinMaxScaler()
2 X_train = scaler.fit_transform(X_train)
3 X_test = scaler.transform(X_test)
```

## Logistic Regression

```
1 logreg = LogisticRegression(max_iter=1000)
2 logreg.fit(X_train, y_train)
```

```
    LogisticRegression(max_iter=1000)
```

```
1 y_pred = logreg.predict(X_test)
2 print('Accuracy of Logistic Regression Classifier on test set: {:.3f}'.format(logreg.score(X_test, y_test)))
```

```
    Accuracy of Logistic Regression Classifier on test set: 0.890
```

## Confusion Matrix -

```
1 confusion_matrix = confusion_matrix(y_test, y_pred)
2 print(confusion_matrix)
```

```
    [[85365   523]
     [11131  9337]]
```

## Classification Report -

```
1 print(classification_report(y_test, y_pred))
```

```
              precision    recall  f1-score   support

           0       0.88      0.99      0.94     85888
           1       0.95      0.46      0.62     20468

    accuracy                           0.89    106356
   macro avg       0.92      0.73      0.78    106356
weighted avg       0.90      0.89      0.87    106356
```

## ROC Curve -

An ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds. This curve plots two parameters:

- True Positive Rate
- False Positive Rate

True Positive Rate (TPR) is a synonym for recall and is therefore defined as follows:

- TPR=(TP)/(TP+FN)

False Positive Rate (FPR) is defined as follows:
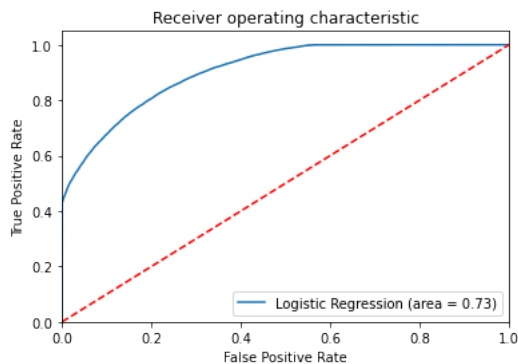
- FPR=(FP)/(FP+TN)

An ROC curve plots TPR vs. FPR at different classification thresholds. Lowering the classification threshold classifies more items as positive, thus increasing both False Positives and True Positives. The following figure shows a typical ROC curve.
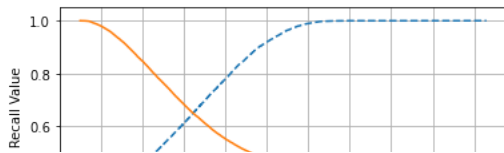
## ⌄ AUC (Area under the ROC Curve) -

AUC stands for "Area under the ROC Curve." That is, AUC measures the entire two-dimensional area underneath the entire ROC curve (think integral calculus) from (0,0) to (1,1).

AUC provides an aggregate measure of performance across all possible classification thresholds. One way of interpreting AUC is as the probability that the model ranks a random positive example more highly than a random negative example. For example, given the following examples, which are arranged from left to right in ascending order of logistic regression predictions:

```
1 logit_roc_auc = roc_auc_score(y_test, logreg.predict(X_test))
2 fpr, tpr, thresholds = roc_curve(y_test, logreg.predict_proba(X_test)[:,1])
3 plt.figure()
4 plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % logit_roc_auc)
5 plt.plot([0, 1], [0, 1],'r--')
6 plt.xlim([0.0, 1.0])
7 plt.ylim([0.0, 1.05])
8 plt.xlabel('False Positive Rate')
9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver operating characteristic')
11 plt.legend(loc="lower right")
12 plt.savefig('Log_ROC')
13 plt.show()
```



```
1 def precision_recall_curve_plot(y_test, pred_proba_c1):
2     precisions, recalls, thresholds = precision_recall_curve(y_test, pred_proba_c1)
3
4     threshold_boundary = thresholds.shape[0]
5     # plot precision
6     plt.plot(thresholds, precisions[0:threshold_boundary], linestyle='--', label='precision')
7     # plot recall
8     plt.plot(thresholds, recalls[0:threshold_boundary], label='recalls')
9
10    start, end = plt.xlim()
11    plt.xticks(np.round(np.arange(start, end, 0.1), 2))
12
13    plt.xlabel('Threshold Value'); plt.ylabel('Precision and Recall Value')
14    plt.legend(); plt.grid()
15    plt.show()
16
17 precision_recall_curve_plot(y_test, logreg.predict_proba(X_test)[:,1])
```

## Multicollinearity check using Variance Inflation Factor (VIF) -

Multicollinearity occurs when two or more independent variables are highly correlated with one another in a regression model. Multicollinearity can be a problem in a regression model because we would not be able to distinguish between the individual effects of the independent variables on the dependent variable.

Multicollinearity can be detected via various methods. One such method is Variance Inflation Factor aka VIF. In VIF method, we pick each independent feature and regress it against all of the other independent features. VIF score of an independent variable represents how well the variable is explained by other independent variables.

VIF = 1/1-R2

```
1 def calc_vif(X):
2     # Calculating the VIF
3     vif = pd.DataFrame()
4     vif['Feature'] = X.columns
5     vif['VIF'] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
6     vif['VIF'] = round(vif['VIF'], 2)
7     vif = vif.sort_values(by='VIF', ascending = False)
8     return vif
9
10 calc_vif(X)[:5]
```

|    | Feature | VIF |
|----|---------|-----|
| 43 | application_type_INDIVIDUAL | 156.97 |
| 2  | int_rate | 122.82 |
| 14 | purpose_debt_consolidation | 51.00 |
| 1  | term | 27.30 |
| 13 | purpose_credit_card | 18.48 |

```
1 X.drop(columns=['application_type_INDIVIDUAL'], axis=1, inplace=True)
2 calc_vif(X)[:5]
```

|    | Feature | VIF |
|----|---------|-----|
| 2  | int_rate | 103.43 |
| 14 | purpose_debt_consolidation | 27.49 |
| 1  | term | 24.31 |
| 5  | open_acc | 13.75 |
| 9  | total_acc | 12.69 |

```
1 X.drop(columns=['int_rate'], axis=1, inplace=True)
2 calc_vif(X)[:5]
```

|    | Feature | VIF |
|----|---------|-----|
| 1  | term | 23.35 |
| 13 | purpose_debt_consolidation | 22.35 |
| 4  | open_acc | 13.64 |
| 8  | total_acc | 12.69 |
| 7  | revol_util | 9.06 |