

```
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import import cross_val_score
from sklearn.metrics import average_precision_score, recall_score
from sklearn.metrics import accuracy_score
from sklearn.datasets import make_regression
from sklearn.ensemble import GradientBoostingRegressor
import joblib
import pickle

In [2]:
# Update sklearn to prevent version mismatches
!pip install sklearn --upgrade

Requirement already up-to-date: sklearn in c:\users\glori\anaconda3\lib\site-packages (0.0)
Requirement already satisfied, skipping upgrade: scikit-learn in c:\users\glori\anaconda3\lib\site-packages (from sklearn) (0.23.2)
Requirement already satisfied, skipping upgrade: numpy>=1.13.3 in c:\users\glori\anaconda3\lib\site-packages (from scikit-learn->sklearn) (1.19.2)
Requirement already satisfied, skipping upgrade: threadpoolctl>=2.0.0 in c:\users\glori\anaconda3\lib\site-packages (from scikit-learn->sklearn) (2.1.0)
Requirement already satisfied, skipping upgrade: scipy>=0.19.1 in c:\users\glori\anaconda3\lib\site-packages (from scikit-learn->sklearn) (1.5.2)
Requirement already satisfied, skipping upgrade: joblib>=0.11 in c:\users\glori\anaconda3\lib\site-packages (from scikit-learn->sklearn) (0.17.0)

In [3]:
# Install joblib. This will be used to save your model.
# Restart your kernel after installing
!pip install joblib

Requirement already satisfied: joblib in c:\users\glori\anaconda3\lib\site-packages (0.17.0)

In [4]:
import re
import os

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import datetime as dt

In [5]:
file_path = "Data/model_1_combine_data.csv"

In [6]:
my_data = pd.read_csv(file_path, error_bad_lines=False, low_memory=False)
my_data

Out[6]:
   Unnamed: 0  Amount_Requested  Risk_Score  debt_to_income_ratio  State  Employment_Length
0      421097             5000.0         669.0                2180    OK                8
1      421098             15000.0         704.0                1829    FL                2  debt_
2      421099             11200.0         669.0                4397    NH                0
3      421100             25000.0         669.0                1289    AL                10  debt_
4      421101             3000.0         764.0                0.58    WA                9  ma
...         ...                ...         ...                ...    ...                ...
1081515  19698926             3000.0         535.0                1.54    CT                5  Med
1081516  19698936             6000.0         580.0                5.55    TX                0
1081517  19698996             3000.0         561.0                0.00    OR                0  Debt
1081518  19699052             2500.0         573.0                8.55    NY                0
1081519  19699073             2500.0         567.0                0.00    AL                2

1081520 rows x 8 columns

In [7]:
my_data['Loan_Status'].value_counts()

Out[7]:
Rejected    628344
Accepted    453176
Name: Loan_Status, dtype: int64

In [8]:
# Understanding the split the reject and accept applications
print(my_data['Loan_Status'].value_counts(normalize=True))

Rejected    0.580982
Accepted    0.419018
Name: Loan_Status, dtype: float64

In [9]:
#predicting loan status creting countplot
sns.countplot(x='Loan_Status', data=my_data)

Out[9]:
<AxesSubplot:xlabel='Loan_Status', ylabel='count'>



In [10]:
#histogram of the loan amount
plt.figure(figsize=(12,4))
sns.distplot(my_data['Amount_Requested'], kde=False, bins=40)
plt.xlim(0,45000)

C:\Users\glori\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning:
'sdistplot' is a deprecated function and will be removed in a future version. Please
adapt your code to use either 'distplot' (a figure-level function with similar flexibili
ty) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[10]:
(0.0, 45000.0)



In [11]:
#what is the correlation b/t the variables
my_data.corr()

Out[11]:
   Unnamed: 0  Amount_Requested  Risk_Score  debt_to_income_ratio  Employment_Length
0      1.000000             -0.236809             -0.561502             0.019285             -0.618999
Amount_Requested             -0.236809             1.000000             0.310587             0.002001             0.190206
Risk_Score                  -0.561502             0.310587             1.000000             -0.002608             0.449296
debt_to_income_ratio         0.019285             0.002001             -0.002608             1.000000             -0.021967
Employment_Length           -0.618999             0.190206             0.449296             -0.021967             1.000000

In [12]:
#correlation using heatmap
plt.figure(figsize=(30, 6))
sns.heatmap(my_data.corr(), annot=True, cmap='viridis')
plt.xlim(10,0)

Out[12]:
(10.0, 0.0)



In [13]:
#boxplot showing relationship between the loan amount and loan status
sns.boxplot(x='Loan_Status', y='Amount_Requested')

Out[13]:
<AxesSubplot:xlabel='Loan_Status', ylabel='Amount_Requested'>



In [14]:
#loan Amount, grouped by Loan Status
my_data.groupby('Loan_Status')['Amount_Requested'].describe()

Out[14]:
   Loan_Status  count      mean      std      min      25%      50%      75%      max
Accepted    453176  16304.636995  10200.954582  1000.0  8500.0  14400.0  22825.0  40000.0
Rejected    628344  11385.468589  10647.183241  525.0  3000.0  8000.0  15500.0  90725.0

In [15]:
#split of all categorical columns
my_data['State'].value_counts().plot.bar(figsize=(20,5), title='State')
plt.show()
my_data['Purpose'].value_counts().plot.bar(figsize=(20,5), title='Purpose of Loan')
plt.show()

State
Debt consolidation  10000
Car financing      9000
Credit card refinancing  8000
Other              7000
Home improvement   6000
Medical expenses   5000
Moving and relocat  4000
Major purchase     3000
Vacation           2000
Business           1000
Wedding            500
Other              500
Purpose of Loan
Debt consolidation  10000
Car financing      9000
Credit card refinancing  8000
Other              7000
Home improvement   6000
Medical expenses   5000
Moving and relocat  4000
Major purchase     3000
Vacation           2000
Business           1000
Wedding            500
Other              500

In [16]:
#split of Numeric columns
sns.distplot(my_data['Amount_Requested'])
plt.show()

C:\Users\glori\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning:
'sdistplot' is a deprecated function and will be removed in a future version. Please
adapt your code to use either 'distplot' (a figure-level function with similar flexibili
ty) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[16]:
(0.0, 45000.0)



In [17]:
#split of numeric columns
sns.distplot(my_data['Risk_Score'])
plt.show()

C:\Users\glori\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning:
'sdistplot' is a deprecated function and will be removed in a future version. Please
adapt your code to use either 'distplot' (a figure-level function with similar flexibili
ty) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[17]:
(300, 1000)



In [18]:
# understanding split of all Numeric Columns
sns.distplot(my_data['debt_to_income_ratio'])
plt.show()

print(my_data['debt_to_income_ratio'].max())
print('-----')
print(my_data['debt_to_income_ratio'].min())
print('-----')
print(my_data[['debt_to_income_ratio', 'Loan_Status']].sort_values(by='debt_to_income_

C:\Users\glori\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning:
'sdistplot' is a deprecated function and will be removed in a future version. Please
adapt your code to use either 'distplot' (a figure-level function with similar flexibili
ty) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[18]:
(0.0, 70000.0)



In [19]:
#define categorical and numerical features
categorical_features = ['State', 'Purpose', 'Loan_Status']
numerical_features = ['Amount_Requested', 'Risk_Score', 'debt_to_income_ratio', 'Empl

features = categorical_features + numerical_features

#define the features
print('Categorical Features:\t', len(categorical_features))
print('Numerical Features:\t', len(numerical_features))
print('Total Features:\t\t', len(features))

Categorical Feature: 3
Numerical Feature: 4
Total Features: 7

In [20]:
my_data

Out[20]:
   Unnamed: 0  Amount_Requested  Risk_Score  debt_to_income_ratio  State  Employment_Length
0      421097             5000.0         669.0                2180    OK                8
1      421098             15000.0         704.0                1829    FL                2  debt_
2      421099             11200.0         669.0                4397    NH                0
3      421100             25000.0         669.0                1289    AL                10  debt_
4      421101             3000.0         764.0                0.58    WA                9  ma
...         ...                ...         ...                ...    ...                ...
1081515  19698926             3000.0         535.0                1.54    CT                5  Med
1081516  19698936             6000.0         580.0                5.55    TX                0
1081517  19698996             3000.0         561.0                0.00    OR                0  Debt
1081518  19699052             2500.0         573.0                8.55    NY                0
1081519  19699073             2500.0         567.0                0.00    AL                2

1081520 rows x 8 columns

In [21]:
loan_df=my_data.drop(['Unnamed: 0'], axis=1).fillna(-1)
loan_df

Out[21]:
   Amount_Requested  Risk_Score  debt_to_income_ratio  State  Employment_Length  Purpose
0      5000.0         669.0                2180    OK                8      other
1     15000.0         704.0                1829    FL                2  debt_consolidation
2     11200.0         669.0                4397    NH                0      medical
3     25000.0         669.0                1289    AL                10  debt_consolidation
4      3000.0         764.0                0.58    WA                9  major_purchase
...         ...                ...         ...                ...    ...                ...
1081515      3000.0         535.0                1.54    CT                5  Medical expenses
1081516      6000.0         580.0                5.55    TX                0    Car financing
1081517      3000.0         561.0                0.00    OR                0  Debt consolidation
1081518      2500.0         573.0                8.55    NY                0    Other
1081519      2500.0         567.0                0.00    AL                2    Car financing

1081520 rows x 7 columns

In [22]:
# df[df.b > 10]
dti_values_over100 = loan_df[loan_df.debt_to_income_ratio > 100]

In [23]:
print(dti_values_over100.info())

<class 'pandas.core.frame.DataFrame'>
Int64Index: 22081 entries, 493 to 1081423
Data columns (total 7 columns):
 #   Column                Non-Null Count  Dtype
---  ---
 0   Amount_Requested      22081 non-null  float64
 1   Risk_Score            22081 non-null  float64
 2   debt_to_income_ratio  22081 non-null  float64
 3   State                 22081 non-null  object
 4   Employment_Length     22081 non-null  object
 5   Purpose               22081 non-null  object
 6   Loan_Status           22081 non-null  object
dtypes: float64(3), int64(1), object(3)
memory usage: 1.3+ MB
None

In [24]:
dti_values_over100.head()

Out[24]:
   Amount_Requested  Risk_Score  debt_to_income_ratio  State  Employment_Length  Purpose
493      28000.0         704.0                112.20    NV                0  debt_consolidation
643      40000.0         709.0                101.46    OH                4  debt_consolidation
1229     10000.0         714.0                114.27    TX                9  debt_consolidation
1683      9000.0         684.0                106.66    IL                4  credit_card
1879     35000.0         714.0                102.47    UT                0  debt_consolidation

In [25]:
dti_values_over100['Loan_Status'].value_counts()

Out[25]:
Rejected    2193
Accepted    888
Name: Loan_Status, dtype: int64

In [26]:
clean_dti_data=loan_df[loan_df.debt_to_income_ratio <= 100]

In [27]:
clean_dti_data

Out[27]:
   Amount_Requested  Risk_Score  debt_to_income_ratio  State  Employment_Length  Purpose
0      5000.0         669.0                2180    OK                8      other
1     15000.0         704.0                1829    FL                2  debt_consolidation
2     11200.0         669.0                4397    NH                0      medical
3     25000.0         669.0                1289    AL                10  debt_consolidation
4      3000.0         764.0                0.58    WA                9  major_purchase
...         ...                ...         ...                ...    ...                ...
1081515      3000.0         535.0                1.54    CT                5  Medical expenses
1081516      6000.0         580.0                5.55    TX                0    Car financing
1081517      3000.0         561.0                0.00    OR                0  Debt consolidation
1081518      2500.0         573.0                8.55    NY                0    Other
1081519      2500.0         567.0                0.00    AL                2    Car financing

1034236 rows x 7 columns

In [28]:
sns.distplot(clean_dti_data['debt_to_income_ratio'])
plt.show()



In [29]:
#plotting the clean data
sns.countplot(clean_dti_data['Loan_Status'], value_counts(normalize=True))

Out[29]:
Rejected    0.562685
Accepted    0.437315
Name: Loan_Status, dtype: float64

In [30]:
#Dependent Variable - Loan_Status
sns.countplot(x='Loan_Status', data=clean_dti_data)

Out[30]:
<AxesSubplot:xlabel='Loan_Status', ylabel='count'>



In [31]:
clean_dti_data['State'].value_counts().plot.bar(figsize=(20,5), title='State')
plt.show()
clean_dti_data['Purpose'].value_counts().plot.bar(figsize=(20,5), title='Purpose of Loan')
plt.show()

State
Debt consolidation  10000
Car financing      9000
Credit card refinancing  8000
Other              7000
Home improvement   6000
Medical expenses   5000
Moving and relocat  4000
Major purchase     3000
Vacation           2000
Business           1000
Wedding            500
Other              500
Purpose of Loan
Debt consolidation  10000
Car financing      9000
Credit card refinancing  8000
Other              7000
Home improvement   6000
Medical expenses   5000
Moving and relocat  4000
Major purchase     3000
Vacation           2000
Business           1000
Wedding            500
Other              500

In [32]:
sns.distplot(clean_dti_data['Amount_Requested'])
plt.show()

C:\Users\glori\anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning:
'sdistplot' is a deprecated function and will be removed in a future version. Please
adapt your code to use either 'distplot' (a figure-level function with similar flexibili
ty) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[32]:
(0.0, 45000.0)



In [33]:
sns.distplot(clean_dti_data['Risk_Score'])
plt.show()



In [34]:
#Loan purpose recategorized
clean_dti_data.loc[clean_dti_data.Purpose == "debt_consolidation", "Purpose"] = "Debt consolidation"
clean_dti_data.loc[clean_dti_data.Purpose == "credit_card_refinancing", "Purpose"] = "Credit card refinancing"
clean_dti_data.loc[clean_dti_data.Purpose == "car_finance", "Purpose"] = "Car Financing"
clean_dti_data.loc[clean_dti_data.Purpose == "home_improvement", "Purpose"] = "Home Improvement"
clean_dti_data.loc[clean_dti_data.Purpose == "moving", "Purpose"] = "Relocation Expense"
clean_dti_data.loc[clean_dti_data.Purpose == "major_purchase", "Purpose"] = "Major purchase"
clean_dti_data.loc[clean_dti_data.Purpose == "medical", "Purpose"] = "Medical Expense"
clean_dti_data.loc[clean_dti_data.Purpose == "house", "Purpose"] = "Home Buying"
clean_dti_data.loc[clean_dti_data.Purpose == "vacation", "Purpose"] = "Vacation"
clean_dti_data.loc[clean_dti_data.Purpose == "small_business", "Purpose"] = "Business"
clean_dti_data.loc[clean_dti_data.Purpose == "renewable_energy", "Purpose"] = "Energy"

C:\Users\glori\anaconda3\lib\site-packages\pandas\core\indexing.py:1765: SettingWithCopyWarning:
A value is being assigned to a slice of a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/10min.html#returning-a-view-versus-a-copy
r_guide/indexing.html#returning-a-view-versus-a-copy
isetter(loc, value)

In [35]:
#Cleaned Purpose Column
clean_dti_data['Purpose'].value_counts()

Out[35]:
Debt consolidation    463923
credit_card_refinanci  179281
Credit card refinanci  118031
Car financing         77889
Home Improvement      31333
Home improvement       29363
Home buying           19056
Medical expenses      18315
Moving and relocation  17612
Major purchase        16669
Home improvement       13046
Business              13010
Major purchase        10754
Vacation              7780
Medical Expense       5783
Home Buying           4921
Car financing         4546
Relocation Expense    2825
Green loan            808
Energy loan            252
Learning and training  34
wedding               5
Name: Purpose, dtype: int64

In [36]:
# plot the heatmap
corr = clean_dti_data.corr()
sns.heatmap(corr, xticklabels=corr.columns, yticklabels=corr.columns, annot=True, cma

Out[36]:
<AxesSubplot>



In [37]:
clean_dti_data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1034236 entries, 0 to 1081519
Data columns (total 7 columns):
 #   Column                Non-Null Count  Dtype
---  ---
 0   Amount_Requested      1034236 non-null  float64
 1   Risk_Score            1034236 non-null  float64
 2   debt_to_income_ratio  1034236 non-null  float64
 3   State                 1034236 non-null  object
 4   Employment_Length     1034236 non-null  object
 5   Purpose               1034236 non-null  object
 6   Loan_Status           1034236 non-null  object
dtypes: float64(3), int64(1), object(3)
memory usage: 63.1+ MB

In [38]:
Loan_Purpose = pd.crosstab(clean_dti_data['Purpose'], clean_dti_data['Loan_Status'])
Employment_Length = pd.crosstab(clean_dti_data['Employment_Length'], clean_dti_data['Loan_Status'])
State = pd.crosstab(clean_dti_data['State'], clean_dti_data['Loan_Status'])

Loan_Purpose.div(Loan_Purpose.sum(1).astype(float), axis = 0).plot(kind='bar', stacked=True)
Employment_Length.div(Employment_Length.sum(1).astype(float), axis = 0).plot(kind='bar', stacked=True)
State.div(State.sum(1).astype(float), axis = 0).plot(kind='bar', stacked=True, figsize=(10, 10))



In [39]:
clean_dti_data.groupby('Loan_Status')['Amount_Requested'].mean().plot.bar()

bins=(0,10000,20000,30000)
group=['Low', 'Average', 'High']
clean_dti_data['Amount_Requested_bins']=pd.cut(clean_dti_data['Amount_Requested'], bins=bins)
Amount_Requested_bins=pd.crosstab(clean_dti_data['Amount_Requested_bins'], clean_dti_data['Loan_Status'])
Amount_Requested_bins.div(Amount_Requested_bins.sum(1).astype(float), axis=0).plot(kind='bar', stacked=True, label='Loan Amount')
plt.ylabel('Percentage')

<ipython-input-39-70b0929e1dd6>:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/10min.html#returning-a-view-versus-a-copy
clean_dti_data['Risk_Score_bins']=pd.cut(clean_dti_data['Risk_Score'], bins, labels=gr

Out[39]:
Text(0, 0.5, 'Percentage')



In [40]:
#plot risk scores
bins=(0,629,689,719,850)
group=['Bad', 'Fair', 'Good', 'Excellent']
clean_dti_data['Risk_Score_bins']=pd.cut(clean_dti_data['Risk_Score'], bins, labels=group)
Risk_Score_bins=pd.crosstab(clean_dti_data['Risk_Score_bins'], clean_dti_data['Loan_Status'])
Risk_Score_bins.div(Risk_Score_bins.sum(1).astype(float), axis=0).plot(kind='bar', stacked=True, label='Risk Score')
plt.ylabel('Percentage')

<ipython-input-40-fdf65cc71f5>:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/10min.html#returning-a-view-versus-a-copy
clean_dti_data['Risk_Score_bins']=pd.cut(clean_dti_data['Risk_Score'], bins, labels=gr

Out[40]:
Text(0, 0.5, 'Percentage')


```


