

Project Flow

- Load Dataset
- EDA
- Create Baseline Models
- Feature Engineering, Transformations & Feature Selections
- Cross Validation of Models
- Finalise Dataset
- Model Selection via CV
- HyperParameter Tuning of Selected Model
- Train Tuned Model using Entire new Data
- Save & Deploy Data
- Predict

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 import warnings
7 warnings.filterwarnings("ignore")
8
9 pd.set_option('display.max_columns', None)
```

In [2]:

```
1 train = pd.read_csv('train.csv')
2 train.shape
```

Out[2]:

(614, 13)

In [3]:

```
1 test = pd.read_csv('test.csv')
2 test.shape
3
4 test_backup = test.copy()
```

In [4]:

```
1 train.head()
```

Out[4]:

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	Coappl
0	LP001002	Male	No	0	Graduate	No	5849	
1	LP001003	Male	Yes	1	Graduate	No	4583	
2	LP001005	Male	Yes	0	Graduate	Yes	3000	
3	LP001006	Male	Yes	0	Not Graduate	No	2583	
4	LP001008	Male	No	0	Graduate	No	6000	

In []:

```
1
```

Data Validation

since we have categorical features, we would encode that, but need to check if all categories present in testing set are present in training data as well, otherwise we need to handle that separately

In [5]:

```

1 categorical_features = ['Gender', 'Married', 'Dependents', 'Education', 'Self_Employed']
2
3 for col in categorical_features:
4     print('=> ', col)
5     print('Train -> ', train[col].unique())
6     print('Test -> ', test[col].unique())
7     print('*'*33)

```

=> Gender

Train -> ['Male' 'Female' nan]

Test -> ['Male' 'Female' nan]

=> Married

Train -> ['No' 'Yes' nan]

Test -> ['Yes' 'No']

=> Dependents

Train -> ['0' '1' '2' '3+' nan]

Test -> ['0' '1' '2' '3+' nan]

=> Education

Train -> ['Graduate' 'Not Graduate']

Test -> ['Graduate' 'Not Graduate']

=> Self_Employed

Train -> ['No' 'Yes' nan]

Test -> ['No' 'Yes' nan]

=> Property_Area

Train -> ['Urban' 'Rural' 'Semiurban']

Test -> ['Urban' 'Semiurban' 'Rural']

There are no unique categories present in Test set, which are not in Training data => Validated

EDA has not been performed yet. Covered in another notebook

Outliers Detection

IQR

In [6]:

```
1 def IQR(df, col):
2
3     # calculate quantiles
4     q25, q75 = np.quantile(df[col], 0.25), np.quantile(df[col], 0.75)
5
6     # calculate IQR
7     iqr = q75 - q25
8
9     # defining boundaries
10    lower, upper = q25 - (1.5 * iqr), q75 + (1.5 * iqr)
11
12    print('IQR is ',iqr)
13    print('lower bound is ', lower)
14    print('upper bound is ', upper)
15
16    # calculate records count beyond range
17    df_upper = df[df[col] > upper]
18    df_lower = df[df[col] < lower]
19
20    print('Total number of outliers are', df_upper.shape[0] + df_lower.shape[0])
21
22    # Visualising
23    plt.figure(figsize = (10,6))
24    sns.distplot(df[col], kde=False)
25    plt.axvspan(xmin = lower,xmax= df[col].min(),alpha=0.2, color='red')
26    plt.axvspan(xmin = upper,xmax= df[col].max(),alpha=0.2, color='red')
```

In [7]:

```
1 numerical_features = [ 'ApplicantIncome', 'CoapplicantIncome', 'LoanAmount', 'Loan_Amount'
```

In [8]:

```

1 for col in numerical_features:
2     print('=> ', col)
3     IQR(train, col)
4     print('***33)

```

```

=> ApplicantIncome
IQR is 2917.5
lower bound is -1498.75
upper bound is 10171.25
Total number of outliers are 50
*****

=> CoapplicantIncome
IQR is 2297.25
lower bound is -3445.875
upper bound is 5743.125
Total number of outliers are 18
*****

=> LoanAmount
IQR is nan
lower bound is nan
upper bound is nan
Total number of outliers are 0
*****

=> Loan_Amount_Term
nan .

```

Z Score

In [9]:

```

1 def zscore(col):
2
3     outliers = []
4     z_scores = []
5     threshold = 3
6
7     # calculate the mean and standard deviation of the data frame
8     data_mean, data_std = col.mean(), col.std()
9
10    for i in col:
11        z_score = (i - data_mean) / data_std
12        z_scores.append(z_score)
13
14        if np.abs(z_score) > threshold:
15            outliers.append(i)
16
17    print('Total Outliers are ', len(outliers))
18
19    # Visualising
20    plt.figure(figsize = (10,6))
21    sns.distplot(z_scores)
22    plt.axvspan(xmin = 3 ,xmax= max(z_scores),alpha=0.2, color='red')

```

In [10]:

```
1 for col in numerical_features:
2     print('=> ', col)
3     zscore(train[col])
4     print('*'*33)
```

=> ApplicantIncome

Total Outliers are 8

=> CoapplicantIncome

Total Outliers are 6

=> LoanAmount

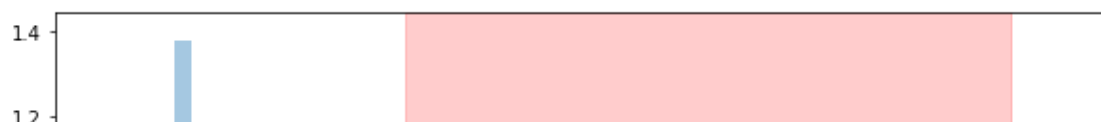
Total Outliers are 14

=> Loan_Amount_Term

Total Outliers are 12

=> Credit_History

Total Outliers are 0



Standard Deviation

In [11]:

```

1  def std(df, col):
2
3      # calculate the mean and standard deviation of the data frame
4      data_mean, data_std = df[col].mean(), df[col].std()
5
6      # initialise the cutoff value
7      cut_off = data_std * 3
8
9      # calculate the lower and upper bound value
10     lower, upper = data_mean - cut_off, data_mean + cut_off
11
12     print('The lower bound value is', lower)
13     print('The upper bound value is', upper)
14
15     # calculate records count beyond range
16     df_upper = df[df[col] > upper]
17     df_lower = df[df[col] < lower]
18
19     print('Total number of outliers are', df_upper.shape[0] + df_lower.shape[0])
20
21     # Visualising
22     plt.figure(figsize = (10,6))
23     sns.distplot(df[col], kde=False)
24     plt.axvspan(xmin = lower,xmax= df[col].min(),alpha=0.2, color='red')
25
26
27

```

In [12]:

```

1  for col in numerical_features:
2      print('=> ', col)
3      std(train, col)
4      print('*'*33)

```

```

=> ApplicantIncome
The lower bound value is -12923.665736773899
The upper bound value is 23730.584303549145
Total number of outliers are 8
*****
=> CoapplicantIncome
The lower bound value is -7157.499309645475
The upper bound value is 10399.990905699677
Total number of outliers are 6
*****
=> LoanAmount
The lower bound value is -110.34981354495417
The upper bound value is 403.1741378692785
Total number of outliers are 14
*****
=> Loan_Amount_Term
The lower bound value is 146.6387704361623
The upper bound value is 537.3612295638377
Total number of outliers are 12
*****

```

In [13]:

```
1 train.dtypes
```

Out[13]:

```
Loan_ID          object
Gender           object
Married          object
Dependents       object
Education        object
Self_Employed    object
ApplicantIncome  int64
CoapplicantIncome float64
LoanAmount       float64
Loan_Amount_Term float64
Credit_History   float64
Property_Area     object
Loan_Status       object
dtype: object
```

In [14]:

```
1 # Removing outliers from ApplicantIncome considering Standard Deviation
2 train = train[(train['ApplicantIncome'] < 23730) & (train['ApplicantIncome'] > -12923)]
```

In [15]:

```
1 # Removing outliers from CoApplicantIncome considering Standard Deviation
2 train = train[(train['CoapplicantIncome'] < 10399) & (train['CoapplicantIncome'] > -715)]
```

In [16]:

```
1 train.shape
```

Out[16]:

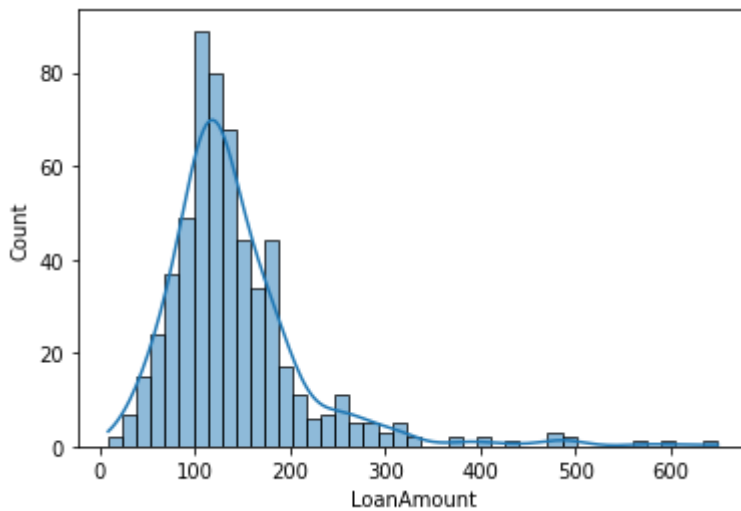
```
(600, 13)
```


In [17]:

```
1 sns.histplot(data= train['LoanAmount'], kde= True)
```

Out[17]:

<AxesSubplot: xlabel='LoanAmount', ylabel='Count'>



In [18]:

```
1 train.info()
```

<class 'pandas.core.frame.DataFrame'>

Int64Index: 600 entries, 0 to 613

Data columns (total 13 columns):

#	Column	Non-Null Count	Dtype
0	Loan_ID	600 non-null	object
1	Gender	589 non-null	object
2	Married	597 non-null	object
3	Dependents	585 non-null	object
4	Education	600 non-null	object
5	Self_Employed	570 non-null	object
6	ApplicantIncome	600 non-null	int64
7	CoapplicantIncome	600 non-null	float64
8	LoanAmount	578 non-null	float64
9	Loan_Amount_Term	586 non-null	float64
10	Credit_History	551 non-null	float64
11	Property_Area	600 non-null	object
12	Loan_Status	600 non-null	object

dtypes: float64(4), int64(1), object(8)

memory usage: 65.6+ KB

In [19]:

```
1 test.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 367 entries, 0 to 366
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID                367 non-null    object
1   Gender                 356 non-null    object
2   Married                367 non-null    object
3   Dependents             357 non-null    object
4   Education              367 non-null    object
5   Self_Employed          344 non-null    object
6   ApplicantIncome        367 non-null    int64
7   CoapplicantIncome      367 non-null    int64
8   LoanAmount             362 non-null    float64
9   Loan_Amount_Term       361 non-null    float64
10  Credit_History         338 non-null    float64
11  Property_Area          367 non-null    object
dtypes: float64(3), int64(2), object(7)
memory usage: 34.5+ KB
```

In [20]:

```
1 train[train.duplicated(keep= 'first')].shape
```

Out[20]:

(0, 13)

In [21]:

```
1 train.isnull().sum()
```

Out[21]:

```
Loan_ID          0
Gender           11
Married           3
Dependents       15
Education         0
Self_Employed    30
ApplicantIncome   0
CoapplicantIncome 0
LoanAmount       22
Loan_Amount_Term  14
Credit_History   49
Property_Area     0
Loan_Status       0
dtype: int64
```

In [22]:

```
1 test.isnull().sum()
```

Out[22]:

```
Loan_ID          0
Gender           11
Married          0
Dependents       10
Education         0
Self_Employed    23
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount        5
Loan_Amount_Term  6
Credit_History   29
Property_Area     0
dtype: int64
```

Splitting Train data into Train & Val

In [23]:

```
1 from sklearn.model_selection import train_test_split
2
3 X = train.drop(['Loan_Status'], axis= 1)
4 y = train['Loan_Status']
5
6 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=9,
7
8 X_train.shape, X_val.shape, y_train.shape, y_val.shape
```

Out[23]:

```
((480, 12), (120, 12), (480,), (120,))
```

In [24]:

```
1 X_train['Gender'].value_counts(normalize= True)
```

Out[24]:

```
Male      0.817021
Female     0.182979
Name: Gender, dtype: float64
```

Why this is important : Suppose if one of the values were dominant (covering > 80% suppose), then we could have used that value to impute in these missing values

if suppose both were too close, and if we assign there missing values to a particular class (using MODE), then the distribution will significantly change, in that case we need to think other way around

check column transformer to do all tasks together

In [25]:

```
1 from sklearn.impute import SimpleImputer
2 from sklearn.preprocessing import OneHotEncoder
3 from sklearn.compose import ColumnTransformer
4 from sklearn.preprocessing import OrdinalEncoder
```

In [26]:

```
1 Imputer = ColumnTransformer(transformers = [
2     ('mode', SimpleImputer(strategy= 'most_frequent'), ['Gender', 'Married', 'Dependent
3     ('mean', SimpleImputer(strategy= 'mean'), ['LoanAmount'])
4 ], remainder= 'drop')
```

In [27]:

```
1 Imputed_data_train = Imputer.fit_transform(X_train)
2 Imputed_data_train
```

Out[27]:

```
array([[ 'Male', 'Yes', '2', ..., 360.0, 1.0, 155.0],
       [ 'Male', 'No', '0', ..., 360.0, 1.0, 311.0],
       [ 'Male', 'Yes', '0', ..., 180.0, 0.0, 90.0],
       ...,
       [ 'Male', 'No', '0', ..., 360.0, 1.0, 144.1459694989107],
       [ 'Female', 'No', '1', ..., 360.0, 1.0, 144.1459694989107],
       [ 'Male', 'Yes', '0', ..., 360.0, 1.0, 144.1459694989107]],
      dtype=object)
```

In [28]:

```
1 Imputer.named_transformers_[ 'mode' ].get_feature_names_out()
```

Out[28]:

```
array([ 'Gender', 'Married', 'Dependents', 'Self_Employed',
       'Loan_Amount_Term', 'Credit_History'], dtype=object)
```

In [29]:

```
1 Imputer.named_transformers_[ 'mean' ].get_feature_names_out()
```

Out[29]:

```
array([ 'LoanAmount'], dtype=object)
```

In [30]:

```
1 new_cols = Imputer.get_feature_names_out()
2 new_cols
```

Out[30]:

```
array([ 'mode__Gender', 'mode__Married', 'mode__Dependents',
       'mode__Self_Employed', 'mode__Loan_Amount_Term',
       'mode__Credit_History', 'mean__LoanAmount'], dtype=object)
```

In [31]:

```
1 Imputed_df_train = pd.DataFrame(Imputed_data_train, columns=new_cols)
2 Imputed_df_train.tail()
```

Out[31]:

	mode__Gender	mode__Married	mode__Dependents	mode__Self_Employed	mode__Loan_A
475	Male	Yes	0	No	
476	Male	Yes	2	No	
477	Male	No	0	No	
478	Female	No	1	Yes	
479	Male	Yes	0	No	

In [32]:

```
1 Imputed_df_train.isnull().sum()
```

Out[32]:

```
mode__Gender      0
mode__Married     0
mode__Dependents  0
mode__Self_Employed  0
mode__Loan_Amount_Term  0
mode__Credit_History  0
mean__LoanAmount  0
dtype: int64
```

In [33]:

```
1 Imputed_df_train.shape
```

Out[33]:

```
(480, 7)
```

In [34]:

```

1 X_train = X_train.reset_index(drop= True)
2 Imputed_df_train = Imputed_df_train.reset_index(drop=True)
3
4 X_train = pd.concat([X_train, Imputed_df_train], axis=1)
5 X_train.tail()

```

Out[34]:

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome
475	LP001578	Male	Yes	0	Graduate	No	2439	
476	LP002832	Male	Yes	2	Graduate	No	8799	
477	LP002697	Male	No	0	Graduate	No	4680	
478	LP001392	Female	No	1	Graduate	Yes	7451	
479	LP001356	Male	Yes	0	Graduate	No	4652	

In [35]:

```
1 X_train.isnull().sum()
```

Out[35]:

```

Loan_ID          0
Gender           10
Married          3
Dependents       14
Education        0
Self_Employed    22
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount       21
Loan_Amount_Term 11
Credit_History   40
Property_Area     0
mode__Gender      0
mode__Married     0
mode__Dependents  0
mode__Self_Employed 0
mode__Loan_Amount_Term 0
mode__Credit_History 0
mean__LoanAmount  0
dtype: int64

```

One Hot Encoding

In [36]:

```
1 encoder = ColumnTransformer(transformers = [  
2     ('encOH', OneHotEncoder(sparse=False, drop='first'), ['mode__Gender', 'mode__Married',  
3     ('encOR', OrdinalEncoder(categories =[[ 'Not Graduate', 'Graduate']]), ['Education'  
4 ], remainder='drop')
```

In [37]:

```
1 encoded_data_train = encoder.fit_transform(X_train)  
2 encoded_data_train
```

Out[37]:

```
array([[1., 1., 0., ..., 1., 0., 1.],  
       [1., 0., 0., ..., 0., 0., 1.],  
       [1., 1., 0., ..., 0., 0., 1.],  
       ...,  
       [1., 0., 0., ..., 1., 0., 1.],  
       [0., 0., 1., ..., 1., 0., 1.],  
       [1., 1., 0., ..., 1., 0., 1.]])
```

In [38]:

```
1 encoder.named_transformers_['encOH'].get_feature_names_out()
```

Out[38]:

```
array(['mode__Gender_Male', 'mode__Married_Yes', 'mode__Dependents_1',  
       'mode__Dependents_2', 'mode__Dependents_3+',  
       'mode__Self_Employed_Yes', 'Property_Area_Semiurban',  
       'Property_Area_Urban'], dtype=object)
```

In [39]:

```
1 new_cols = encoder.get_feature_names_out()  
2 new_cols
```

Out[39]:

```
array(['encOH__mode__Gender_Male', 'encOH__mode__Married_Yes',  
       'encOH__mode__Dependents_1', 'encOH__mode__Dependents_2',  
       'encOH__mode__Dependents_3+', 'encOH__mode__Self_Employed_Yes',  
       'encOH__Property_Area_Semiurban', 'encOH__Property_Area_Urban',  
       'encOR__Education'], dtype=object)
```

In [40]:

```
1 encoded_df_train = pd.DataFrame(encoded_data_train, columns=new_cols, dtype= int)
2 encoded_df_train.head()
```

Out[40]:

	encOH__mode__Gender_Male	encOH__mode__Married_Yes	encOH__mode__Dependents_1	enc
0	1	1	0	
1	1	0	0	
2	1	1	0	
3	1	0	0	
4	1	1	0	

In [41]:

```
1 X_train = X_train.reset_index(drop= True)
2 encoded_df_train = encoded_df_train.reset_index(drop=True)
3
4 X_train = pd.concat([X_train, encoded_df_train], axis=1)
5 X_train.head()
```

Out[41]:

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	Coappl
0	LP002197	Male	Yes	2	Graduate	No	5185	
1	LP002652	Male	No	0	Graduate	No	5815	
2	LP002449	Male	Yes	0	Graduate	No	2483	
3	LP001761	Male	No	0	Graduate	Yes	6400	
4	LP001926	Male	Yes	0	Graduate	No	3704	

In [42]:

```
1 X_train.isnull().sum()
```

Out[42]:

```
Loan_ID          0
Gender           0
Married          3
Dependents       14
Education        0
Self_Employed    22
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount       21
Loan_Amount_Term 11
Credit_History  40
Property_Area    0
mode__Gender     0
mode__Married    0
mode__Dependents 0
mode__Self_Employed 0
mode__Loan_Amount_Term 0
mode__Credit_History 0
```

Imputing & Encoding Validation & Test Dataset

Imputing Val

In [43]:

```
1 Imputed_data_val = Imputer.transform(X_val)
2 new_cols = Imputer.get_feature_names_out()
3 Imputed_df_val = pd.DataFrame(Imputed_data_val, columns=new_cols)
4
5 X_val = X_val.reset_index(drop=True)
6 Imputed_df_val = Imputed_df_val.reset_index(drop=True)
7
8 X_val = pd.concat([X_val, Imputed_df_val], axis=1)
```

Encoding Val

In [44]:

```
1 encoded_data_val = encoder.transform(X_val)
2 new_cols = encoder.get_feature_names_out()
3 encoded_df_val = pd.DataFrame(encoded_data_val, columns=new_cols, dtype=int)
4
5 X_val = X_val.reset_index(drop=True)
6 encoded_df_val = encoded_df_val.reset_index(drop=True)
7
8 X_val = pd.concat([X_val, encoded_df_val], axis=1)
```

Imputing Test

In [45]:

```
1 Imputed_data_test = Imputer.transform(test)
2 new_cols = Imputer.get_feature_names_out()
3 Imputed_df_test = pd.DataFrame(Imputed_data_test, columns=new_cols)
4
5 test = test.reset_index(drop=True)
6 Imputed_df_test = Imputed_df_test.reset_index(drop=True)
7
8 test = pd.concat([test, Imputed_df_test], axis=1)
```

Encoding Test

In [46]:

```
1 encoded_data_test = encoder.transform(test)
2 new_cols = encoder.get_feature_names_out()
3 encoded_df_test = pd.DataFrame(encoded_data_test, columns=new_cols, dtype=int)
4
5 test = test.reset_index(drop=True)
6 encoded_df_test = encoded_df_test.reset_index(drop=True)
7
8 test = pd.concat([test, encoded_df_test], axis=1)
```

Dropping Index Identifier : loan_id

In [47]:

```
1 X_train.drop(['Loan_ID'], axis=1, inplace=True)
2 X_val.drop(['Loan_ID'], axis=1, inplace=True)
3 test.drop(['Loan_ID'], axis=1, inplace=True)
```

Validating Data Types of Features

In [48]:

```
1 # Train
2 X_train['mode__Loan_Amount_Term'] = X_train['mode__Loan_Amount_Term'].astype('int')
3 X_train['mode__Credit_History'] = X_train['mode__Credit_History'].astype('int')
4 X_train['mean__LoanAmount'] = X_train['mean__LoanAmount'].astype('float')
5
6 # Test
7 X_val['mode__Loan_Amount_Term'] = X_val['mode__Loan_Amount_Term'].astype('int')
8 X_val['mode__Credit_History'] = X_val['mode__Credit_History'].astype('int')
9 X_val['mean__LoanAmount'] = X_val['mean__LoanAmount'].astype('float')
10
11 # Val
12 test['mode__Loan_Amount_Term'] = test['mode__Loan_Amount_Term'].astype('int')
13 test['mode__Credit_History'] = test['mode__Credit_History'].astype('int')
14 test['mean__LoanAmount'] = test['mean__LoanAmount'].astype('float')
```

In []:

```
1
```

Separating Numerical Features for Further Analysis

In [49]:

```
1 numerical_features = [ 'ApplicantIncome', 'CoapplicantIncome', 'mean__LoanAmount', 'mode__Loan_Amount_Term', 'mode__Credit_History' ]
```

Detecting Multi-Collinearity

In [50]:

```
1 X_train[numerical_features].corr().round(2)
```

Out[50]:

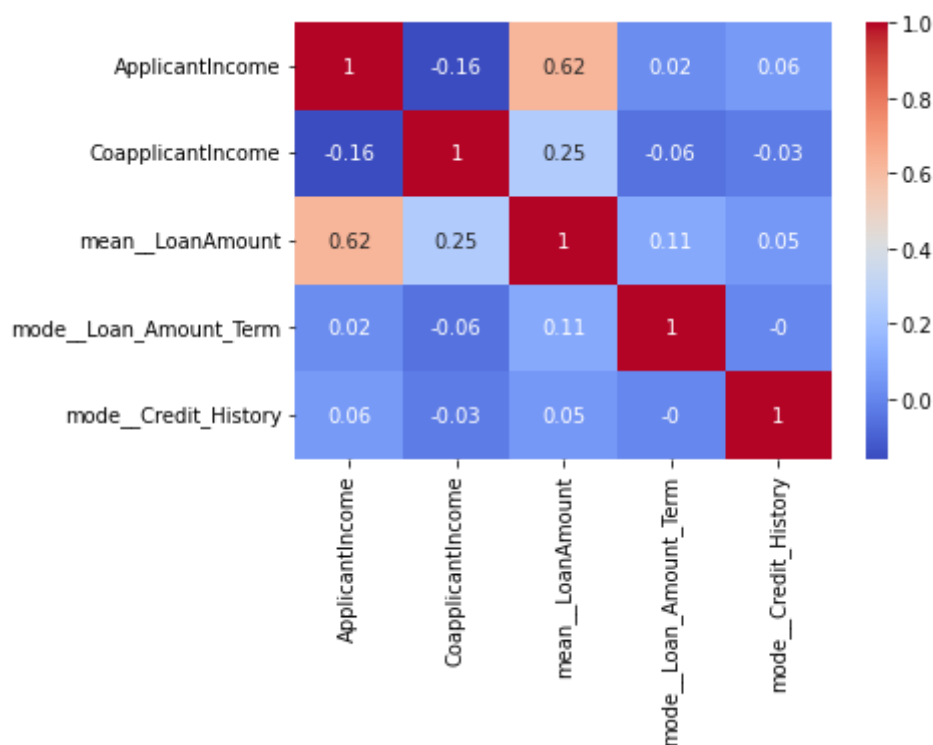
	ApplicantIncome	CoapplicantIncome	mean__LoanAmount	mode__Loan_Amount_Term	mode__Credit_History
ApplicantIncome	1.00	-0.16	0.62	0.02	0.06
CoapplicantIncome	-0.16	1.00	0.25	-0.06	-0.03
mean__LoanAmount	0.62	0.25	1.00	0.11	0.05
mode__Loan_Amount_Term	0.02	-0.06	0.11	1.00	0.05
mode__Credit_History	0.06	-0.03	0.05	0.05	1.00

In [51]:

```
1 sns.heatmap(X_train[numerical_features].corr().round(2), annot=True, cmap='coolwarm')
```

Out[51]:

<AxesSubplot: >



Quick Statistical Summary of Data

In [52]:

```
1 X_train[numerical_features].describe().round(2)
```

Out[52]:

	ApplicantIncome	CoapplicantIncome	mean_LoanAmount	mode_Loan_Amount_Term	mode_Credit_History
count	480.00	480.00	480.00	480.00	480.00
mean	5008.53	1441.32	144.15	342.00	0.00
std	3555.31	1723.58	77.87	66.39	0.00
min	150.00	0.00	9.00	12.00	0.00
25%	2912.75	0.00	100.75	360.00	0.00
50%	3807.00	1128.50	129.50	360.00	0.00
75%	5818.25	2259.25	165.00	360.00	0.00
max	20833.00	8980.00	650.00	480.00	0.00

- verify the minimum values of feature, if they are negative.
- verify if max values of features make sense and are logical
- verify if it contains outliers, an idea could be generated by looking min, max & Mean

Building Base Model

In [53]:

```
1 final_processed_features = [ 'encOH__mode__Gender_Male', 'encOH__mode__Married_Yes',  
2     'encOH__mode__Dependents_1', 'encOH__mode__Dependents_2',  
3     'encOH__mode__Dependents_3+', 'encOH__mode__Self_Employed_Yes',  
4     'encOH__Property_Area_Semiurban', 'encOH__Property_Area_Urban',  
5     'encOR__Education' , 'mode__Loan_Amount_Term', 'mode__Credit_History', 'mean__Lo
```

In [54]:

```
1 ### Let's try out a few more models  
2  
3 from sklearn.preprocessing import StandardScaler  
4 from sklearn.preprocessing import MinMaxScaler  
5  
6 from sklearn import model_selection  
7 from sklearn.metrics import classification_report  
8 from sklearn.metrics import confusion_matrix  
9 from sklearn.metrics import accuracy_score  
10 from sklearn.pipeline import Pipeline  
11 from sklearn.model_selection import GridSearchCV  
12 from sklearn.linear_model import LogisticRegression  
13 from sklearn.tree import DecisionTreeClassifier  
14 from sklearn.neighbors import KNeighborsClassifier  
15 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis  
16 from sklearn.naive_bayes import GaussianNB  
17 from sklearn.svm import SVC  
18  
19 from sklearn.ensemble import AdaBoostClassifier  
20 from sklearn.ensemble import GradientBoostingClassifier  
21 from sklearn.ensemble import RandomForestClassifier  
22
```

In [55]:

```
1 models = []
2
3 models.append(('LR', LogisticRegression(class_weight='balanced', random_state=9)))
4 models.append(('LDA', LinearDiscriminantAnalysis()))
5 models.append(('KNN', KNeighborsClassifier()))
6 models.append(('CART', DecisionTreeClassifier(random_state=9)))
7 models.append(('NB', GaussianNB()))
8 models.append(('SVM', SVC( random_state=9)))
9 models.append(('RF', RandomForestClassifier(random_state=9)))
10 models.append(('GB', GradientBoostingClassifier(random_state=9)))
11 models.append(('AdaB', AdaBoostClassifier(random_state=9)))
12
13 models
```

Out[55]:

```
[('LR', LogisticRegression(class_weight='balanced', random_state=9)),
 ('LDA', LinearDiscriminantAnalysis()),
 ('KNN', KNeighborsClassifier()),
 ('CART', DecisionTreeClassifier(random_state=9)),
 ('NB', GaussianNB()),
 ('SVM', SVC(random_state=9)),
 ('RF', RandomForestClassifier(random_state=9)),
 ('GB', GradientBoostingClassifier(random_state=9)),
 ('AdaB', AdaBoostClassifier(random_state=9))]
```

In [56]:

```
1 train_base_scores = []
2 val_base_scores = []
3 model_names = []
4
5 for name, model in models:
6     model.fit(X_train[final_processed_features], y_train)
7     model_names.append(name)
8     train_base_scores.append(model.score(X_train[final_processed_features], y_train ))
9     val_base_scores.append(model.score(X_val[final_processed_features], y_val ))
```

In [57]:

```

1 for model, train_score, val_score in zip(model_names, train_base_scores, val_base_scores):
2     print('Model -> ', model)
3     print('Train Score -> ', round(train_score,2))
4     print('Val Score -> ', round(val_score,2))
5     print('***30)

```

Model -> NB
Train Score -> 0.82

Val Score -> 0.79

Model -> SVM

Train Score -> 0.69

Val Score -> 0.69

Model -> RF

Train Score -> 0.99

Val Score -> 0.73

Model -> GB

Train Score -> 0.88

Val Score -> 0.79

Model -> AdaB

Train Score -> 0.84

Val Score -> 0.78

```

1 Findings
2 -----
3 Logistic - fine
4 LDA - fine
5 KNN - fine
6 CART - Overfitting, could be used if tuned
7 NB - fine
8 SVM - fine
9 RF - Overfitting, could be used if tuned
10 GB - Overfitting, could be used if tuned
11 AdaB - Overfitting, could be used if tuned
12

```

```

1 NOTE
2 ====
3
4 we simply can't proceed with the model performing best with test data, as it's not
  guaranteed that with new test data, the model performance would be consistent given
  the pattern is same.
5
6 answer is CROSSVALIDATION
7 as it is repeated training and testing, we can see the variation in the models
  performance.
8
9 So we need to check how close are models predicted values among themselves on similar
  dataset.
10 Even if is consistently predictiong wrong, it is highly precise.
11 a model which is sometimes good and sometimes bad, thats not good.
12

```

```
13 THIS IS THE REAL USE CASE OF CROSS VALIDATION
```

In [58]:

```
1 y_train.value_counts(normalize=True)
```

Out[58]:

```
Y    0.691667
```

```
N    0.308333
```

```
Name: Loan_Status, dtype: float64
```


In [59]:

```

1  from sklearn.model_selection import StratifiedKFold, cross_val_score
2
3  results = []
4  model_names = []
5  n_splits= 5
6
7  for name, model in models:
8      kfold = StratifiedKFold(n_splits= 5, shuffle= True, random_state= 9)
9      cv_results = cross_val_score(model, X_train[final_processed_features], y_train, cv=
10     results.append(cv_results)
11     model_names.append(name)
12
13     msg = "%s: %5.2f (%5.2f)" % (name, cv_results.mean()*100, cv_results.std()*100)
14     print(msg)
15
16 results_df = pd.DataFrame(results, index = model_names, columns = 'CV1 CV2 CV3 CV4 CV5
17 results_df['CV Mean'] = results_df.iloc[:, 0:n_splits].mean(axis=1)
18 results_df['CV Std Dev'] = results_df.iloc[:, 0:n_splits].std(axis=1)
19 results_df.sort_values(by= 'CV Mean', ascending= False)*100
20
21
22

```

LR: 75.83 (2.02)
 LDA: 81.88 (1.06)
 KNN: 60.42 (1.98)
 CART: 69.79 (5.02)
 NB: 81.04 (0.78)
 SVM: 69.17 (0.51)
 RF: 79.38 (2.83)
 GB: 77.08 (2.95)
 AdaB: 78.96 (1.21)

Out[59]:

	CV1	CV2	CV3	CV4	CV5	CV Mean	CV Std Dev
LDA	83.333333	80.208333	82.291667	81.250000	82.291667	81.875000	1.187683
NB	82.291667	80.208333	81.250000	80.208333	81.250000	81.041667	0.871521
RF	84.375000	78.125000	80.208333	78.125000	76.041667	79.375000	3.159531
AdaB	78.125000	78.125000	81.250000	79.166667	78.125000	78.958333	1.358167
GB	78.125000	71.875000	80.208333	79.166667	76.041667	77.083333	3.294039
LR	72.916667	77.083333	78.125000	77.083333	73.958333	75.833333	2.258280
CART	62.500000	72.916667	69.791667	77.083333	66.666667	69.791667	5.609547
SVM	69.791667	69.791667	68.750000	68.750000	68.750000	69.166667	0.570544
KNN	60.416667	63.541667	60.416667	60.416667	57.291667	60.416667	2.209709

```

1  Note
2  =====
3
4  We should use whole Training Data for Cross Validation (not X_Train or any splitted
   part.

```

```
5 Reason being, the data is already pslitted in Cross Validation in n folds, and here
also if we are passing Splitted data then the overall % of data used to train
individual models becomes insufficient to get descent scores and understand real
predictive power of our model.
```

Scaling Features

In [60]:

```
1 from sklearn.preprocessing import StandardScaler
2
3 scaler = StandardScaler()
```

In [61]:

```
1 cols_to_scale = ['mode__Loan_Amount_Term', 'mean__LoanAmount']
```

In [62]:

```
1 scaler.fit(X_train[cols_to_scale])
```

Out[62]:

```
▼ StandardScaler
StandardScaler()
```

In [63]:

```
1 X_train[cols_to_scale] = scaler.transform(X_train[cols_to_scale])
2 X_val[cols_to_scale] = scaler.transform(X_val[cols_to_scale])
3 test[cols_to_scale] = scaler.transform(test[cols_to_scale])
```

Re Training & Evaluating Models

In [64]:

```

1  from sklearn.model_selection import StratifiedKFold, cross_val_score
2
3  results = []
4  model_names = []
5  n_splits= 5
6
7  for name, model in models:
8      kfold = StratifiedKFold(n_splits= 5, shuffle= True, random_state= 9)
9      cv_results = cross_val_score(model, X_train[final_processed_features], y_train, cv=
10     results.append(cv_results)
11     model_names.append(name)
12
13     msg = "%s: %5.2f (%5.2f)" % (name, cv_results.mean()*100, cv_results.std()*100)
14     print(msg)
15
16 results_df = pd.DataFrame(results, index = model_names, columns = 'CV1 CV2 CV3 CV4 CV5
17 results_df['CV Mean'] = results_df.iloc[:, 0:n_splits].mean(axis=1)
18 results_df['CV Std Dev'] = results_df.iloc[:, 0:n_splits].std(axis=1)
19 results_df.sort_values(by= 'CV Mean', ascending= False)*100
20
21
22

```

LR: 76.04 (2.47)
 LDA: 81.88 (1.06)
 KNN: 74.17 (2.83)
 CART: 69.79 (5.02)
 NB: 81.04 (0.78)
 SVM: 81.04 (2.02)
 RF: 79.38 (2.83)
 GB: 77.29 (2.59)
 AdaB: 79.17 (1.14)

Out[64]:

	CV1	CV2	CV3	CV4	CV5	CV Mean	CV Std Dev
LDA	83.333333	80.208333	82.291667	81.250000	82.291667	81.875000	1.187683
NB	82.291667	80.208333	81.250000	80.208333	81.250000	81.041667	0.871521
SVM	83.333333	80.208333	83.333333	80.208333	78.125000	81.041667	2.258280
RF	84.375000	78.125000	80.208333	78.125000	76.041667	79.375000	3.159531
AdaB	78.125000	79.166667	81.250000	79.166667	78.125000	79.166667	1.275776
GB	78.125000	72.916667	80.208333	79.166667	76.041667	77.291667	2.890508
LR	71.875000	79.166667	77.083333	77.083333	75.000000	76.041667	2.755991
KNN	71.875000	77.083333	76.041667	76.041667	69.791667	74.166667	3.159531
CART	62.500000	72.916667	69.791667	77.083333	66.666667	69.791667	5.609547

```

1  Naive Bayes is performing better with appropriate Standard Deviation
2

```

```
3 Let's perform Hyperparameter Tuning for it.
```

In [65]:

```
1 from sklearn.model_selection import GridSearchCV
2
3 clf = GaussianNB()
4
5 # define parameter grid for grid search
6 param_grid = {'var_smoothing': [1e-9, 1e-8, 1e-7, 1e-6, 1e-5]}
```

In [66]:

```
1 # perform grid search using 5-fold cross-validation
2 grid_search = GridSearchCV(clf, param_grid=param_grid, cv=5, scoring='accuracy')
3 grid_search.fit(X_train[final_processed_features], y_train)
4
5 # print the best parameters and corresponding accuracy score
6 print("Best parameters: ", grid_search.best_params_)
7 print("Best accuracy score: ", grid_search.best_score_)
```

Best parameters: {'var_smoothing': 1e-09}

Best accuracy score: 0.7979166666666667

In [67]:

```
1 clf = GaussianNB(var_smoothing= 1e-09)
```

Training Final Model

In [68]:

```
1 clf.fit(X_train[final_processed_features], y_train)
```

Out[68]:

```
▼ GaussianNB
GaussianNB()
```

In [69]:

```
1 # Training Score
2 clf.score(X_train[final_processed_features], y_train)
```

Out[69]:

0.8166666666666667

In [70]:

```
1 # Validation Score
2 clf.score(X_val[final_processed_features], y_val)
```

Out[70]:

0.7916666666666666

In [71]:

```

1 #checking predictions son Validation set
2 pred_test = clf.predict(test[final_processed_features])
3 pred_test

```

Out[71]:

```

array(['Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N',
       'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'N',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'N',
       'Y', 'N', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y',
       'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'N', 'N', 'N', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'N', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y',
       'Y', 'Y', 'Y', 'Y', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'N', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'Y',
       'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'N', 'Y', 'N', 'Y', 'Y',
       'N', 'N', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N',
       'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y',
       'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y',
       'Y', 'N', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y',
       'Y', 'Y', 'Y', 'N', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'Y', 'N', 'Y',
       'Y', 'Y', 'Y'], dtype='<U1')

```

In [72]:

```

1 pred_test_df = pd.DataFrame({'Loan_ID' : test_backup['Loan_ID'], 'Loan_Status' : pred_t
2 pred_test_df.head()

```

Out[72]:

	Loan_ID	Loan_Status
0	LP001015	Y
1	LP001022	Y
2	LP001031	Y
3	LP001035	Y
4	LP001051	Y

In [73]:

```
1 # saving File
2 pred_test_df.to_csv("predicted_Loan_Defaulter.csv", index=False)
```

In []:

```
1
```

Trying Lazy Predictor

```
1 !pip install lazypredict
```

In [77]:

```
1 import lazypredict
2 from lazypredict.Supervised import LazyClassifier
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split
```


In []:

1	
---	--