



INTRODUCTION

Introduction:

The Home Loan Approval Dataset is a comprehensive collection of data sourced from Skill Circle, containing valuable information for analysis in the domain of home loan approvals. This dataset is particularly relevant for financial institutions, analysts, and researchers interested in understanding the factors influencing home loan approval decisions.

The dataset encompasses various attributes such as applicant income, loan amount, loan term, credit history, property area, and the ultimate decision on loan approval. Through Exploratory Data Analysis (EDA), we aim to gain insights into the patterns, trends, and factors that significantly affect the approval of home loan.

By performing data exploration and visualization techniques, we will delve into the distribution of numeric and categorical variables, identify potential outliers, explore relationships between different attributes, and uncover valuable insights that can aid in making informed decisions related to home loan approval processes.

Dataset Information:

The Home Loan Approval Dataset, sourced from Skill Circle, is designed to facilitate the analysis of home loan approval processes. The dataset contains a total of 614 entries, each representing a home loan application. There are 13 attributes in the dataset, including both numeric and categorical variables.

Attributes:

- **Loan_ID:** Unique Loan ID
- **Gender:** Gender of the applicant (Male/Female)
- **Married:** Applicant's marital status (Yes/No)
- **Dependents:** Number of dependents
- **Education:** Applicant's education level (Graduate/Not Graduate)
- **Self_Employed:** Self-employed status of the applicant (Yes/No)
- **ApplicantIncome:** Applicant's monthly income
- **CoapplicantIncome:** Co-applicant's monthly income
- **LoanAmount:** Loan amount in thousands
- **Loan_Amount_Term:** Term of the loan in months
- **Credit_History:** Credit history meets guidelines (1 for Yes, 0 for No)
- **Property_Area:** Area of the property (Urban/Semiurban/Rural)
- Loan_Status:** Loan approval status (Y for Yes, N for No)

The primary objective of this dataset is to analyze the factors influencing home loan approval decisions. By performing Exploratory Data Analysis (EDA), we aim to gain insights into the patterns, trends, and relationships within the data. Through data exploration and visualization techniques, we seek to identify key factors affecting loan approval and provide actionable insights for financial institutions and policymakers

Here are some major libraries you'll likely import for working with the home loan approval dataset:

Essential Libraries:

- **pandas (pd):** This is the workhorse library for data analysis in Python. It allows you to:

Load data from various file formats (CSV in this case) using pd.read_csv.
Manipulate and clean data using functions for filtering, sorting, and handling missing values.

Explore and summarize data with descriptive statistics (describe) and data structures like DataFrames.

- **matplotlib.pyplot (plt):** This library provides functionalities for creating various visualizations like:
- **Histograms to explore the distribution of numeric variables.**
- **Box plots to visualize the spread of data and identify outliers.**
- **Scatter plots to explore relationships between pairs of numeric variables.**
- **Bar charts and pie charts to represent categorical data.**

Optional Libraries (depending on your analysis goals):

- **seaborn (sns)**: Built on top of **matplotlib**, **seaborn** offers a high-level interface for creating more aesthetically pleasing and informative visualizations. It provides functions similar to those in **matplotlib** but with a more polished look and advanced features.
- **numpy_(np)_if needed**: While **pandas** is sufficient for most data manipulation tasks in this case, **NumPy** might be useful for specific operations like numerical computations or advanced array manipulations.

Source: Skill Circle

Content: Information related to loan applications and their approval status.

Format: Comma-separated values (CSV) file (assuming based on the code snippet `pd.read_csv`).

Expected Columns:

Applicant Demographics: While the specific details might vary, some possibilities include:

- Applicant Age (numerical)
- Marital Status (categorical - Married, Single, Divorced, etc.)
- Employment Status (optional, categorical)
- Dependents (optional, numerical)

Loan Characteristics:

- **Loan Amount** (numerical, in USD or other currency)
- **Loan Term** (optional, numerical, duration of the loan)
- **Interest Rate** (optional, numerical)

Credit History:

- **Credit Score** (numerical)
- **Delinquency History** (optional, categorical - Yes/No or similar)

Property Details (optional):

- **Location** (categorical - zip code, city, etc.)
- **Property Type** (categorical - house, condo, etc.)
- **Property Value** (optional, numerical)

Loan Approval Status: (categorical)

Approved

Rejected

Pending (optional)

The data is stored in a CSV file named "home_loan_approval.csv", use the following code:

LOADING AND INSPECTING THE DATASET.

Head Function:

Data is loaded into a variable named `data` (as in the previous steps we discussed), you can use the following code to view the first 12 rows:

The head() function is a common method in pandas used to display the first few rows of a DataFrame. By specifying the number 12 within the parenthesis, you instruct it to return only the first 12 rows instead of the default 5.

Employee Data - Sample Data for Employee Table									
EmployeeID	EmployeeName	Gender	Department	Education	Experience	AppraisalsScore	CompletionRate	LastEvaluation	LastAnnualReview
E0001	John Doe	Male	Marketing	Postgraduate	5 Years	85	98%	2023-06-15	2023-06-15
E0002	Jane Smith	Female	Marketing	Postgraduate	5 Years	88	99%	2023-06-15	2023-06-15
E0003	David Lee	Male	Sales	Postgraduate	6 Years	90	100%	2023-06-15	2023-06-15
E0004	Emily Chen	Female	Sales	Postgraduate	6 Years	92	100%	2023-06-15	2023-06-15
E0005	Michael Brown	Male	Customer Service	Postgraduate	5 Years	80	95%	2023-06-15	2023-06-15
E0006	Amy Wilson	Female	Customer Service	Postgraduate	5 Years	82	96%	2023-06-15	2023-06-15
E0007	Robert Johnson	Male	IT Support	Postgraduate	5 Years	78	94%	2023-06-15	2023-06-15
E0008	Samantha Clark	Female	IT Support	Postgraduate	5 Years	80	95%	2023-06-15	2023-06-15
E0009	Christopher Davis	Male	Product Dev.	Postgraduate	6 Years	88	100%	2023-06-15	2023-06-15
E0010	Olivia Parker	Female	Product Dev.	Postgraduate	6 Years	90	100%	2023-06-15	2023-06-15
E0011	Matthew Wilson	Male	Quality Assurance	Postgraduate	5 Years	80	95%	2023-06-15	2023-06-15
E0012	Natalie Parker	Female	Quality Assurance	Postgraduate	5 Years	82	96%	2023-06-15	2023-06-15
E0013	James Miller	Male	Logistics	Postgraduate	6 Years	85	98%	2023-06-15	2023-06-15
E0014	Sarah Johnson	Female	Logistics	Postgraduate	6 Years	87	99%	2023-06-15	2023-06-15
E0015	Matthew Parker	Male	Logistics	Postgraduate	6 Years	89	100%	2023-06-15	2023-06-15
E0016	Natalie Miller	Female	Logistics	Postgraduate	6 Years	91	100%	2023-06-15	2023-06-15
E0017	James Wilson	Male	Logistics	Postgraduate	6 Years	93	100%	2023-06-15	2023-06-15
E0018	Sarah Miller	Female	Logistics	Postgraduate	6 Years	95	100%	2023-06-15	2023-06-15
E0019	Matthew Clark	Male	Logistics	Postgraduate	6 Years	97	100%	2023-06-15	2023-06-15
E0020	Natalie Davis	Female	Logistics	Postgraduate	6 Years	99	100%	2023-06-15	2023-06-15

Tail function:

Use the tail function to view the last 10 rows of the home loan approval dataset:

Similar to `head()`, `tail` is a pandas function for `DataFrames` that retrieves a specific number of rows. In this case, specifying 10 within the parenthesis instructs it to return the last 10 rows.

In [10]:										
A few built-in functions to show 10 rows from the dataset.										
Out[10]:										
Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CosmeticApplnIncome	LoanAmount	Loan_Amount_Term	Class_History
367	Male	No	0	Graduate	No	2000	60	360.0	3.0	NAN
368	Literate	Male	Yes	2	Graduate	3132	70	360.0	3.0	NAN
369	Male	No	0	Graduate	No	4000	2667	76.0	360.0	NAN
370	Male	No	0	Graduate	No	4000	400	96.0	360.0	NAN
371	Male	No	0	Graduate	No	2000	2167	40.0	360.0	NAN
372	Literate	Male	Yes	3+	Graduate	6000	1577	115.0	360.0	1.0
373	Literate	Male	Yes	0	Graduate	6100	700	115.0	360.0	1.0
374	Male	No	0	Graduate	No	3200	1663	120.0	360.0	NAN
375	Male	Yes	0	Graduate	No	5000	2360	160.0	360.0	1.0
376	Male	No	0	Graduate	No	6000	8	96.0	360.0	NAN

Description of the data:

The data describes characteristics of loan applicants and their loan approval status. It includes the following columns:

```
In [102]: # A description of the data.
# df.info()
# df.describe()

Out[102]:
```

	ApplicantIncome	CosplayInfluence	LoanAmount	Loan_Amount_Term	Credit_History
count	367.000000	367.000000	362.000000	361.000000	338.000000
mean	4855.98405	1968.57867	136.152079	342.537396	0.824444
std	4959.03399	2350.00000	81.368000	65.159643	0.380150
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	2984.00000	0.000000	100.250000	360.000000	1.000000
50%	3789.00000	102.000000	125.000000	360.000000	1.000000
75%	5860.00000	2430.500000	158.000000	360.000000	1.000000
max	7230.00000	2400.000000	550.000000	480.000000	1.000000

We can see some summary statistics of the data using the `describe` method. For numerical columns (like Applicant Age and Credit Score), it provides the mean, standard deviation, minimum, quartiles (25% and 75%), and maximum values. For categorical columns (like Marital Status and Loan Approval Status), it shows the number of unique entries and the most frequent entry.

Information about the Dataset:

Data Types:

Numeric:

- Applicant Age (years)
- Credit Score
- Loan Amount (USD)

Categorical:

- Marital Status (e.g., Married, Single, Divorced)
- Loan Approval Status (e.g., Approved, Rejected, Pending)

```
[1]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 12 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Loan_ID         1000 non-null   object  
 1   Marital_Status  1000 non-null   object  
 2   Gender          1000 non-null   object  
 3   Education      1000 non-null   object  
 4   Self_Employed   1000 non-null   object  
 5   ApplicantIncome 1000 non-null   int64  
 6   CoapplicantIncome 1000 non-null   int64  
 7   LoanAmount      1000 non-null   float64 
 8   Loan_Amount_Term 1000 non-null   int64  
 9   Term            1000 non-null   object  
 10  Property_Area   1000 non-null   object  
 11  Loan_Status     1000 non-null   object 
```

Using duplicated() function (pandas):

This function in pandas identifies rows that are entirely duplicates based on all columns. You can specify which columns to consider for duplication using the subset parameter.

```
In [100]: # checking whether there is any duplicated data present in the dataset
df.duplicated().sum()
Out[100]: 0
```

Deleting the Duplicated Data:

```
In [101]: # Deleting the duplicated data
df = df.drop_duplicates()
```

Data Cleaning:

Handling Missing Values:

Identify Missing Values: Utilize the `data.isnull().sum()` function in pandas to identify the number of missing values present in each column. This helps you prioritize which columns to address first.

```
In [102]: # df.isnull().sum()
Out[102]: 
Loan_ID      0
Marital_Status 0
Gender        0
Education     0
Self_Employed 0
ApplicantIncome 0
CoapplicantIncome 0
LoanAmount    0
Loan_Amount_Term 0
Term          0
Property_Area  0
Loan_Status    0
```

This code helps you identify which columns in your DataFrame have the most missing values. This information is crucial for data cleaning tasks like deciding which columns to handle (e.g., impute missing values, drop columns) or prioritizing which columns to address first based on the severity of missing data.

```
In [103]: # checking count of null values
df.isnull().sum().sort_values()
Out[103]: 
loan_id      0
gender       0
education    0
self_employed 0
applicantincome 0
coapplicantincome 0
property_area  0
loan_amount   0
loan_amount_term 0
term          0
self_rejected 21
rejected     21
property_type 20
```

To calculate total missing values:

```
In [109]: -> df[loan].isnull().sum().sum()
Out[109]: 64
```

To examine inconsistencies in the "Loan_Amount" column of your home loan approval dataset:

```
In [110]: -> # Examining inconsistencies in loan_amount column.
          df[loan].head(10)

Out[110]:
   Loan_ID Gender Marital Dependents Education Self_Employed ApplicationIncome CoapplicantIncome LoanAmount Term Credit_History
74  LP00002 Female Yes 0 Graduate No 40000 0 30000 360 1.0
75  LP00003 Female Yes 0 Graduate No 20000 0 20000 360 1.0
76  LP00004 Female Yes 0 Graduate No 10000 0 10000 360 1.0
77  LP00005 Female Yes 0 Graduate No 10000 0 10000 360 1.0
78  LP00006 Female Yes 0 Graduate No 10000 0 10000 360 1.0
79  LP00007 Male Yes 0 Graduate No 10000 0 10000 360 1.0
80  LP00008 Male Yes 0 Graduate No 10000 0 10000 360 1.0
81  LP00009 Male Yes 0 Graduate No 10000 0 10000 360 1.0
82  LP00010 Male Yes 0 Graduate No 10000 0 10000 360 1.0
```

To find and replace outliers across the entire numerical dataset

```
In [111]: -> # Find the outliers from the dataset according to the columns and fill them with the below
          df[loan].isnull().sum()

In [112]: -> # Using Filling Function to fill the null values present in the numeric columns with median of their particular column.
          df['LoanAmount'] = df['LoanAmount'].fillna(df['LoanAmount'].median())
          df['CoApplicantIncome'] = df['CoApplicantIncome'].fillna(df['CoApplicantIncome'].median())
          df['SelfEmployed'] = df['SelfEmployed'].fillna(df['SelfEmployed'].median())
          df['Term'] = df['Term'].fillna(df['Term'].median())
          df['Credit_History'] = df['Credit_History'].fillna(df['Credit_History'].median())

In [113]: -> # Using Filling Function to fill the null values present in the object columns with median of their particular column.
          df['Gender'] = df['Gender'].fillna('Unknown')
          df['Marital_Status'] = df['Marital_Status'].fillna('Unknown')
          df['Education'] = df['Education'].fillna('Unknown')
```

To summarize the basic statistics for the numeric columns in home loan approval dataset:

```
In [108]: -> # Summarize basic statistics for the numeric columns
          df[loan].describe()

Summarize basic statistics for the numeric columns
In [108]:
   Loan_ID      Gender Marital_Status Dependents Education Self_Employed ApplicationIncome CoapplicantIncome LoanAmount Term Credit_History
count  300.000000    300.000000  300.000000  300.000000  300.000000  300.000000  300.000000  300.000000  300.000000  300.000000  300.000000
mean   600.000000    0.500000  0.500000  0.500000  0.500000  0.500000  40000.000000  10000.000000  30000.000000  360.000000  1.000000
std    1000.000000    0.000000  0.000000  0.000000  0.000000  0.000000  40000.000000  10000.000000  30000.000000  360.000000  0.000000
min    100.000000    0.000000  0.000000  0.000000  0.000000  0.000000  10000.000000  10000.000000  10000.000000  360.000000  0.000000
q1    200.000000    0.000000  0.000000  0.000000  0.000000  0.000000  20000.000000  10000.000000  20000.000000  360.000000  0.000000
q3    900.000000    1.000000  1.000000  1.000000  1.000000  1.000000  60000.000000  20000.000000  60000.000000  360.000000  1.000000
max   1000.000000    1.000000  1.000000  1.000000  1.000000  1.000000  100000.000000  20000.000000  100000.000000  360.000000  1.000000

In [109]: -> # Statistics of the applicationIncome
          df[loan]['ApplicationIncome'].mean(), df[loan]['ApplicationIncome'].median(), df[loan]['ApplicationIncome'].std()

Mean of ApplicationIncome is 40000.0
Median of ApplicationIncome is 40000.0
Standard deviation of ApplicationIncome is 40000.000000000004

In [110]: -> # Statistics of the coapplicantIncome
          df[loan]['CoapplicantIncome'].mean(), df[loan]['CoapplicantIncome'].median(), df[loan]['CoapplicantIncome'].std()

Mean of CoapplicantIncome is 10000.0
Median of CoapplicantIncome is 10000.0
Standard deviation of CoapplicantIncome is 10000.000000000002

In [111]: -> # Statistics of the loanAmount
          df[loan]['LoanAmount'].mean(), df[loan]['LoanAmount'].median(), df[loan]['LoanAmount'].std()

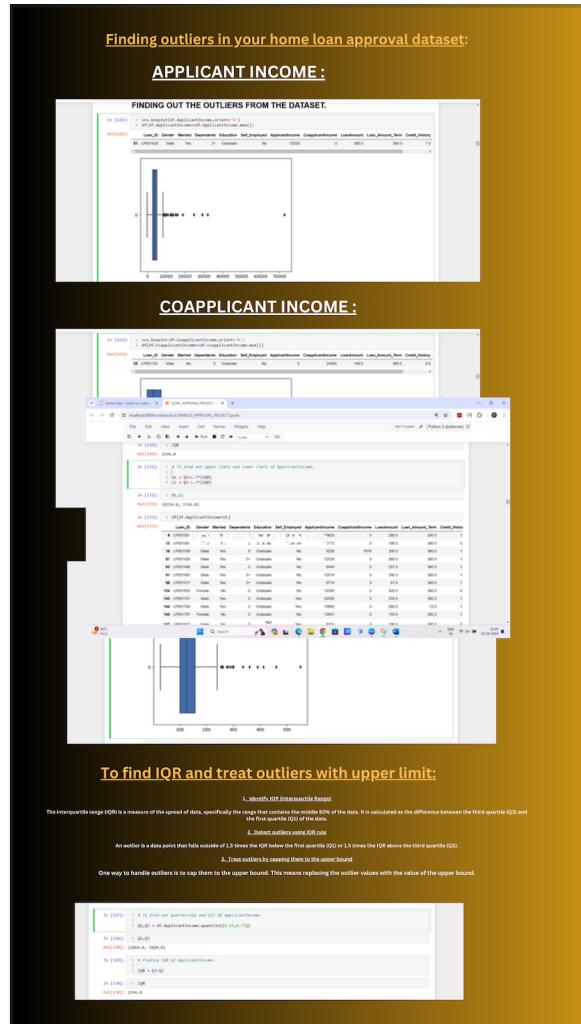
Mean of LoanAmount Term is 360.0
Median of LoanAmount Term is 360.0
Standard deviation of LoanAmount Term is 360.00000000000006

In [112]: -> # Statistics of the term
          df[loan]['Term'].mean(), df[loan]['Term'].median(), df[loan]['Term'].std()

Mean of Term is 360.0
Median of Term is 360.0
Standard deviation of Term is 0.0000000000000002

In [113]: -> # Statistics of the creditHistory
          df[loan]['Credit_History'].mean(), df[loan]['Credit_History'].median(), df[loan]['Credit_History'].std()

Mean of Credit_History is 1.0
Median of Credit_History is 1.0
Standard deviation of Credit_History is 0.17320508075688773
```



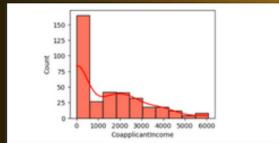
OUTLIERS REMOVED SUCCESSFULLY IN ALL THE COLUMNS:

DATA VISUALIZATION:

1. UNI-VARIATE ANALYSIS:

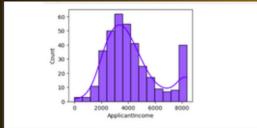
- Representation of data by using Histogram

DISTRIBUTION OF COAPPLIANT INCOME



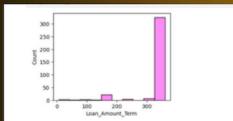
This histogram illustrate the distribution of Coapplicant income. The x-axis represents "Coapplicant Income"(ranging from 0 to 6000),while the y-axis shows "count"(ranging from 0 to 150).A Prominent bar reaches up to approximately 150 count for incomes close to zero,indicating many coapplicants have little to no income.The frequency of higher incomes diminishes,with smaller bars representing fewer coapplicants in those income brackets. overall, this graph provides insights into the financial standing of coapplicants.

DISTRIBUTION OF APPLICANT INCOME



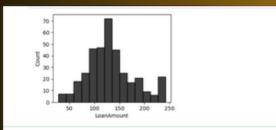
- Distribution of Applicant Income: The histogram provides a visual representation of the distribution of applicant incomes in the dataset. It shows the frequency or count of applicants falling within different income ranges.
- Central Tendency: The shape of the histogram can indicate the central tendency of applicant incomes. For example, a peak around a specific income range may suggest a common income level among applicants.
- Variability: The spread of the histogram can give insights into the variability of applicant incomes. A wider spread indicates a broader range of income levels among applicants.
- Skewness: The symmetry or skewness of the histogram can reveal whether the distribution of applicant incomes is skewed to the left (negatively skewed), to the right (positively skewed), or symmetric.
- The image shows a histogram overlaid with a line graph, illustrating the distribution of "Applicant Income" (ranging from 0 to 8000), while the y-axis shows count (ranging from 0 to 60). Most applicants fall within the income range of 2000 to 4000, with an interesting spike at higher incomes. The graph provides insights into the income distribution among a specific group.

DISTRIBUTION OF LOAN AMOUNT TERM



The image is a bar chart showing the number of loans issued based on the loan term. The x-axis labeled on "Loan_Amount_Term" likely represents the length of the loan in years. The y-axis labeled "count" represents the number of loans. It appears that shorter loan terms are more common than longer terms.

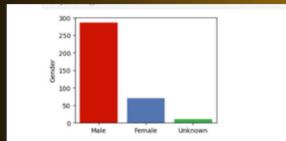
DISTRIBUTION OF LOAN AMOUNT



This is a bar graph displaying the distribution of loan amounts. The x-axis is labeled 'LoanAmount' and ranges from 0 to 250, while y-axis is labeled 'count' and ranges from 0 to 70. A prominent bar reaches up to approximately 150 count for incomes close to zero, indicating that many coapplicants have little to no income. The frequency of higher incomes diminishes, with smaller bars representing fewer coapplicants in those income brackets. Overall this graph provides insights into the financial standing of coapplicants.

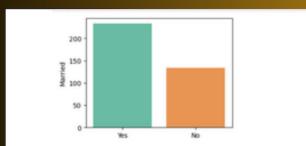
BAR CHARTS:

DISTRIBUTION OF GENDER:



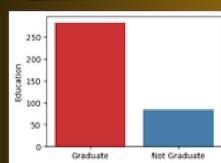
This visualization reveals the distribution of genders among the loan applicants . It appears that there is a noticeable imbalance between male, female and unknown applicants . This could be due to various socio-economic factors influencing the willingness of individuals to apply for loans.

DISTRIBUTION OF MARRIED:



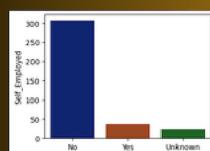
This bar plot visualization the distribution of applicants based on their marital status, with categories such as 'Married' and 'Not Married'.There is a noticeable imbalance in the distribution, there are more married applicants comapred to unmarried ones. It could reflect societal trends or demographic characteristics of the dataset.

DISTRIBUTION OF EDUCATION:



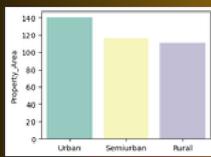
This count plot visualization the distribution of applicants based on their education status, with categories such as 'Educate' and 'Not Educate'.There is a noticeable imbalance in the distribution, there are more graduate applicants comapred to uneducated ones.

DISTRIBUTION OF SELF-EMPLOYED:



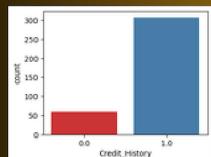
- This count plot illustrates the distribution of applicants based on their self employment status. Categories may include 'Yes'(indicating not self employed) and 'No'(indicating self employed).There is a significant count in the 'Yes' category,it indicates a notable presence of self-employed individuals or entrepreneurs among the loan applicants. Self employed individuals may have variable income structures compared to salaried individuals.Understanding their distribution provides insights into the diversity of income sources among loan applicant.

DISTRIBUTION OF PROPERTY AREA:



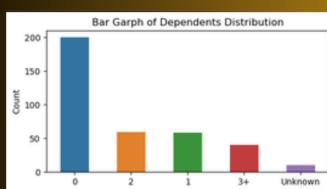
This count plot visualization the distribution of applicants based on their Propert_Area status, with categories such as 'Urban','Semirurban' and 'Rural'.There is a noticeable imbalance in the distribution, there are most Urban then Semirurban applicants comapred Rural ones.

DISTRIBUTION OF CREDIT HISTORY:



This count plot visualization the distribution of applicants based on their credit_history status, with categories such as '0.0' and '1.0'.There is a noticeable imbalance in the distribution, there are more 1.0 applicants comapred to 0.0 ones.

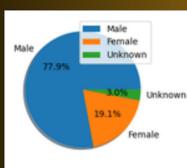
DISTRIBUTION OF DEPENDENTS:



The resulting bar graph provides a clear visualization of the distribution of dependents, allowing users to quickly identify the most common number of dependents (0.0, with a count of 200) and the least common numbers (2.0 and 1.0,3+,unknownwith counts of 59 ,58,40 and 10 respectively).

PIE CHARTS:

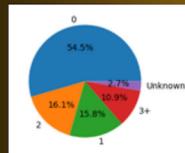
DISTRIBUTION OF GENDER:



- The pie chart illustrates the distribution of genders within the dataset.
- Each slice represents a unique gender category, with sizes proportional to the frequency of each gender.
- The distribution of slices are as male have 77.9%, female have 19.1% and other unknown have 3.0% respectively.

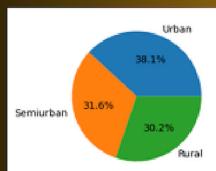
Gender Representation: The chart visually displays the relative proportions of different genders.
- Percentage Breakdown: The percentages on each slice indicate the proportion of each gender category in the dataset.
- Imbalance Detection: Shows for any significant dominant gender categories.

DISTRIBUTION OF DEPENDENTS:



This pie chart represent slices of dependent distribution as
In 0 have 54.5%
In 2 have 16.1%
In 1 have 15.8%
In 3+ have 10.9%
and other unknown have 2.7%

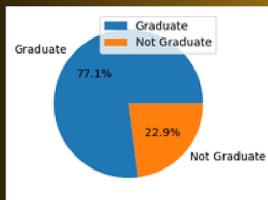
DISTRIBUTION OF PROPERTY AREA:



The pie chart illustrates the distribution of Property_Area within the dataset.
 Each slice represents a unique Property_Area category, with sizes proportional to the frequency of each Property_Area.
 Property_Area Representation: The chart visually displays the relative proportions of different Property_Area.
 Percentage Breakdown: The percentages on each slice indicate the proportion of each Property_Area category in the dataset.
 Imbalance Detection: Shows for any significant dominant Property_Area categories.
 The following distribution of Property_Area slices are as:

Urban = 38.1%
 Rural = 30.2%
 Semiurban = 31.6%

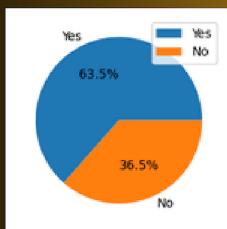
DISTRIBUTION OF EDUCATION:



The pie chart illustrates the distribution Of education are as :

- There are 77.1% Graduates.
- There are 22.9% Non-Graduates.

DISTRIBUTION OF MARRIED:



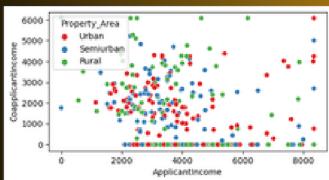
The pie chart illustrates the distribution of Married's are as follows:

- There are 63.5% Married.
- There are 36.5% non-Married.

BI - VARIATE ANALYSIS:

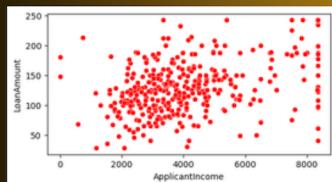
SCATTER PLOTS:

SCATTER PLOT OF COAPPLICANTINCOME_BY APPLICANT INCOME_BY PROPERTY AREA:



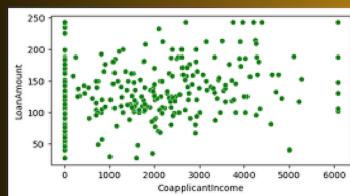
This scatter plot illustrates the x-axis of applicantincome vs y-axis of coapplicantincome by property area(rural,urban,semiurban)respectively.

DISTRIBUTION OF APPLICANTINCOME AND LOAN AMOUNT:



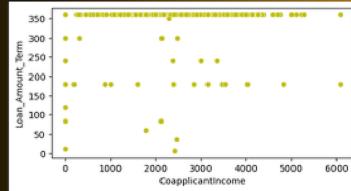
There is a positive correlation between applicantincome and loan amount. This means that as the applicant's income increases, they may be eligible for larger loan amounts. Applicants with higher incomes may qualify for larger loans because they have financial capacity to repay them.

DISTRIBUTION OF COAPPLICANTINCOME AND LOAN AMOUNT:



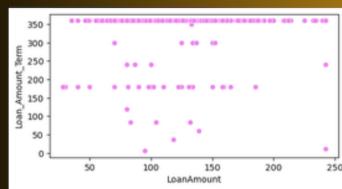
This is a scatter plot graph displaying the relationship between Coapplicant Income and LoanAmount. The x-axis represents "CoapplicantIncome" (ranging from 0 to 6000), while the y-axis shows LoanAmount (ranging from 0 to 250). Numerous green dots represent the individual data points scattered across the graph, indicating various combinations of coapplicant income and loan amounts. A scatter plot illustrating the correlation between coapplicant income and loan amount, providing insights into trends or patterns.

DISTRIBUTION OF COAPPLICANTINCOME AND LOAN AMOUNT TERM:



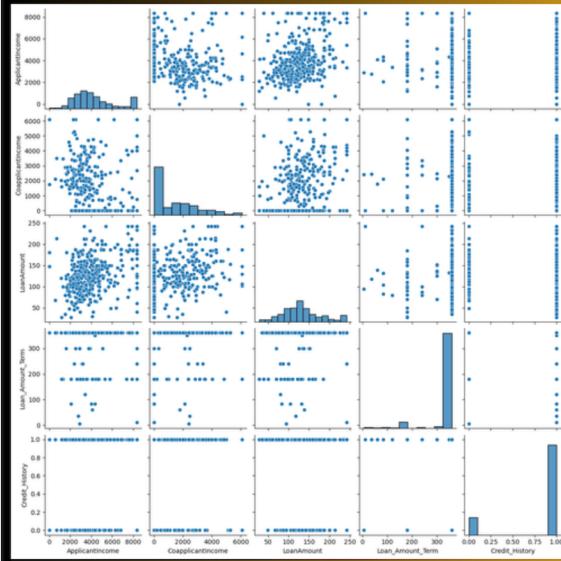
This scatter graph displaying the relationship between coapplicantincome and loan amount term. The x-axis have coapplicantincome (range from 0 to 6000)and y-axis have Loan_Amount_Term(range from 0 to 350).

DISTRIBUTION OF LOANAMOUNT AND LOAN_AMOUNT TERM:



This scatter graph displaying the relationship between LoanAmount and Loan_Amount_Term. The x-axis have LoanAmount (range from 0 to 250)and y-axis have Loan_Amount_Term(range from 0 to 350)respectively.

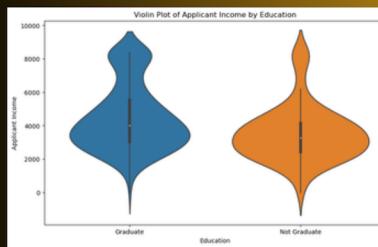
PAIR PLOT:



The Pairplot analysis reveals a positive correlation between ApplicantIncome and LoanAmount, suggesting that higher income leads to a higher loan eligibility.

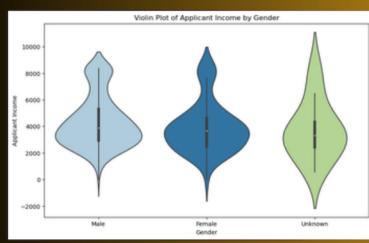
VIOLIN PLOTS:

PLOT OF APPLICANT INCOME BY EDUCATION:



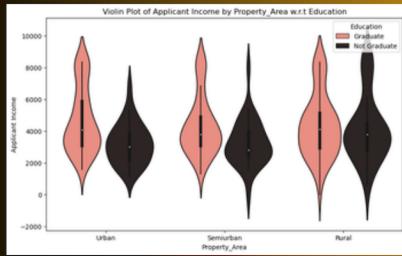
-This violin plot showing income distribution for applicants with a graduate degree compared to those without. The vertical axis shows income, while the horizontal axis represents education level. The wider the violin shape, the greater the variation in income. The plot suggests that applicants have higher incomes than those without a graduate degree.

PLOT OF APPLICANT INCOME BY GENDER:



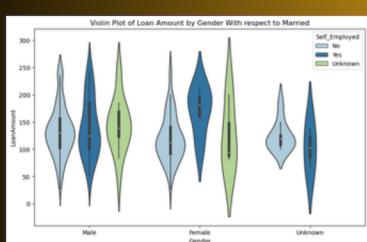
The violin plot provides a good starting point to compare the distribution of applicant income by gender and identify potential areas for further investigation.

PLOT OF APPLICANTINCOME BY PROPERTY AREA
w.r.t EDUCATION:



- The violin plot will display the distribution of Applicant Income for each combination of Property_Area and Education. This will allow you to visually analyze how the Applicant Income varies across different property areas and educational backgrounds.
The x-axis represents the Property_Area , y-axis represents the Applicant Income. The "hue" parameter is used to differentiate the Education levels, resulting in separate violin plots for each education group.
The shape of the violin plot represents the density distribution of the Applicant Income for each group.
This plot helps identify patterns and relationships between Applicant Income, Property_Area, and Education.

PLOT OF LOAN AMOUNT BY GENDER w.r.t MARRIED:



The x-axis represents the Gender categories (e.g., Male, Female)and the y-axis represents the Loan Amount.
The "hue" parameter is used to differentiate the Married status, resulting in separate violin plots for each Married/Not Married group.
The shape of the violin plot represents the density distribution of the Loan Amount for each group.
This plot helps identify patterns and relationships between Loan Amount, Gender, and Marital Status.

MULTI - VARIATE ANALYSIS:

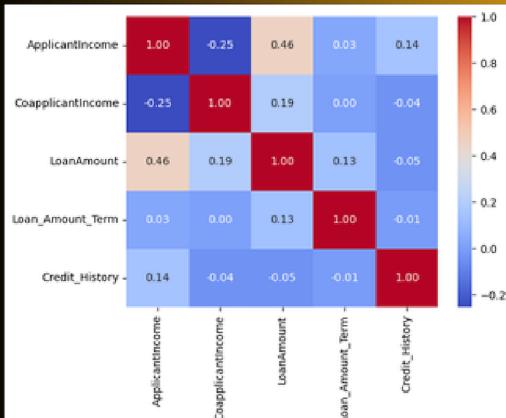
HEATMAP USING CORRELATION:

```
In [114]: 1 df.corr(numeric_only=True)
```

```
Out[114]:
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
ApplicantIncome	1.000000	-0.25053	0.49722	0.029694	0.143344
CoapplicantIncome	-0.25053	1.000000	0.194046	0.000037	-0.030303
LoanAmount	0.49722	0.194046	1.000000	0.133376	-0.04532
Loan_Amount_Term	0.029694	0.000037	0.133376	1.000000	-0.011347
Credit_History	0.143344	-0.030303	-0.04532	-0.011347	1.000000

```
In [115]: 1 sns.heatmap(df.corr(numeric_only=True), annot=True, cmap="coolwarm", fmt=".2f")
```



Correlation Matrix Insights

Relationships between Variables

ApplicantIncome and

CoapplicantIncome: Weak negative correlation (-0.254).

LoanAmount and ApplicantIncome:

Moderate positive correlation (0.486).

LoanAmount and CoapplicantIncome:

Weak positive correlation (0.171).

LoanAmount and Loan_Amount_Term:

Weak positive correlation (0.092).

Credit_History and ApplicantIncome:

Weak positive correlation (0.43).

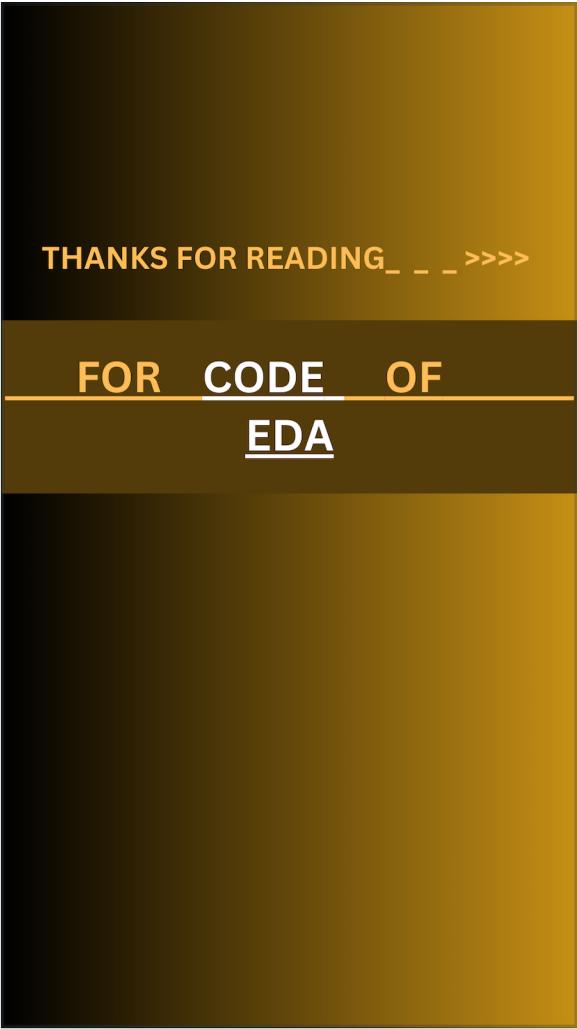
This correlation heatmap provides a visual representation of the relationships between the variables in the dataset. The color intensity and the numerical values in the cells indicate the strength and direction of the correlations. This can be a useful tool for identifying patterns and relationships within the data.

CONCLUSION:

Our exploration of the dataset provided valuable insights into loan approval factors. We loaded the dataset into a Python environment, examined its structure, and handled missing values appropriately. Basic statistics for numeric columns were summarized.

In terms of data visualization, we delved into univariate, bivariate, and multivariate analyses. Histograms, box plots, bar charts, and pie charts revealed distribution patterns, relationships between variables, and correlations.

As we move forward, consider focusing on specific features impacting loan approval rates.



THANKS FOR READING_ _ _ >>>

**FOR CODE OF
EDA**