



## Assignment 2

1. Classify the email using the binary classification method. Email Spam detection has two states: a) Normal State – Not Spam, b) Abnormal State – Spam.  
Use K-Nearest Neighbors and Support Vector Machine for classification. Analyze their performance. Dataset link: The emails.csv dataset on the Kaggle  
<https://www.kaggle.com/datasets/balaka18/email-spam-classification-dataset-csv>

```
In [19]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn import metrics

In [20]: df=pd.read_csv('emails.csv')

In [21]: df.head()

Out[21]:
   Email No. the to ect and for of a you hou ... connevey jay valued lay infrastructure military allowing ff dry Prediction
0 Email 1 0 0 1 0 0 0 2 0 0 ... 0 0 0 0 0 0 0 0 0 0
1 Email 2 8 13 24 6 6 2 102 1 27 ... 0 0 0 0 0 0 0 0 1 0 0
2 Email 3 0 0 1 0 0 0 8 0 0 ... 0 0 0 0 0 0 0 0 0 0 0
3 Email 4 0 5 22 0 5 1 51 2 10 ... 0 0 0 0 0 0 0 0 0 0 0
4 Email 5 7 6 17 1 5 2 57 0 9 ... 0 0 0 0 0 0 0 0 1 0 0

5 rows × 3002 columns

In [22]: df.columns

Out[22]: Index(['Email No.', 'the', 'to', 'ect', 'and', 'for', 'of', 'a', 'you', 'hou',
       ...
       'connevey', 'jay', 'valued', 'lay', 'infrastructure', 'military',
       'allowing', 'ff', 'dry', 'Prediction'],
       dtype='object', length=3002)

In [23]: df.isnull().sum()

Out[23]:
Email No.      0
the            0
to             0
ect            0
and            0
...
military       0
allowing       0
ff             0
dry            0
Prediction     0
Length: 3002, dtype: int64

In [24]: df.dropna(inplace = True)

In [25]: df.drop(['Email No.'],axis=1,inplace=True)
X = df.drop(['Prediction'],axis = 1)
y = df['Prediction']

In [26]: from sklearn.preprocessing import scale
X = scale(X)
# split into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 42)
```

### KNN classifier

```
In [35]: from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=7)

knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)

In [36]: print("Prediction",y_pred)
Prediction [0 0 1 ... 1 1 1]

In [37]: print("KNN accuracy = ",metrics.accuracy_score(y_test,y_pred))
KNN accuracy =  0.8009020618556701

In [39]: print("Confusion matrix",metrics.confusion_matrix(y_test,y_pred))
Confusion matrix [[804 293]
 [ 16 439]]
```

### SVM classifier

```
In [27]: # cost C = 1
model = SVC(C = 1)

# fit
model.fit(X_train, y_train)

# predict
y_pred = model.predict(X_test)

In [28]: metrics.confusion_matrix(y_true=y_test, y_pred=y_pred)

Out[28]: array([[1091,    6],
       [ 90, 365]])

In [29]: print("SVM accuracy = ",metrics.accuracy_score(y_test,y_pred))

SVM accuracy =  0.9381443298969072
```

## Given a bank customer, build a neural network-based classifier that can determine whether they will leave or not in the next 6 months.

Dataset Description: The case study is from an open-source dataset from Kaggle. The dataset contains 10,000 sample points with 14 distinct features such as CustomerId, CreditScore, Geography, Gender, Age, Tenure, Balance, etc. Link to the Kaggle project: <https://www.kaggle.com/barelydedicated/bank-customer-churn-modelling>. Perform following steps:

1. Read the dataset.
2. Distinguish feature and target set and divide the data set into training and test sets.
3. Normalize the column and test data.
4. Initialize and build the model. Identify the points of improvement and implement the same.
5. Print the accuracy score and confusion matrix.

```
In [46]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt #Importing the libraries
```

```
In [47]: df = pd.read_csv("Churn_Modelling.csv")
```

### Preprocessing.

```
In [48]: df.head()
```

```
Out[48]:
```

RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	
0	1	Herrgave	619	France	Female	42	2	0.00	1	1	1	101348.88	
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	1	0	1	112542.58
2	3	15019304	Onio	502	France	Female	42	8	159660.80	3	1	0	113931.57
3	4	15701954	Boni	699	France	Female	39	1	0.00	2	0	0	93826.63
4	5	15737688	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10

```
In [49]: df.shape
```

```
Out[49]:
```

(10000, 14)

```
In [50]: df.describe()
```

```
Out[50]:
```

RowNumber	CustomerId	CreditScore	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
count	10000.000000	1.000000e+04	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.500000	1.56904e+03	38.921800	2.892174	63485.889298	1.532000	0.70550	0.515100	10090.239881
std	2886.95658	7.193619e+04	501652399	10.487806	581654	0.45584	0.499797	5710.492818	
min	1.000000	1.55657e+07	350.000000	18.000000	0.000000	0.000000	1.000000	0.000000	11.580000
25%	2500.500000	1.56285e+07	658400000	32.000000	3.000000	1.000000	0.000000	5102.110000	
50%	5000.500000	1.56904e+07	652000000	37.000000	5.000000	97198.540000	1.000000	1.000000	100193.915000
75%	7500.500000	1.57532e+07	718.000000	44.000000	7.000000	127644.240000	2.000000	1.000000	149388.247500
max	10000.000000	1.58156e+07	850.000000	92.000000	10.000000	25988.090000	4.000000	1.000000	19992.480000

```
In [51]: df.isnull()
```

```
Out[51]:
```

RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	False	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...	...	...	...	...
9995	False	False	False	False	False	False	False	False	False	False	False	False
9996	False	False	False	False	False	False	False	False	False	False	False	False
9997	False	False	False	False	False	False	False	False	False	False	False	False
9998	False	False	False	False	False	False	False	False	False	False	False	False
9999	False	False	False	False	False	False	False	False	False	False	False	False

10000 rows x 14 columns

```
In [52]: df.isnull().sum()
```

```
Out[52]:
```

RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...	...	...	...
9995	0	0	0	0	0	0	0	0	0	0	0	0
9996	0	0	0	0	0	0	0	0	0	0	0	0
9997	0	0	0	0	0	0	0	0	0	0	0	0
9998	0	0	0	0	0	0	0	0	0	0	0	0
9999	0	0	0	0	0	0	0	0	0	0	0	0

```
In [53]: df.info()
```

```
Out[53]:
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 10000 entries, 0 to 9999  
Data columns (total 14 columns):  
 # Column Non-Null Count Dtype  
---  
 0 RowNumber 10000 non-null int64  
 1 CustomerId 10000 non-null int64  
 2 Surname 10000 non-null object  
 3 CreditScore 10000 non-null int64  
 4 Geography 10000 non-null object  
 5 Gender 10000 non-null object  
 6 Age 10000 non-null int64  
 7 Tenure 10000 non-null int64  
 8 Balance 10000 non-null float64  
 9 NumOfProducts 10000 non-null int64  
 10 HasCrCard 10000 non-null int64  
 11 IsActiveMember 10000 non-null int64  
 12 EstimatedSalary 10000 non-null float64  
 13Exited 10000 non-null int64  
 dtypes: float64(4), int64(9), object(3)  
 memory usage: 1.14+ MB

```
In [54]: df.dtypes
```

```
Out[54]:
```

RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
int64	int64	object	int64	object	object	int64	int64	float64	int64	int64	int64	float64

```
In [55]: df.columns
```

```
Out[55]:
```

Index(['RowNumber', 'CustomerId', 'Surname', 'CreditScore', 'Geography', 'Gender', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary', 'Exited'], dtype='object')

```
In [56]: df = df.drop(['RowNumber', 'Surname', 'CustomerId'], axis=1) #Dropping the unnecessary columns
```

```
In [57]: df.head()
```

```
Out[57]:
```

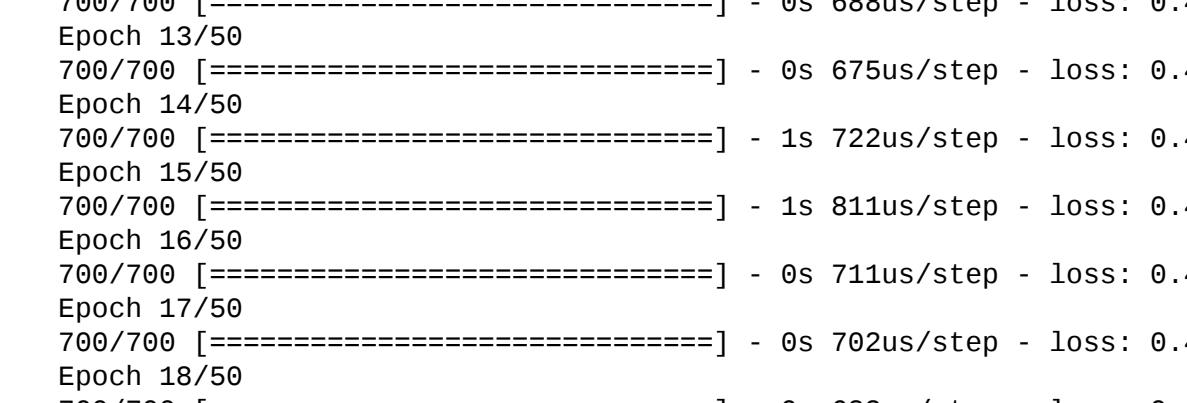
CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	France	Female	42	2	0.00	1	1	101348.88	1
1	608	Spain	Female	41	1	83807.86	1	0	112542.58	0
2	502	France	Female	42	8	159660.80	3	1	113931.57	1
3	699	France	Female	39	1	0.00	2	0	93826.63	0
4	850	Spain	Female	43	2	125510.82	1	1	79084.10	0

### Visualization

```
In [101]: def visualization(x, y, xlabel):  
    plt.figure(figsize=(10,5))  
    plt.hist([x,y], colors=['red', 'green'], label = ['exit', 'not_exit'])  
    plt.xlabel(xlabel, fontsize=20)  
    plt.ylabel("No. of customers", fontsize=20)  
    plt.legend()
```

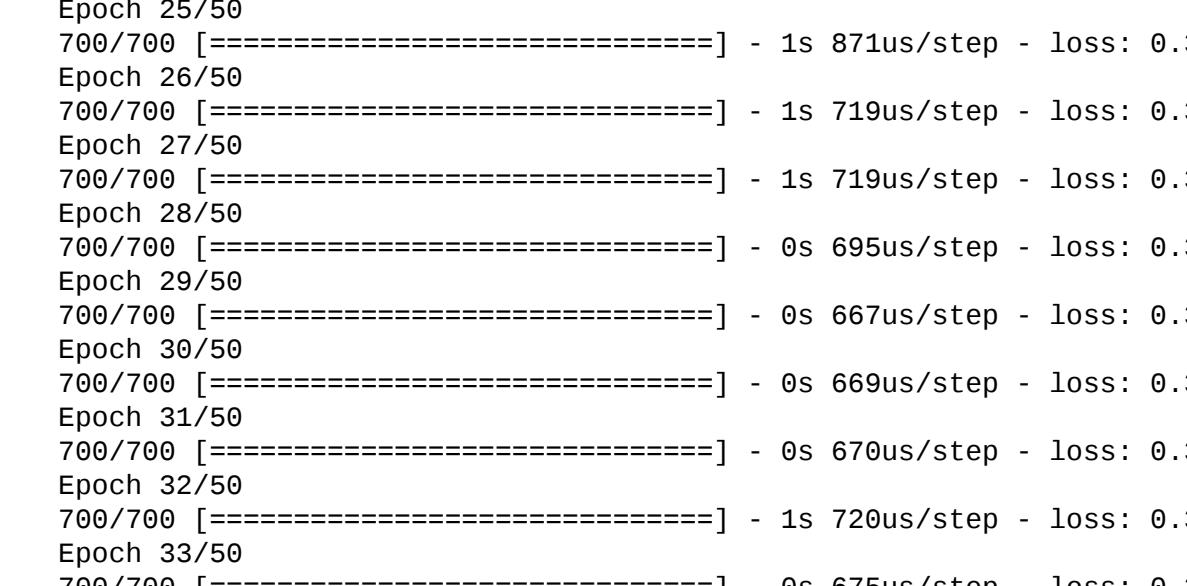
```
In [102]: df_churn_exited = df[df['Exited']==1]['Tenure']  
df_churn_not_exited = df[df['Exited']==0]['Tenure']
```

```
In [103]: visualization(df_churn_exited, df_churn_not_exited, "Tenure")
```



```
In [105]: df_churn_exited2 = df[df['Exited']==1]['Age']  
df_churn_not_exited2 = df[df['Exited']==0]['Age']
```

```
In [106]: visualization(df_churn_exited2, df_churn_not_exited2, "Age")
```



### Converting the Categorical Variables

```
In [59]: X = df[['CreditScore', 'Gender', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary', 'y']]  
states = pd.get_dummies(df['Geography'], drop_first = True)  
gender = pd.get_dummies(df['Gender'], drop_first = True)
```

```
In [61]: df = pd.concat([df, gender, states], axis = 1)
```

### Splitting the training and testing Dataset

```
In [62]: df.head()
```

```
Out[62]:
```

CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	France	Female	42	2	0.00	1	1	101348.88	1
1	608	Spain	Female	41	1	83807.86	1	0	112542.58	0
2	502	France	Female	42	8	159660.80	3	1	113931.57	1
3	699	France	Female	39	1	0.00	2	0	93826.63	0
4	850	Spain	Female	43	2	125510.82	1	1	79084.10	0

```
In [63]: X = df[['CreditScore', 'Age', 'Tenure
```

**Prcatical No. 4: Manual Content**



**Guru Gobind Singh Foundation's  
Guru Gobind Singh College of Engineering and  
Research Center, Nashik**



**Experiment No: 04**

**Title of Experiment:** Implement Gradient Descent Algorithm to find the local minima of a function. For example, find the local minima of the function  $y=(x+3)^2$  starting from the point  $x=2$ .

<b>Student Name:</b>				
<b>Class:</b>	<b>BE (Computer)</b>			
<b>Div:</b>	-	<b>Batch:</b>	<b>CO</b>	
<b>Roll No.:</b>				
<b>Date of Attendance (Performance):</b>				
<b>Date of Evaluation:</b>				
<b>Marks (Grade) Attainment of CO Marks out of 10</b>	<b>A</b>	<b>P</b>	<b>W</b>	<b>T</b>
<b>CO Mapped</b>	CO3: Select and apply appropriately supervised machine learning algorithms for real time applications.			
<b>Signature of Subject Teacher</b>				

**TITLE:** Implement Gradient Descent Algorithm to find the local minima of a function. For example, find the local minima of the function  $y=(x+3)^2$  starting from the point  $x=2$ .

**AIM:** The aim of this practical is to illustrate the implementation of the Gradient Descent algorithm for locating local minima of a given function. This involves understanding the gradient descent method, performing iterations to converge towards a local minimum, and evaluating the effectiveness of the algorithm in the context of the function  $y = (x + 3)^2$  starting from the initial point ( $x = 2$ ).

**OBJECTIVES:** Based on above main aim following are the objectives

1. **Understanding Optimization Techniques:** Develop a solid understanding of optimization methods, with a focus on the Gradient Descent algorithm. Grasp the concept of finding local minima in mathematical functions through iterative updates.
2. **Algorithm Implementation:** Implement the Gradient Descent algorithm practically. Apply it to the function  $y = (x + 3)^2$ , beginning with an initial value of ( $x = 2$ ), and observe how the algorithm converges towards a local minimum.
3. **Analyzing Convergence:** Examine the iterative behavior of the Gradient Descent algorithm. Gain insights into how the algorithm iteratively refines the value of ( $x$ ) by utilizing the gradient of the function.
4. **Algorithm Evaluation:** Evaluate the effectiveness of the Gradient Descent algorithm in successfully identifying the local minimum of the given function. Compare the obtained result with the anticipated theoretical solution.

**Prerequisite:**

Basic of Python, Concept of Gradient Descent Algorithm

**Software Requirements:**

Anaconda with Python 3.7

**Hardware Requirement:**

PIV, 2GB RAM, 500 GB HDD

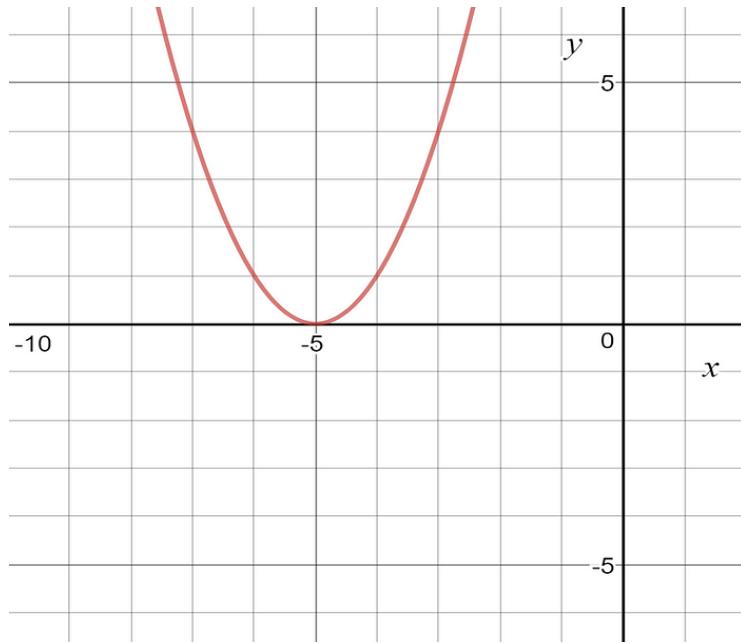
## Theory Concepts:

What is gradient descent ?

It is an optimization algorithm to find the minimum of a function. We start with a random point on the function and move in the **negative direction** of the **gradient of the function** to reach the **local/global minima**.

Example:

**Question :** Find the local minima of the function  $y=(x+5)^2$  starting from the point  $x=3$



**Solution :** We know the answer just by looking at the graph.  $y = (x+5)^2$  reaches it's minimum value when  $x = -5$  (i.e when  $x=-5$ ,  $y=0$ ). Hence  $x=-5$  is the local and global minima of the function.

Now, let's see how to obtain the same numerically using gradient descent.

**Step 1 :** Initialize  $x = 3$ . Then, find the gradient of the function,  $dy/dx = 2*(x+5)$ .

**Step 2 :** Move in the direction of the negative of the gradient ([Why?](#)). But wait, how much to move?

For that, we require a learning rate. Let us assume the **learning rate  $\rightarrow 0.01$**

**Step 3 :** Let's perform 2 iterations of gradient descent

### Initialize Parameters :

$$X_0 = 3$$

$$\text{Learning rate} = 0.01$$

$$\frac{dy}{dx} = \frac{d}{dx} (x + 5)^2 = 2 * (x + 5)$$

### Iteration 1 :

$$X_1 = X_0 - (\text{learning rate}) * \left(\frac{dy}{dx}\right)$$

$$X_1 = 3 - (0.01) * (2 * (3 + 5)) = 2.84$$

### Iteration 2 :

$$X_2 = X_1 - (\text{learning rate}) * \left(\frac{dy}{dx}\right)$$

$$X_2 = 2.84 - (0.01) * (2 * (2.84 + 5)) = 2.6832$$

Step 4 : We can observe that the X value is slowly decreasing and should converge to -5 (the local minima).

## Implementing Gradient Descent in Python

Importing libraries

```
import numpy as np
import matplotlib.pyplot as plt
```

We then define the function  $f(x)$  and its derivative  $f'(x)$  –

```
def f(x):
    return (x+3)**2

def df(x):
    return 2*x + 6
```

$F(x)$  is the function that has to be decreased, and  $df$  is its derivative ( $x$ ). The gradient descent approach makes use of the derivative to guide itself toward the minimum by revealing the function's slope along the way.

The gradient descent function is then defined.

```
def gradient_descent(initial_x, learning_rate, num_iterations):
    x = initial_x
    x_history = [x]

    for i in range(num_iterations):
        gradient = df(x)
        x = x - learning_rate * gradient
        x_history.append(x)

    return x, x_history
```

The starting value of  $x$ , the learning rate, and the desired number of iterations are sent to the gradient descent function. In order to save the values of  $x$  after each iteration, it initializes  $x$  to its original value and generates an empty list. The method then executes the gradient descent for the provided number of iterations, changing  $x$  every iteration in accordance with the equation  $x = x - \text{learning rate} * \text{gradient}$ . The function produces a list of every iteration's  $x$  values together with the final value of  $x$ .

The gradient descent function may now be used to locate the local minimum of  $f(x)$  –

```
initial_x = 2
learning_rate = 0.1
num_iterations = 50

x, x_history = gradient_descent(initial_x, learning_rate, num_iterations)

print("Local minimum: {:.2f}".format(x))
```

## Output

```
Local minimum: -3.00
```

In this illustration,  $x$  is set to 2 at the beginning, with a learning rate of 0.1, and 50 iterations are run. Finally, we publish the value of  $x$ , which ought to be close to the local minimum at  $x=-3$ .

Plotting the function  $f(x)$  and the values of  $x$  for each iteration allows us to see the gradient descent process in action –

```

#Create a range of x values to plot
x_vals = np.linspace(-1, 5, 100)

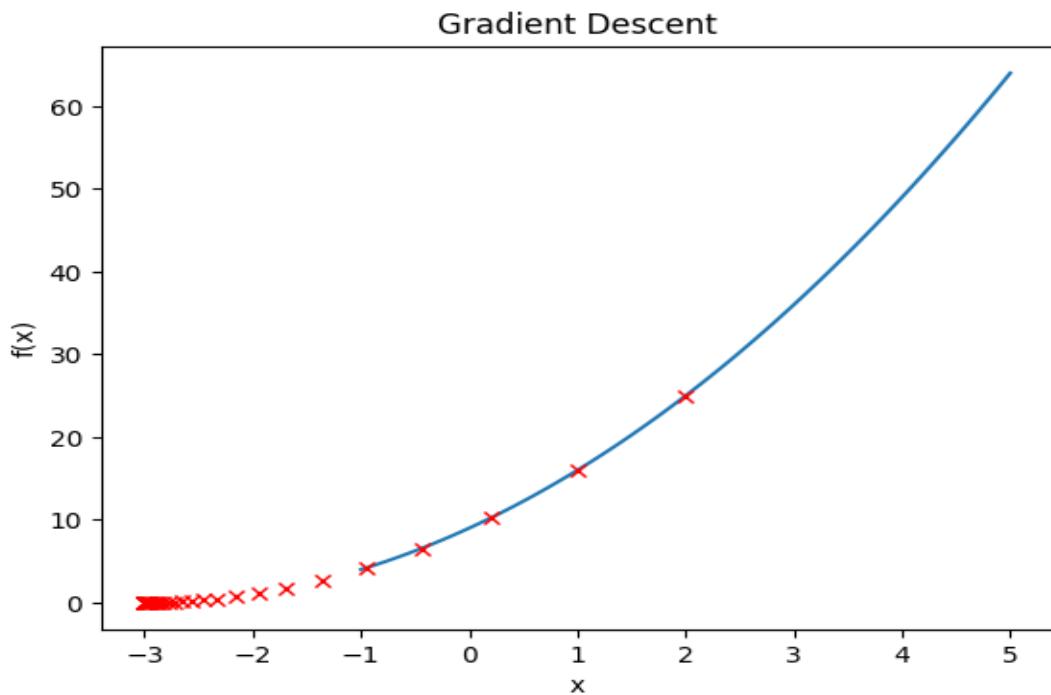
#Plot the function f(x)
plt.plot(x_vals, f(x_vals))

# Plot the values of x at each iteration
plt.plot(x_history, f(np.array(x_history)), 'rx')

#Label the axes and add a title
plt.xlabel('x')
plt.ylabel('f(x)')
plt.title('Gradient Descent')

#Show the plot
plt.show()

```



**Conclusion:** Thus we have implemented Gradient Descent Algorithm to find the local minima of a function  $y=(x+3)^2$  starting from the point  $x=2$  that is  $-3$ .

## Implement K-Means clustering/ hierarchical clustering on sales\_data\_sample.csv dataset. Determine the number of clusters using the elbow method.

```
In [198]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
#Importing the required libraries.

In [199]: from sklearn.cluster import KMeans
from sklearn.decomposition import PCA #Linear Dimensionality reduction.

In [200]: df = pd.read_csv("sales_data_sample.csv") #Loading the dataset.

Preprocessing

In [201]: df.head()
Out[201]:


|   | ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES   | ORDERDATE      | STATUS  | QTR_ID | MONTH_ID | YEAR_ID  | ... ADDRESSLINE1     |
|---|-------------|-----------------|-----------|-----------------|---------|----------------|---------|--------|----------|----------|----------------------|
| 0 | 10107       | 30              | 95.70     | 2               | 2871.00 | 2/24/2003 0:00 | Shipped | 1      | 2        | 2003 ... | 897 Long Aire Aven   |
| 1 | 10121       | 34              | 81.35     | 5               | 2765.90 | 5/7/2003 0:00  | Shipped | 2      | 5        | 2003 ... | 59 Rue Main          |
| 2 | 10134       | 41              | 94.74     | 2               | 3884.34 | 7/1/2003 0:00  | Shipped | 3      | 7        | 2003 ... | 27 Rue Colonel Pie A |
| 3 | 10145       | 45              | 83.26     | 6               | 3746.70 | 8/25/2003 0:00 | Shipped | 3      | 8        | 2003 ... | 78934 Hillside       |


5 rows × 25 columns

In [202]: df.shape
Out[202]: (2823, 25)

In [203]: df.describe()
Out[203]:


|       | ORDERNUMBER  | QUANTITYORDERED | PRICEEACH   | ORDERLINENUMBER | SALES        | QTR_ID      | MONTH_ID    | YEAR_ID     | MSRP        |
|-------|--------------|-----------------|-------------|-----------------|--------------|-------------|-------------|-------------|-------------|
| count | 2823.000000  | 2823.000000     | 2823.000000 | 2823.000000     | 2823.000000  | 2823.000000 | 2823.000000 | 2823.000000 | 2823.000000 |
| mean  | 10258.2515   | 35.092809       | 83.658544   | 6.466171        | 3553.889072  | 2.717476    | 7.092455    | 2003.15059  | 100.715551  |
| std   | 92.085478    | 9.741443        | 20.174277   | 4.225041        | 1841.865106  | 1.203878    | 3.656633    | 0.699667    | 40.187912   |
| min   | 10100.000000 | 6.000000        | 26.800000   | 1.000000        | 482.130000   | 1.000000    | 1.000000    | 2003.000000 | 33.000000   |
| 25%   | 10108.000000 | 27.000000       | 68.800000   | 3.000000        | 2203.430000  | 2.000000    | 4.000000    | 2003.000000 | 68.000000   |
| 50%   | 10262.000000 | 35.000000       | 95.700000   | 6.000000        | 3184.800000  | 3.000000    | 8.000000    | 2004.000000 | 99.000000   |
| 75%   | 10333.500000 | 43.000000       | 100.000000  | 9.000000        | 4508.000000  | 4.000000    | 11.000000   | 2004.000000 | 124.000000  |
| max   | 10425.000000 | 97.000000       | 100.000000  | 18.000000       | 14082.800000 | 4.000000    | 12.000000   | 2005.000000 | 214.000000  |


In [204]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2823 entries, 0 to 2822
Data columns (total 25 columns):
 # Column          Non-Null Count  Dtype  
--- 
 0  ORDERNUMBER      2823 non-null   int64  
 1  QUANTITYORDERED 2823 non-null   int64  
 2  PRICEEACH        2823 non-null   float64 
 3  ORDERLINENUMBER 2823 non-null   int64  
 4  SALES            2823 non-null   float64 
 5  ORDERDATE        2823 non-null   object  
 6  QTR_ID           2823 non-null   object  
 7  QTR_ID           2823 non-null   int64  
 8  MONTH_ID         2823 non-null   int64  
 9  YEAR_ID          2823 non-null   int64  
 10 PRODUCTLINE      2823 non-null   object  
 11 MSRP             2823 non-null   int64  
 12 PRODUCTCODE      2823 non-null   object  
 13 CUSTOMERNAME     2823 non-null   object  
 14 PHONE            2823 non-null   object  
 15 ADDRESSLINE1     2823 non-null   object  
 16 ADDRESSLINE2     2823 non-null   object  
 17 CITY              2823 non-null   object  
 18 STATE             1337 non-null   object  
 19 POSTALCODE       2747 non-null   object  
 20 COUNTRY           70 non-null    object  
 21 TERRITORY         1074 non-null   object  
 22 CONTACTLASTNAME  2823 non-null   object  
 23 CONTACTFIRSTNAME 2823 non-null   object  
 24 DEALSIZE          2823 non-null   object  
dtypes: float64(2), int64(7), object(16)
memory usage: 551.5+ KB

In [205]: df.isnull().sum()
Out[205]:


| ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES | QTR_ID | MONTH_ID | YEAR_ID | MSRP |
|-------------|-----------------|-----------|-----------------|-------|--------|----------|---------|------|
| 0           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 1           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 2           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 3           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 4           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 5           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 6           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 7           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 8           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 9           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 10          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 11          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 12          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 13          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 14          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 15          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 16          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 17          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 18          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 19          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 20          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 21          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 22          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 23          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 24          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 25          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |


In [206]: df.dtypes
Out[206]:


| ORDERNUMBER      | int64   |
|------------------|---------|
| QUANTITYORDERED  | int64   |
| PRICEEACH        | float64 |
| ORDERLINENUMBER  | int64   |
| SALES            | float64 |
| ORDERDATE        | object  |
| STATUS           | object  |
| QTR_ID           | int64   |
| MONTH_ID         | int64   |
| YEAR_ID          | int64   |
| PRODUCTLINE      | object  |
| MSRP             | int64   |
| PRODUCTCODE      | object  |
| CUSTOMERNAME     | object  |
| PHONE            | object  |
| ADDRESSLINE1     | object  |
| ADDRESSLINE2     | object  |
| CITY             | object  |
| STATE            | object  |
| POSTALCODE       | int64   |
| COUNTRY          | object  |
| TERRITORY        | object  |
| CONTACTLASTNAME  | object  |
| CONTACTFIRSTNAME | object  |
| DEALSIZE         | object  |


In [207]: df.drop = ['ADDRESSLINE1', 'ADDRESSLINE2', 'STATUS', 'POSTALCODE', 'CITY', 'TERRITORY', 'PHONE', 'STATE', 'CONTACTFIRSTNAME']
df = df.drop(df.drop, axis=1) #Dropping the categorical unnecessary columns along with columns having null values. Can't fill the null values as there are a lot of null values.

In [208]: df.isnull().sum()
Out[208]:


| ORDERNUMBER | QUANTITYORDERED | PRICEEACH | ORDERLINENUMBER | SALES | QTR_ID | MONTH_ID | YEAR_ID | MSRP |
|-------------|-----------------|-----------|-----------------|-------|--------|----------|---------|------|
| 0           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 1           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 2           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 3           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 4           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 5           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 6           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 7           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 8           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 9           | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 10          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 11          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 12          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 13          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 14          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 15          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 16          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 17          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 18          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 19          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 20          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 21          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 22          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 23          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 24          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |
| 25          | 0               | 0         | 0               | 0     | 0      | 0        | 0       | 0    |


In [209]: df.dtypes
Out[209]:

| QUANTITYORDERED | int64 |
| --- | --- |



<tbl_r cells="2"
```



**GURU GOBIND SINGH FOUNDATION'S**

**Guru Gobind Singh College of Engineering &  
Research Centre, Nashik**

**Department of Computer Engineering**

**Academic Year 2023-24**

**Final year Sem I**

**Machine Learning**

**Miniproject Report**

**Class: BE Computer**

**Guided By:**

**Mr. P. K. Bachhav**

**Title:** Loan Prediction Using Machine Learning

**Group Members:**

<b>Roll No.</b>	<b>Name</b>
B-28	Ruchita Rupesh Raut
B-16	Divesh Sanjay Patil

# **Mini-Project Report**

**Name of Program:** Computer Engineering  
**Semester:** BECO-Sem I  
**Name of Course:** Laboratory Practice-III (Machine Learning)

**Academic Year:** 2023-24  
**Course code:** 410246

## **Title of Mini-Project: "Loan Prediction Using Machine Learning"**

### **1.0 Rationale:**

The rationale for utilizing machine learning in loan prediction is grounded in the need for more accurate, data-driven, and efficient lending practices. Traditional loan approval processes, which often rely heavily on manual assessments and subjective judgment, can be error-prone and time-consuming. Machine learning offers a solution by leveraging historical data to create predictive models that can analyze an applicant's financial history, creditworthiness, and various risk factors in a systematic and unbiased manner. These models can identify patterns and relationships that are not easily discernible through human assessment, leading to more informed lending decisions. By using machine learning in loan prediction, financial institutions can mitigate risks associated with loan defaults, optimize their lending portfolios, and improve the overall quality of their loan portfolios. Additionally, this approach can enhance the accessibility of loans for qualified borrowers, fostering financial inclusion and fairness. In summary, the rationale for loan prediction using machine learning lies in its potential to enhance the accuracy and efficiency of lending processes, reduce risks, and promote a more inclusive and data-driven approach to financial services.

### **2.0 Aim /Benefits of Mini-Project:**

The aim of loan prediction using Machine Learning is to develop a predictive model that can assess the creditworthiness of loan applicants with greater accuracy and efficiency. This model leverages historical data, including information about previous loan applicants, their financial profiles, and whether they repaid their loans on time or defaulted. By analyzing this data, the machine learning model can identify patterns and factors that are indicative of a borrower's likelihood to repay a loan successfully or become a defaulter. The ultimate goal is to assist financial institutions, such as banks and lending agencies, in making more informed lending decisions, reducing the risk of loan defaults, and ensuring that loans are allocated to individuals and businesses who are more likely to meet their repayment obligations. This not only benefits the lenders by minimizing financial losses but also aids borrowers in accessing loans that align with their financial capabilities. Machine learning plays a crucial role in automating and optimizing the loan approval process, making it more efficient and equitable.

### **Course Outcomes achieved (COs):**

- a) Apply preprocessing techniques on datasets (CO1)
- b) Implement and evaluate linear regression and random forest regression models (CO2)

## **4.0 Literature Review: -**

The use of machine learning in loan prediction has gained significant interest in the financial industry. Researchers have developed predictive models using various algorithms, such as decision trees, random forests, support vector machines, and neural networks, to assess creditworthiness. These models use historical data on loan applicants to make informed decisions. The literature highlights the potential benefits of machine learning in improving loan approval processes and identifying non-linear and complex data patterns. Researchers have also investigated the ethical and fairness aspects of machine learning-based loan prediction to prevent biases or discrimination in lending. The literature highlights the growing importance of machine learning in transforming lending practices and addressing risk management concerns in the banking and financial sectors.

## **5.0 Actual Methodology followed:**

In this Loan Prediction using Machine Learning Mini Project, the methodology begins with loading the dataset and conducting an initial exploration to understand its structure and content. Data preprocessing steps follow, addressing missing values and creating new features that categorize passenger ages and titles. These categorical variables are then one-hot encoded for machine learning. Data visualization is employed to visually analyze survival rates across different passenger groups. The dataset is split into features and the target variable for training and validation. Decision tree is a type of supervised education algorithm (having a pre-defined target variable) that is generally used in category problems. In this approach, we disassociate the population or sample into two or added homogeneous sets (or sub-populations) based on the most significant splitter/ differentiator in input variables. Decision trees use multiple algorithms to decide to disunite a bump into two or added sub-knots. The creation of sub-knots increases the unsophistication of attendant sub-knots. In other words, we can say that chasteness of the bump increases with respect to the target variable. This comprehensive methodology covers data analysis, preprocessing, machine learning, and model evaluation, ultimately allowing practical predictions based on historical data.

## **6.0 Actual Code of Program:**

```
In [1]: import pandas as pd  
import numpy as np
```

```
In [2]: df = pd.read_csv("loan_approval_dataset.csv")
```

```
In [3]: df.head()
```

```
Out[3]:   loan_id  no_of_dependents  education  self_employed  income_annum  loan_amount  loan_te  
0          1                  2    Graduate        No      9600000  29900000  
1          2                  0  Not Graduate     Yes      4100000  12200000  
2          3                  3    Graduate        No      9100000  29700000  
3          4                  3    Graduate        No      8200000  30700000  
4          5                  5  Not Graduate     Yes      9800000  24200000
```

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 4269 entries, 0 to 4268  
Data columns (total 13 columns):  
 #   Column           Non-Null Count  Dtype     
---  
 0   loan_id          4269 non-null   int64    
 1   no_of_dependents 4269 non-null   int64    
 2   education        4269 non-null   object    
 3   self_employed    4269 non-null   object    
 4   income_annum    4269 non-null   int64    
 5   loan_amount      4269 non-null   int64    
 6   loan_term        4269 non-null   int64    
 7   cibil_score      4269 non-null   int64    
 8   residential_assets_value 4269 non-null   int64    
 9   commercial_assets_value 4269 non-null   int64    
 10  luxury_assets_value 4269 non-null   int64    
 11  bank_asset_value 4269 non-null   int64    
 12  loan_status      4269 non-null   object    
dtypes: int64(10), object(3)  
memory usage: 433.7+ KB
```

```
In [5]: df.describe()
```

Out[5]:

	<b>loan_id</b>	<b>no_of_dependents</b>	<b>income_annum</b>	<b>loan_amount</b>	<b>loan_term</b>	<b>cibil_score</b>
<b>count</b>	4269.000000	4269.000000	4.269000e+03	4.269000e+03	4269.000000	4269.000000
<b>mean</b>	2135.000000	2.498712	5.059124e+06	1.513345e+07	10.900445	599.936051
<b>std</b>	1232.498479	1.695910	2.806840e+06	9.043363e+06	5.709187	172.430401
<b>min</b>	1.000000	0.000000	2.000000e+05	3.000000e+05	2.000000	300.000000
<b>25%</b>	1068.000000	1.000000	2.700000e+06	7.700000e+06	6.000000	453.000000
<b>50%</b>	2135.000000	3.000000	5.100000e+06	1.450000e+07	10.000000	600.000000
<b>75%</b>	3202.000000	4.000000	7.500000e+06	2.150000e+07	16.000000	748.000000
<b>max</b>	4269.000000	5.000000	9.900000e+06	3.950000e+07	20.000000	900.000000

In [6]: `df.shape`

Out[6]: (4269, 13)

In [7]: `df['Movable_assets'] = df['bank_asset_value'] + df['luxury_assets_value']`

`df['Immovable_assets'] = df['residential_assets_value'] + df['commercial_asset`

In [8]: `df.drop(columns=['bank_asset_value', 'luxury_assets_value', 'residential_asset`

In [9]: `def uniquevals(col):`  
 `print(f'Unique Values in {col} is : {df[col].unique()}')`

`def valuecounts(col):`  
 `print(f'Valuecounts of {col} is: {len(df[col].value_counts())}')`

`for col in df.columns:`  
 `valuecounts(col)`  
`# uniquevals(col)`  
 `print("-"*75)`

```
Valuecounts of loan_id is: 4269
-----
Valuecounts of no_of_dependents is: 6
-----
Valuecounts of education is: 2
-----
Valuecounts of self_employed is: 2
-----
Valuecounts of income_anum is: 98
-----
Valuecounts of loan_amount is: 378
-----
Valuecounts of loan_term is: 10
-----
Valuecounts of cibil_score is: 601
-----
Valuecounts of loan_status is: 2
-----
Valuecounts of Movable_assets is: 484
-----
Valuecounts of Immovable_assets is: 406
```

```
In [10]: catvars = df.select_dtypes(include=['object']).columns
          numvars = df.select_dtypes(include = ['int32','int64','float32','float64']).columns
          catvars, numvars
```

```
Out[10]: (Index(['education', 'self_employed', 'loan_status'], dtype='object'),
           Index(['loan_id', 'no_of_dependents', 'income_anum', 'loan_amount',
                  'loan_term', 'cibil_score', 'Movable_assets', 'Immovable_assets'],
                  dtype='object'))
```

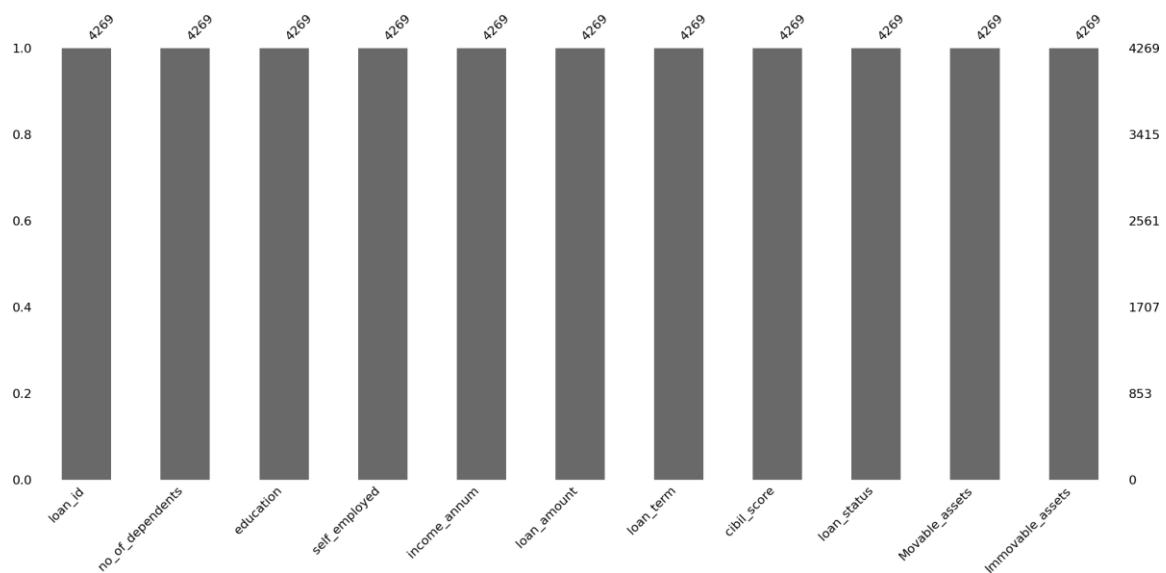
```
In [11]: pip install missingno
```

```
Requirement already satisfied: missingno in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (0.5.2)
Requirement already satisfied: numpy in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from missingno) (1.24.2)
Requirement already satisfied: matplotlib in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from missingno) (3.7.1)
Requirement already satisfied: scipy in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from missingno) (1.10.1)
Requirement already satisfied: seaborn in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from missingno) (0.12.2)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (1.0.7)
Requirement already satisfied: cycler>=0.10 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (4.3.9.3)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (1.4.4)
Requirement already satisfied: packaging>=20.0 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (23.0)
Requirement already satisfied: pillow>=6.2.0 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (9.5.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from matplotlib->missingno) (2.8.2)
Requirement already satisfied: pandas>=0.25 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from seaborn->missingno) (2.0.0)
Requirement already satisfied: pytz>=2020.1 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from pandas>=0.25->seaborn->missingno) (2023.3)
Requirement already satisfied: tzdata>=2022.1 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from pandas>=0.25->seaborn->missingno) (2023.3)
Requirement already satisfied: six>=1.5 in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (from python-dateutil>=2.7->matplotlib->missingno) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
```

```
In [12]: import missingno as msno
import warnings
with warnings.catch_warnings():
    warnings.simplefilter("ignore")
```

```
In [13]: msno.bar(df)
```

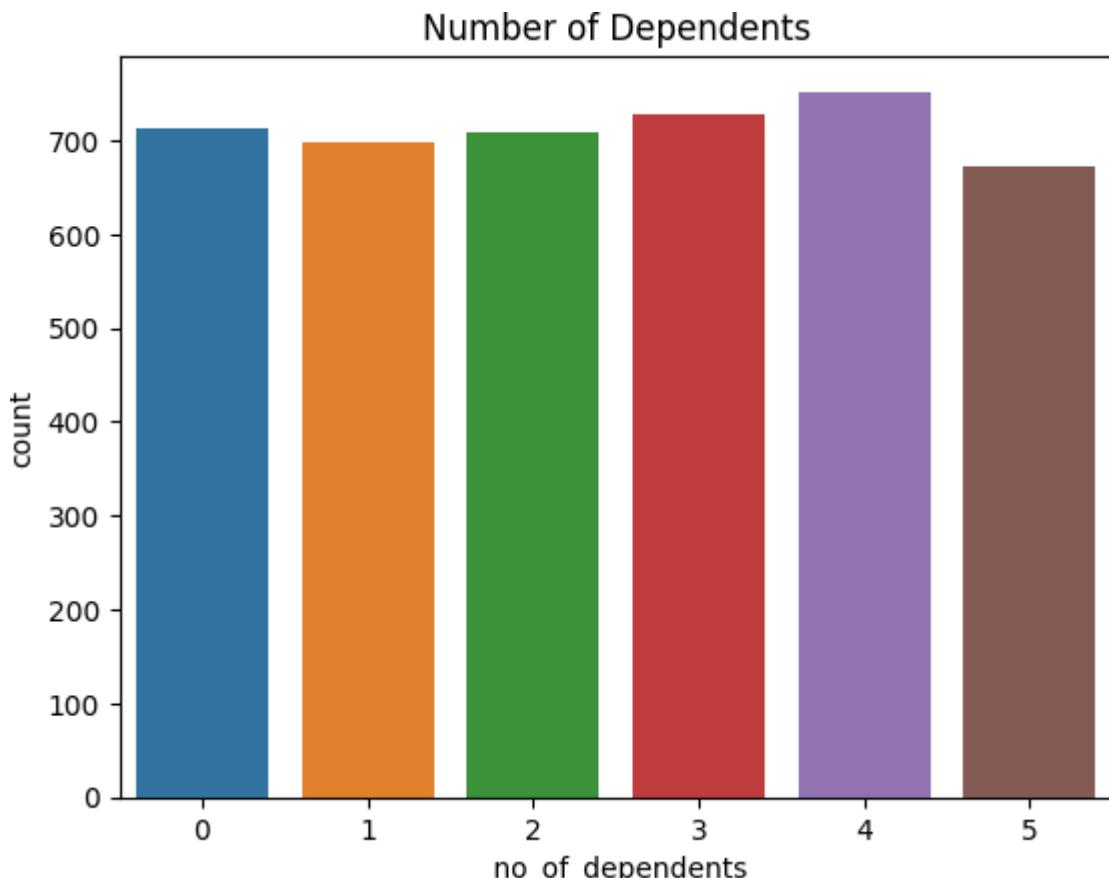
```
Out[13]: <Axes: >
```



```
In [14]: import matplotlib.pyplot as plt
%matplotlib inline

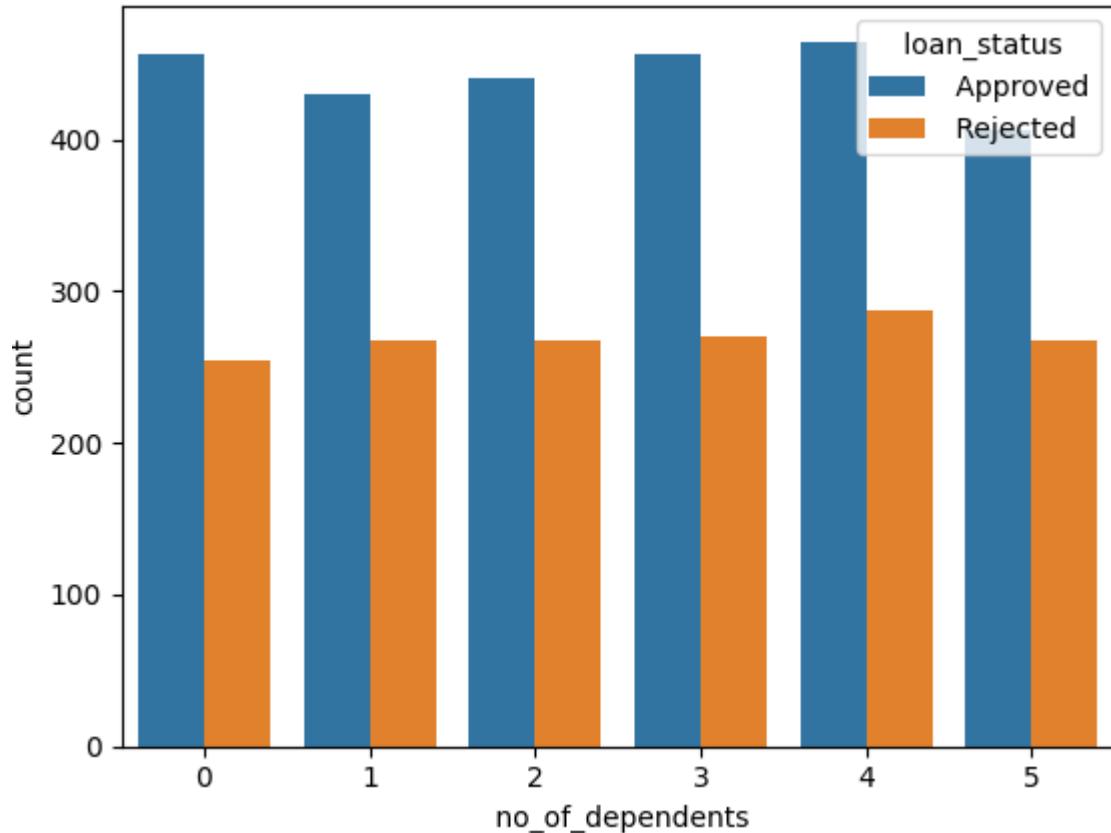
import seaborn as sns
sns.countplot(x = 'no_of_dependents', data = df).set_title('Number of Dependents')
```

Out[14]: Text(0.5, 1.0, 'Number of Dependents')



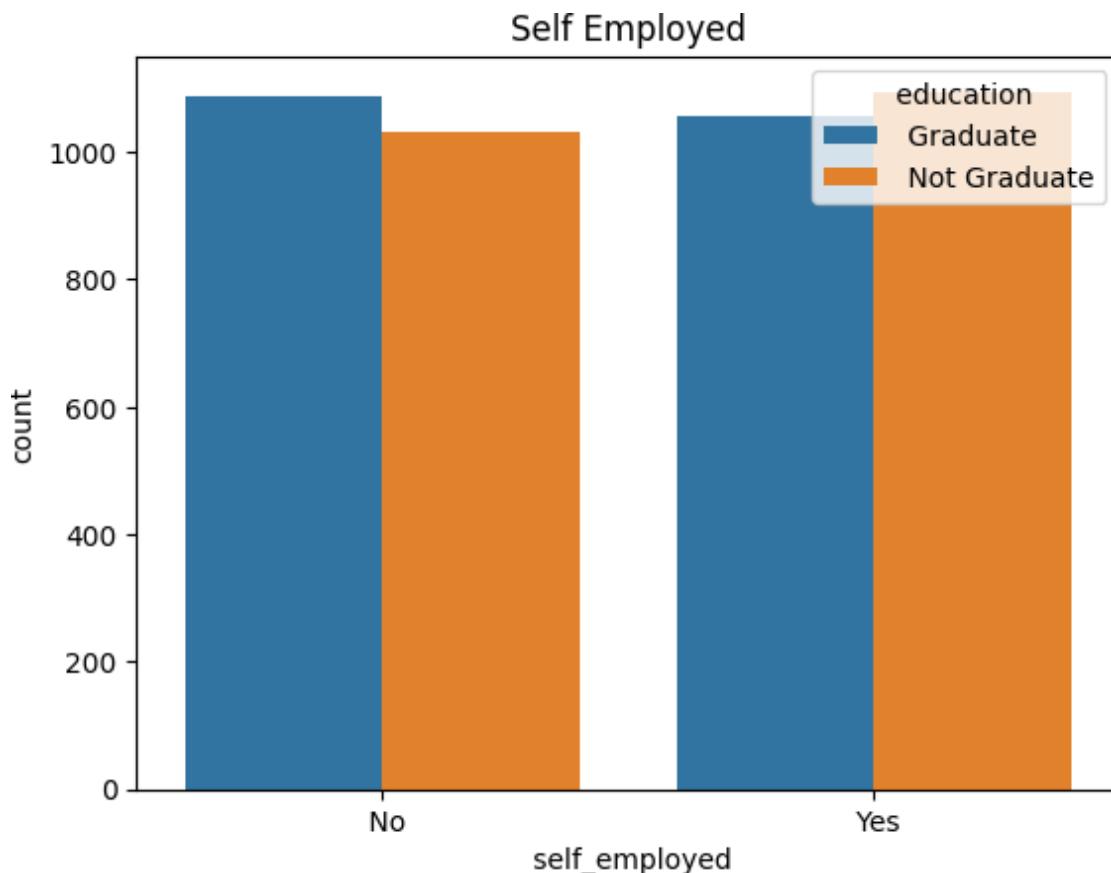
```
In [15]: sns.countplot(x = 'no_of_dependents', data = df, hue = 'loan_status')
```

Out[15]: <Axes: xlabel=' no\_of\_dependents', ylabel='count'>



```
In [16]: sns.countplot(x='self_employed', data = df, hue = 'education').set_title('Self Employed')
```

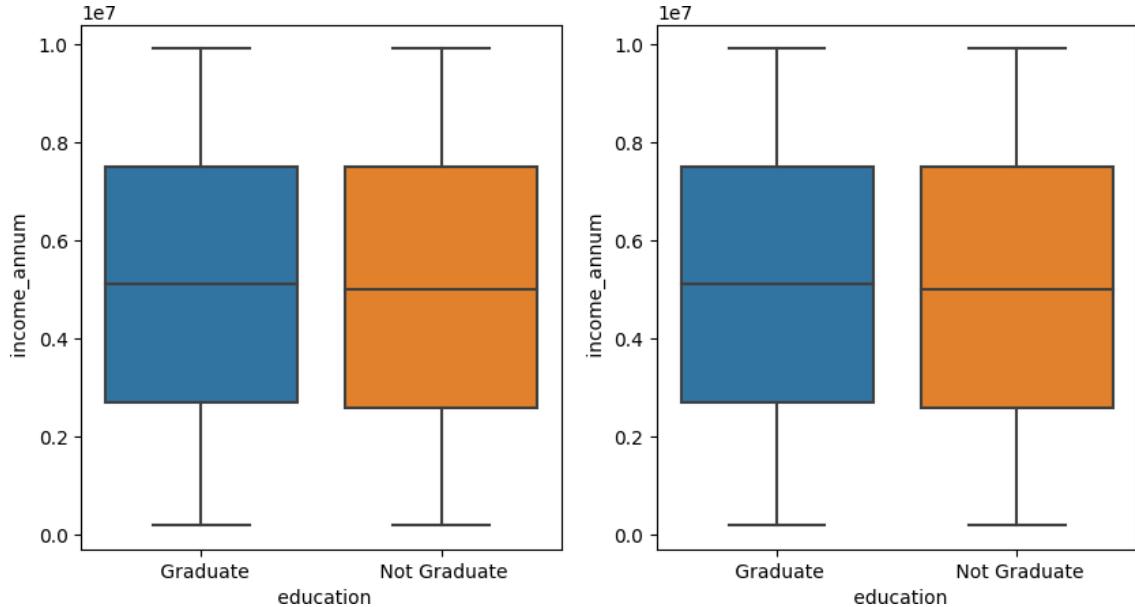
```
Out[16]: Text(0.5, 1.0, 'Self Employed')
```



```
In [17]: fig, ax = plt.subplots(1,2,figsize=(10, 5))  
sns.boxplot(x = 'education', y = 'income_annum', data = df, ax=ax[0])
```

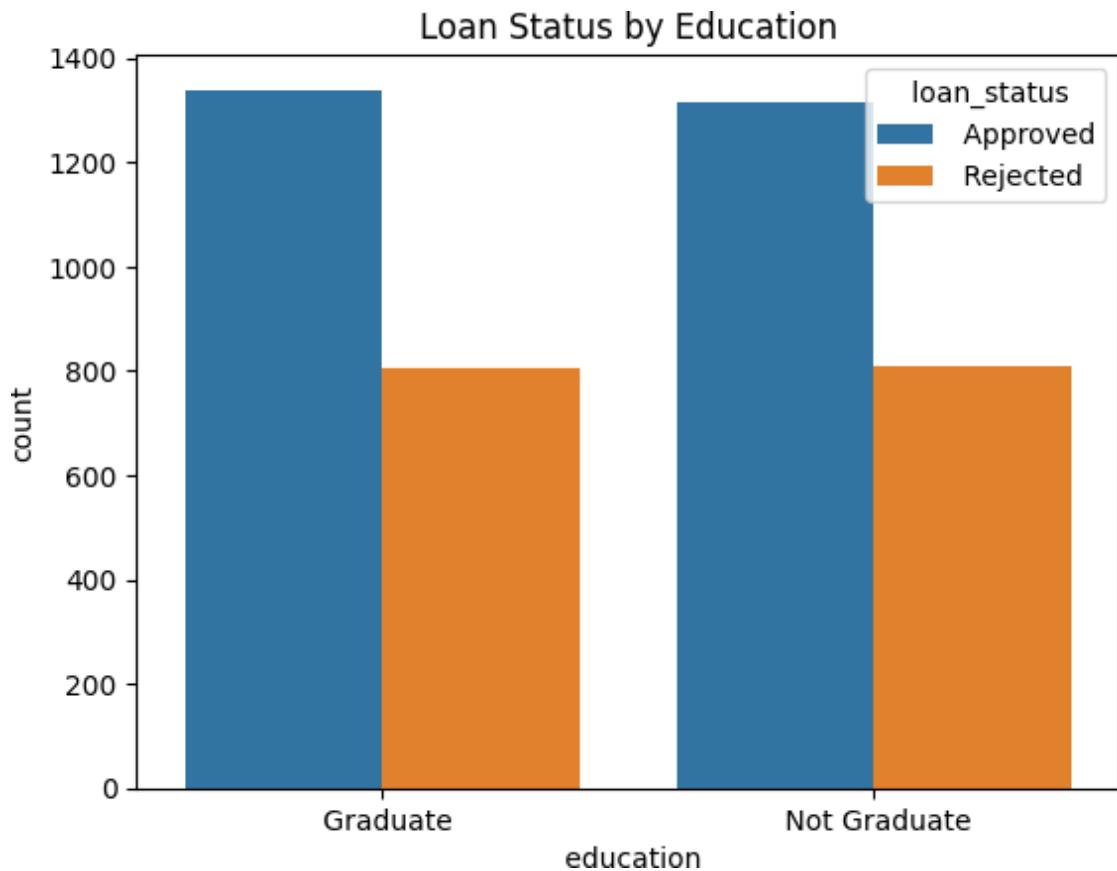
```
sns.boxplot(x = 'education', y = 'income_annum', data = df, ax=ax[1])
```

Out[17]: <Axes: xlabel=' education', ylabel=' income\_annum'>



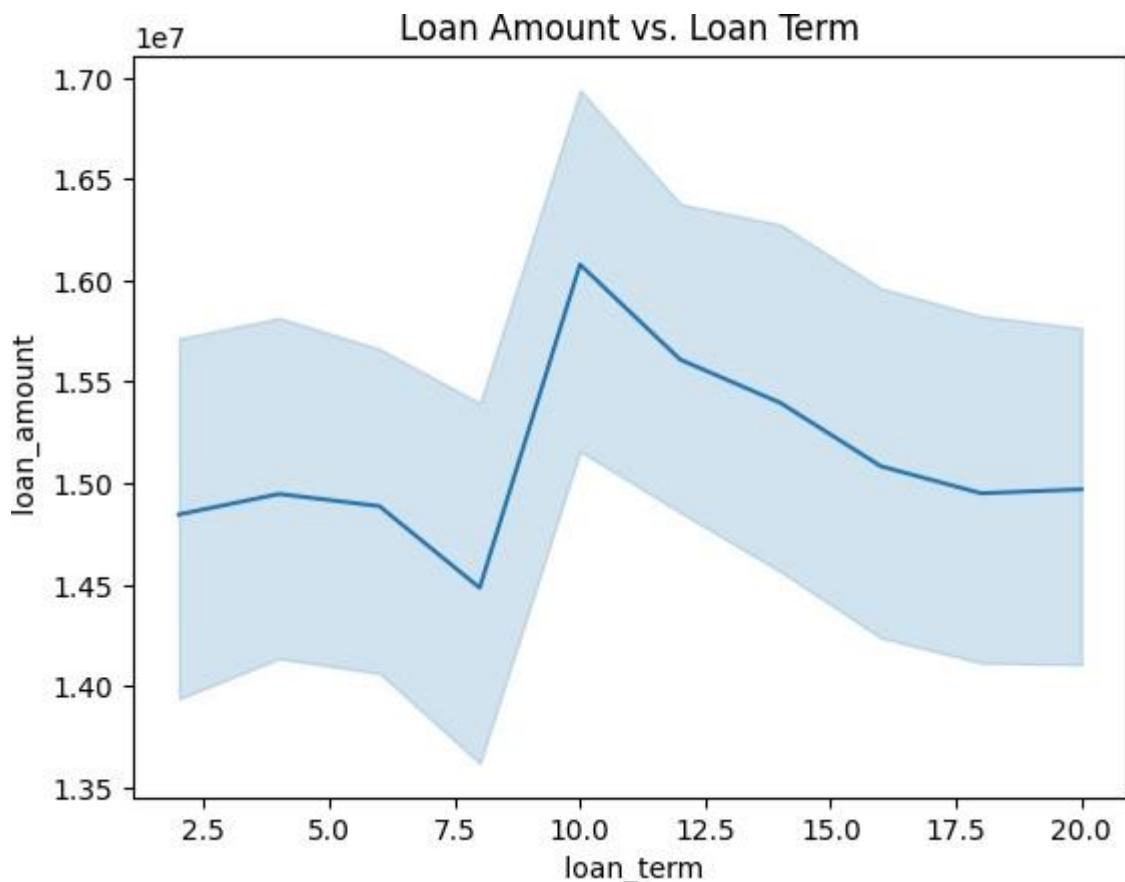
In [18]: sns.countplot(x = 'education', hue = 'loan\_status', data = df).set\_title('Loan Status by Education')

Out[18]: Text(0.5, 1.0, 'Loan Status by Education')



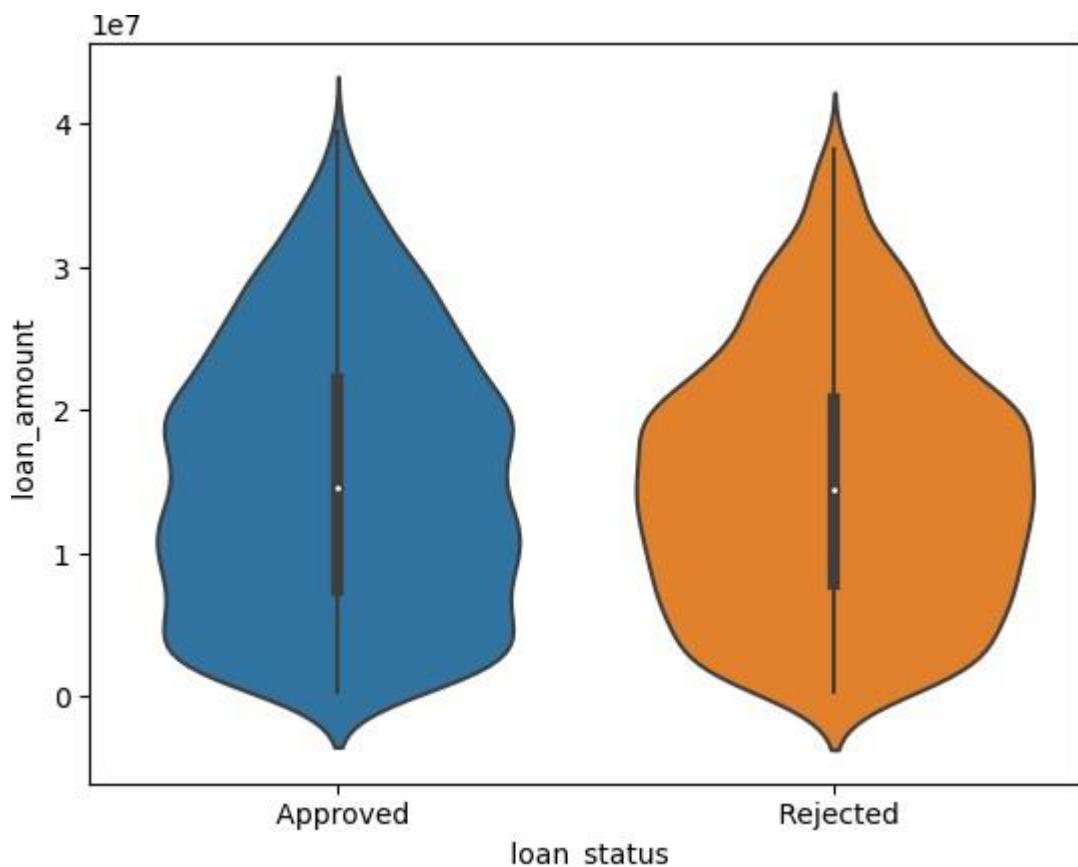
In [19]: sns.lineplot(x = 'loan\_term', y = 'loan\_amount', data = df).set\_title('Loan Amount vs. Loan Term')

Out[19]: Text(0.5, 1.0, 'Loan Amount vs. Loan Term')



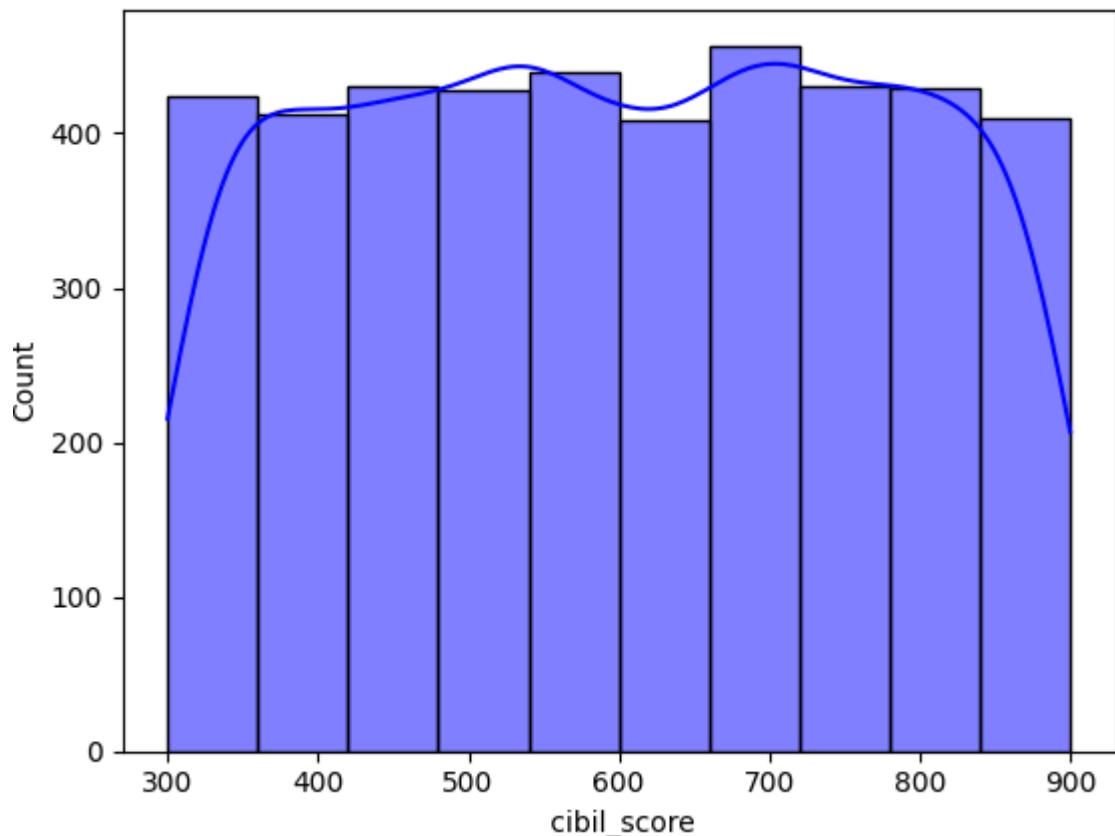
```
In [20]: sns.violinplot(x='loan_status', y='loan_amount', data=df)
```

```
Out[20]: <Axes: xlabel='loan_status', ylabel='loan_amount'>
```



```
In [21]: sns.histplot(df['cibil_score'], bins=10, kde=True, color='blue')
```

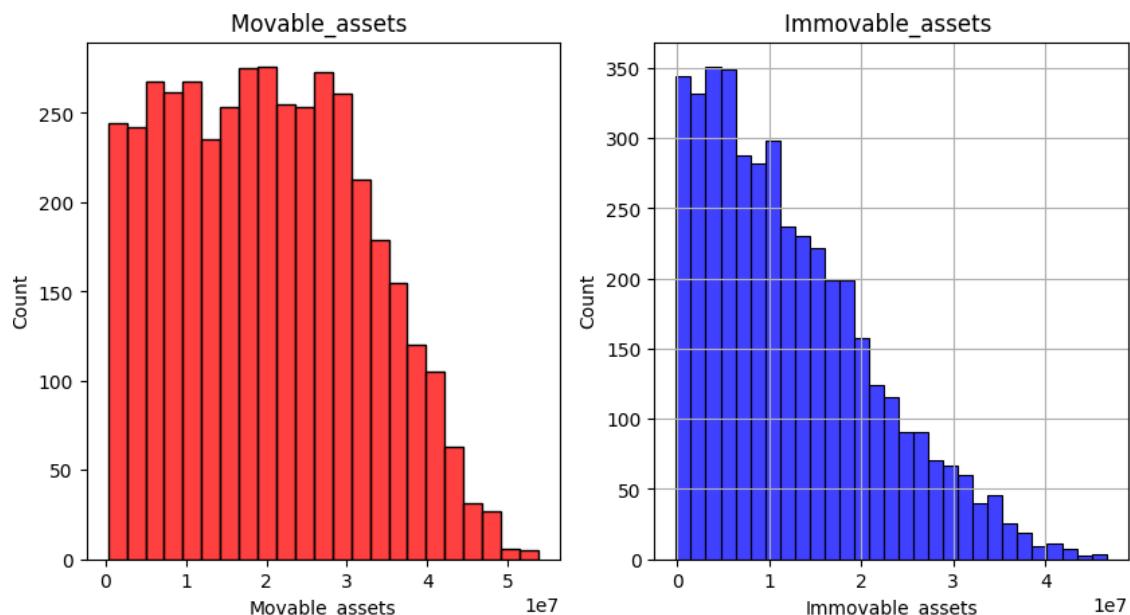
```
Out[21]: <Axes: xlabel=' cibil_score', ylabel='Count'>
```



```
In [22]: fig, ax = plt.subplots(1,2,figsize=(10,5))
plt.subplot(1, 2, 1)
sns.histplot(df['Movable_assets'], ax=ax[0], color='red')
plt.title("Movable_assets ")

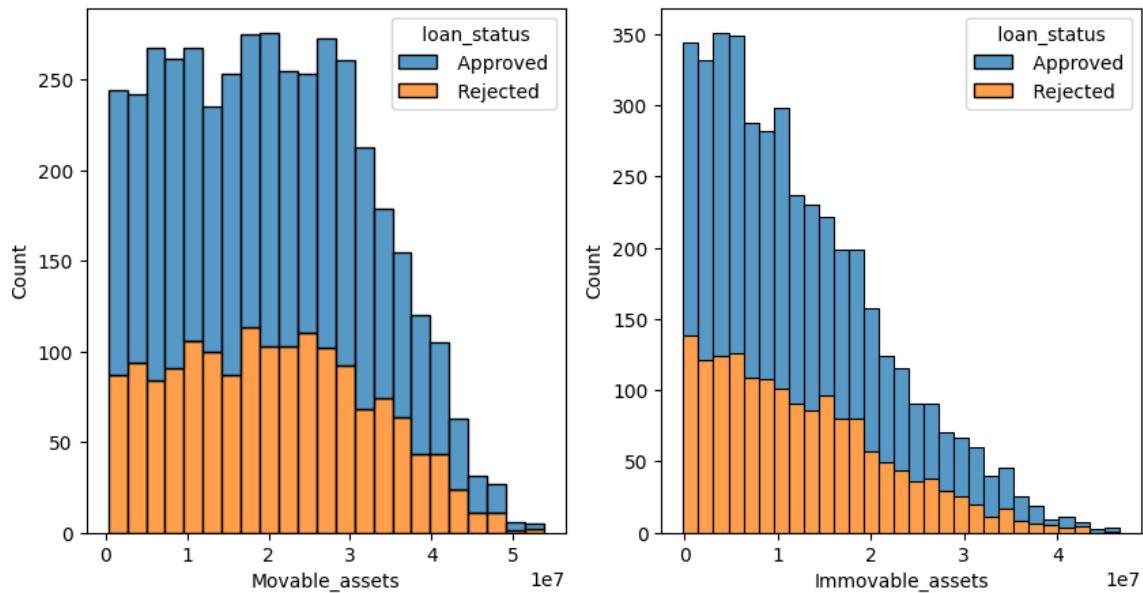
plt.subplot(1, 2, 2)
plt.grid()
sns.histplot(df['Immovable_assets'], ax=ax[1], color='blue')
plt.title("Immovable_assets ")
```

```
Out[22]: Text(0.5, 1.0, 'Immovable_assets ')
```



```
In [23]: fig, ax = plt.subplots(1,2,figsize=(10,5))
sns.histplot(x = 'Movable_assets', data = df, ax=ax[0], hue = ' loan_status', m
sns.histplot(x = 'Immovable_assets', data = df, ax=ax[1], hue = ' loan_status'
```

Out[23]: <Axes: xlabel='Immovable\_assets', ylabel='Count'>



```
In [24]: df[' education'] = df[' education'].map({' Not Graduate':0, ' Graduate':1})
df[' self-employed'] = df[' self-employed'].map({' No':0, ' Yes':1})
df[' loan_status'] = df[' loan_status'].map({' Rejected':0, ' Approved':1})
```

In [25]: pip install sklearn

Requirement already satisfied: sklearn in c:\users\parag\appdata\local\programs\python\python310\lib\site-packages (0.0.post5)  
Note: you may need to restart the kernel to use updated packages.

```
In [26]: from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler,StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.metrics import accuracy_score
from sklearn.compose import ColumnTransformer
from sklearn import tree
```

In [27]: X\_train, X\_test, y\_train, y\_test = train\_test\_split(df.drop(' loan\_status', axis=1), df[' loan\_status'])

```
In [28]: from sklearn.linear_model import LogisticRegression

lgr = LogisticRegression()

lgr.fit(X_train, y_train)

predictions = lgr.predict(X_test)
```

In [29]: accuracy = accuracy\_score(y\_test, predictions)
print("Accuracy:", accuracy)

Accuracy: 0.6276346604215457

In [30]: from sklearn.svm import SVC
model = SVC()

```
model.fit(X_train, y_train)

predictions = model.predict(X_test)
```

In [31]: `from sklearn.ensemble import RandomForestClassifier

# Create a random forest classifier
rfc = RandomForestClassifier()`

In [32]: `rfc.fit(X_train, y_train)`

Out[32]: `▼ RandomForestClassifier
RandomForestClassifier()`

In [33]: `rfc_pred = rfc.predict(X_test)`

In [34]: `accuracy = accuracy_score(y_test, rfc_pred)
print("Accuracy:", accuracy)`

Accuracy: 0.9800936768149883

In [37]: `from sklearn.tree import DecisionTreeClassifier

# Create decision tree object
dtree = DecisionTreeClassifier()`

In [38]: `dtree.fit(X_train, y_train)`

Out[38]: `▼ DecisionTreeClassifier
DecisionTreeClassifier()`

In [39]: `dtree_pred = dtree.predict(X_test)`

In [40]: `dtree.score(X_train, y_train)`

Out[40]: 1.0

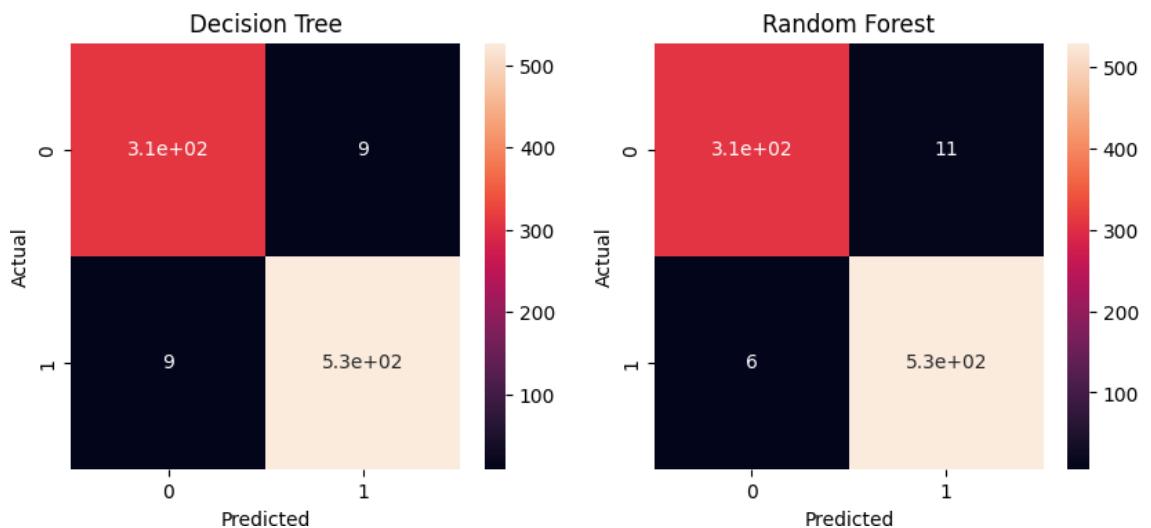
In [41]: `accuracy = accuracy_score(y_test, dtree_pred)
print("Accuracy:", accuracy)`

Accuracy: 0.9789227166276346

In [42]: `from sklearn.metrics import confusion_matrix

fig, ax = plt.subplots(1,2, figsize=(10,4))
sns.heatmap(confusion_matrix(y_test, dtree_pred), annot=True, ax=ax[0]).set_title(
ax[0].set_xlabel('Predicted')
ax[0].set_ylabel('Actual')
sns.heatmap(confusion_matrix(y_test, rfc_pred), annot=True, ax=ax[1]).set_title(
ax[1].set_xlabel('Predicted')
ax[1].set_ylabel('Actual')`

Out[42]: `Text(518.4494949494949, 0.5, 'Actual')`



In [ ]:

In [ ]:

## **7.0. Actual Resources used:**

S. No.	Name of Resource/material	Specifications	Qty	Remarks
1	Computer System	Windows OS,i3 processor ,8 GB RAM	01	-
2	Software	Jupyter Notebook, MS Office word	01	-
3	Printer, Pages	Canon LaserJet	01	-

## **8.0. Skill Developed / Learning outcome from this Mini-Project:**

1. Understanding the concept of random forest classifier.
2. Improved Data Analysis Skills.
3. Improved Data Preprocessing Skills.
4. Understanding the concept of Data splitting.
5. Understanding the concept of Feature Selection and Feature Engineering.

## **9.0. Applications of Mini Project:**

### **1. Auto Financing:**

Machine Learning models can assist in determining whether an applicant is eligible for an auto loan based on factors such as credit score, income, and vehicle value.

### **2. Real Estate:**

Machine Learning models can help assess the creditworthiness of real estate investors applying for investment loans.

### **3. Student Loans:**

Machine Learning can be used by educational institutions and lenders to assess student loan eligibility based on factors like educational history and future earning potential.

### **4. Small Business Loans:**

Machine Learning models can evaluate the creditworthiness of small business owners applying for loans to support their enterprises.