

```
#loading dataset  
df = pd.read_csv("B04_Ho  
  
df.head()  
  
[ 0 :  
    Loan_ID  Gender  Married  
0   LP001002      Male      No
```

```

4 LP001008 Male No 0 Graduate No 6000 0.0

[4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
 #   Column           Non-Null Count Dtype  
--- 
 0   Loan_ID          614 non-null    object  
 1   Gender            601 non-null    object  
 2   Married           611 non-null    object  
 3   Dependents        599 non-null    object  
 4   Education         614 non-null    object  
 5   Self_Employed     582 non-null    object  
 6   ApplicantIncome   614 non-null    int64  
 7   CoapplicantIncome 614 non-null    float64 
 8   LoanAmount        592 non-null    float64 
 9   Loan_Amount_Term  600 non-null    float64 
 10  Credit_History   564 non-null    float64 
 11  Property_Area    614 non-null    object  
 12  Loan_Status       614 non-null    object  
dtypes: float64(4), int64(1), object(8)
memory usage: 62.5+ KB

[5]: #checking null values
df.isnull().sum()

t[5]: Loan_ID      0
Gender        13
Married       3
Dependents    15
Education     0
Self_Employed 32
ApplicantIncome 0
CoapplicantIncome 0
LoanAmount    22
Loan_Amount_Term 14
Credit_History 50
Property_Area 0
Loan_Status    0
dtype: int64



## Data preprocessing



[6]: #To remove the skewness due to outliers, we use log transformation to get a normal distribution
df['loanAmount_log']= np.log(df['LoanAmount'])
df['loanAmount_log'].hist(bins=20)

t[6]: <AxesSubplot:>


A histogram titled 'loanAmount_log' showing the frequency distribution of the logarithm of the loan amount. The x-axis ranges from 2 to 6, and the y-axis ranges from 0 to 140. The distribution is roughly bell-shaped, centered around 4.5, indicating a normal distribution after the log transformation.

[7]: df.head()

t[7]: 

|   | Loan_ID  | Gender | Married | Dependents | Education    | Self_Employed | ApplicantIncome | CoapplicantIncome | Loan_Amount_Term | Credit_History | Property_Area | Loan_Status |
|---|----------|--------|---------|------------|--------------|---------------|-----------------|-------------------|------------------|----------------|---------------|-------------|
| 0 | LP001002 | Male   | No      | 0          | Graduate     | No            | 5849            | 0.0               | 36               | 1              | 0             | Approved    |
| 1 | LP001003 | Male   | Yes     | 1          | Graduate     | No            | 4583            | 1508.0            | 36               | 1              | 0             | Approved    |
| 2 | LP001005 | Male   | Yes     | 0          | Graduate     | Yes           | 3000            | 0.0               | 36               | 1              | 0             | Approved    |
| 3 | LP001006 | Male   | Yes     | 0          | Not Graduate | No            | 2583            | 2358.0            | 36               | 1              | 0             | Approved    |
| 4 | LP001008 | Male   | No      | 0          | Graduate     | No            | 6000            | 0.0               | 36               | 0              | 0             | Approved    |



[8]: df['TotalIncome'] = df['ApplicantIncome'] + df['CoapplicantIncome']
df['TotalIncome_log']= np.log(df['TotalIncome'])
df['TotalIncome_log'].hist(bins=20)

t[8]: <AxesSubplot:>


A histogram titled 'TotalIncome_log' showing the frequency distribution of the logarithm of total income. The x-axis ranges from 7.5 to 11.5, and the y-axis ranges from 0 to 100. The distribution is roughly bell-shaped, centered around 8.5, indicating a normal distribution after the log transformation.

[9]: # filling the null values
df['Gender'].fillna(df['Gender'].mode()[0], inplace = True)
df['Married'].fillna(df['Married'].mode()[0], inplace = True)
df['Self_Employed'].fillna(df['Self_Employed'].mode()[0], inplace = True)
df['Dependents'].fillna(df['Dependents'].mode()[0], inplace = True)

df['LoanAmount'] = df['LoanAmount'].fillna(df['LoanAmount'].mean())
df['loanAmount_log'] = df['loanAmount_log'].fillna(df['loanAmount_log'].mean())


```

```
df.isnull().sum()
```

```
] : Loan_ID          0
Gender           0
Married          0
Dependents       0
Education         0
Self_Employed    0
ApplicantIncome   0
CoapplicantIncome 0
```

```
1 loanAmount_log      0  
2 TotalIncome          0  
3 TotalIncome_log      0  
4 dtype: int64
```

Data Visualization

univariate analysis

```
[10]: df.Gender.value_counts().plot(kind='bar')
```

```
[10]: <AxesSubplot:>
```

A bar chart comparing the counts of 'Male' and 'Female' gender categories. The x-axis labels are 'Male' and 'Female'. The y-axis ranges from 0 to 500 with increments of 100. The bar for 'Male' reaches approximately 500, while the bar for 'Female' reaches approximately 120.

Gender	Count
Male	500
Female	120

```
[11]: df.Education.value_counts().plot(kind='bar')
```

```
[11]: <AxesSubplot:>
```

A bar chart comparing the counts of 'Graduate' and 'Not Graduate' education levels. The x-axis labels are 'Graduate' and 'Not Graduate'. The y-axis ranges from 0 to 500 with increments of 100. The bar for 'Graduate' reaches approximately 480, while the bar for 'Not Graduate' reaches approximately 130.

Education Level	Count
Graduate	480
Not Graduate	130

```
[12]: sns.countplot(x='Property_Area', data=df)
```

```
[12]: <AxesSubplot:xlabel='Property_Area', ylabel='count'>
```

A bar chart showing the count of properties categorized by area: 'Urban', 'Rural', and 'Semiurban'. The x-axis labels are 'Urban', 'Rural', and 'Semiurban'. The y-axis is labeled 'count' and ranges from 0 to 200 with increments of 50. The bar for 'Urban' reaches approximately 200, the bar for 'Rural' reaches approximately 175, and the bar for 'Semiurban' reaches approximately 230.

Property Area	Count
Urban	200
Rural	175
Semiurban	230

A bar chart titled 'Count' on the y-axis and 'Loan_Status' on the x-axis. The y-axis ranges from 0 to 450 with increments of 50. The x-axis has two categories: '0' and '1'. The bar for '0' reaches approximately 420, and the bar for '1' reaches approximately 450.

Loan_Status	Count
0	~420
1	~450

```
[14]: sns.countplot(x=df.Loan_Status, hue=df.Gender)
```

```
[14]: <AxesSubplot:xlabel='Loan_Status', ylabel='count'>
```

```
[15]: df.Loan_Amount_Term.value_counts().plot(kind='bar')
```

```
[15]: <AxesSubplot:>
```

```
[16]: sns.countplot(x = 'Married', data=df, palette = 'Set1')
```

```
[16]: <AxesSubplot:xlabel='Married', ylabel='count'>
```

```
[17]: sns.countplot(x = 'Dependents', data=df, palette = 'Set1')
```

```
[17]: <AxesSubplot:xlabel='Dependents', ylabel='count'>
```

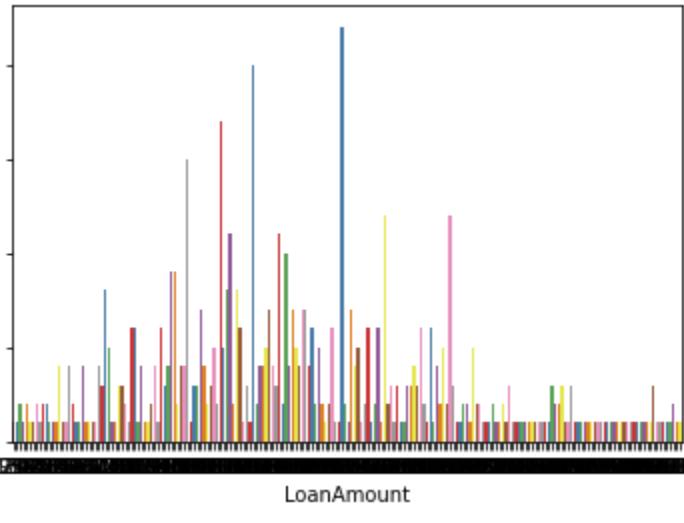
A bar chart titled "Dependents" showing the count of observations for each value of "LoanAmount". The x-axis is labeled "Dependents" and has two categories: 0 and 1. The y-axis is labeled "Count" and ranges from 0 to 150. Category 0 has a red bar reaching approximately 150. Category 1 has a blue bar reaching approximately 100.

Dependents	Count
0	~150
1	~100

[18]:

```
100.000000    15
160.000000    12
...
211.000000    1
250.000000    1
62.000000    1
85.000000    1
436.000000    1
Name: LoanAmount, Length: 204, dtype: int64
```

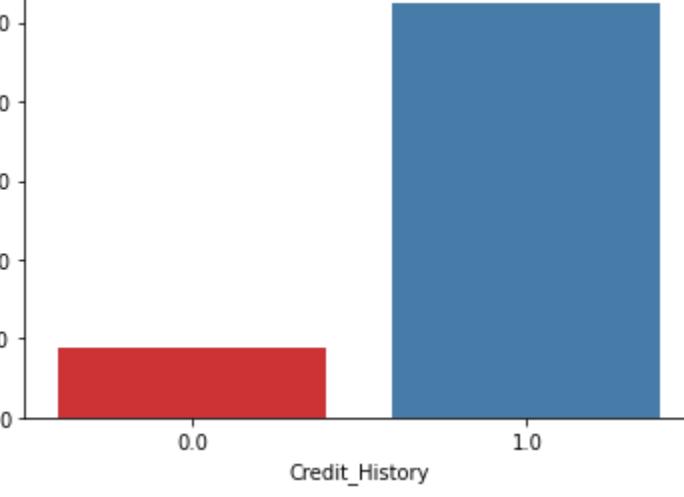
[18]: <AxesSubplot:xlabel='LoanAmount', ylabel='count'>



[19]:

```
sns.countplot(x = 'Credit_History', data=df, palette = 'Set1')
```

[19]: <AxesSubplot:xlabel='Credit_History', ylabel='count'>

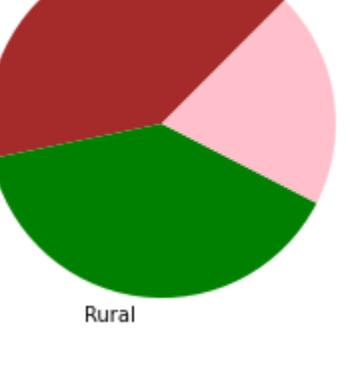


[20]:

```
labels = 'Urban', 'Rural', 'Semiurban'
colors = ['brown', 'green', 'pink']
area = [311, 300, 153]

plt.pie(area, colors = colors, labels = labels,startangle=45)
```

[20]: ([<matplotlib.patches.Wedge at 0x22bad636f40>,
<matplotlib.patches.Wedge at 0x22bad645400>,
<matplotlib.patches.Wedge at 0x22bad645880>],
[Text(-0.5210279395193431, 0.9687775215395059, 'Urban'),
Text(-0.14880903459988948, -1.0898880085685176, 'Rural'),
Text(1.086598312664751, 0.17118442369011413, 'Semiurban')])



[21]:

```
sns.countplot(x = 'Self_Employed', data=df, palette = 'Set1')
```

A bar chart showing the distribution of a variable. The x-axis represents categories numbered 1 to 10. The y-axis represents the count of observations, ranging from 0 to 500. The distribution is highly skewed, with most observations falling into the first few categories.

Category	Count
1	~500
2	~400
3	~300
4	0
5	0
6	0
7	0
8	0
9	0
10	0

Bi-variate Analysis

```
[22]: sns.scatterplot(x=df.ApplicantIncome, y=df.LoanAmount)
```

```
[22]: <AxesSubplot:xlabel='ApplicantIncome', ylabel='LoanAmount'>
```

```
[23]: x = pd.DataFrame(pd.crosstab(df.Education, df.Loan_Status))
x["Total"] = x["Y"] + x["N"]
x["Approval_Rate"] = x["Y"]/x["Total"]*100
x
```

Loan_Status	N	Y	Total	Approval_Rate
Education				
Graduate	140	340	480	70.833333
Not Graduate	52	82	134	61.194030

```
[24]: x = pd.DataFrame(pd.crosstab(df.Gender, df.Loan_Status))
x["Total"] = x["Y"] + x["N"]
x["Approval_Rate"] = x["Y"]/x["Total"]*100
x
```

Loan_Status	N	Y	Total	Approval_Rate
Gender				
Female	37	75	112	66.964286
Male	155	347	502	69.123506

```
[25]: x = pd.DataFrame(pd.crosstab(df.Married, df.Loan_Status))
x["Total"] = x["Y"] + x["N"]
x["Approval_Rate"] = x["Y"]/x["Total"]*100
x
```

Loan_Status	N	Y	Total	Approval_Rate
Married				
No	79	134	213	62.910798
Yes	113	288	401	71.820449

```
[26]: x = pd.DataFrame(pd.crosstab(df.Credit_History, df.Loan_Status))
x["Total"] = x["Y"] + x["N"]
x["Approval_Rate"] = x["Y"]/x["Total"]*100
x
```

Loan_Status	N	Y	Total	Approval_Rate
Credit_History				
0.0	82	7	89	7.865169

```
[27]: x = pd.DataFrame(pd.crosstab(df.Self_Employed, df.Loan_Status))
x["Total"] = x["Y"] + x["N"]
x["Approval_Rate"] = x["Y"] / x["Total"] * 100
x
```

Loan_Status	N	Y	Total	Approval_Rate
Self_Employed				
No	166	366	532	68.796992
Yes	26	56	82	68.292683


```
[28]: x = pd.DataFrame(pd.crosstab(df.Property_Area, df.Loan_Status))
x["Total"] = x["Y"] + x["N"]
x["Approval_Rate"] = x["Y"] / x["Total"] * 100
x
```

Loan_Status	N	Y	Total	Approval_Rate
Property_Area				
Rural	69	110	179	61.452514
Semiurban	54	179	233	76.824034
Urban	69	133	202	65.841584


```
[49]: sns.boxplot(x='Property_Area', y='ApplicantIncome', data=df, palette='rainbow')
```



```
[49]: <AxesSubplot:xlabel='Property_Area', ylabel='ApplicantIncome'>
```

ApplicantIncome	-0.12	0.57	-0.047	-0.019	0.44	0.89	0.72	
CoapplicantIncome	1	0.19	-0.059	0.011	0.21	0.34	0.38	
LoanAmount	0.57	1	0.036	-0.0014	0.9	0.62	0.69	
Loan_Amount_Term	-0.047	-0.059	1	-0.0047	0.085	-0.071	-0.056	
Credit_History	-0.019	0.011	-0.0014	-0.0047	1	-0.019	-0.013	0.021
loanAmount_log	0.44	0.21	0.9	0.085	-0.019	1	0.51	0.66
TotalIncome	0.89	0.34	0.62	-0.071	-0.013	0.51	1	0.85
TotalIncome_log	0.72	0.38	0.69	-0.056	0.021	0.66	0.85	1