

##Name - Anjali N

##Project name - Loan Status Prediction

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
df=pd.read_csv("loan.csv")
df.head()
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
0	LP001002	Male	No	0	Graduate	No	5849	0.0	NaN	360.0	1.0	Urban	Y
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	1.0	Rural	N
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0	1.0	Urban	Y
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0	1.0	Urban	Y
4	LP001008	Male	No	0	Graduate	No	6000	0.0	141.0	360.0	1.0	Urban	Y

```
In [2]: df.tail()
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
609	LP002070	Female	No	0	Graduate	No	2200	0.0	71.0	360.0	1.0	Rural	Y
610	LP002071	Male	Yes	3+	Graduate	No	4106	0.0	40.0	180.0	1.0	Rural	Y
611	LP002091	Male	Yes	1	Graduate	No	8072	240.0	253.0	360.0	1.0	Urban	Y
612	LP002084	Male	Yes	2	Graduate	No	7583	0.0	187.0	360.0	1.0	Urban	Y
613	LP002090	Female	No	0	Graduate	Yes	4553	0.0	123.0	360.0	0.0	Semurban	N

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

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.00000	564.000000
mean	5403.459283	1821.245798	146.412162	342.00000	0.842199
std	6109.041673	2926.248369	85.837325	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.00000	0.000000
25%	2877.500000	0.000000	100.250000	360.00000	1.000000
50%	3812.500000	1188.500000	128.000000	360.00000	1.000000
75%	5795.000000	2297.250000	164.750000	360.00000	1.000000
max	81000.000000	41667.000000	700.000000	480.00000	1.000000

```
In [4]: df['Property_Area'].value_counts()
```

```
Name: Gender, dtype: int64
```

```
df['Married'].value_counts()
```

```
Yes    398
No     213
```

```
In [5]: df['Loan_Status'].value_counts()
```

Y	422
N	192
Name: Loan_Status, dtype: int64	

```
In [6]: df['Self_Employed'].value_counts()
```

No	599
Yes	82
Name: Self_Employed, dtype: int64	

```
In [7]: df['Gender'].value_counts()
```

Male	419
Female	122
Name: Gender, dtype: int64	

```
In [8]: df['Married'].value_counts()
```

Yes	398
No	213
Name: Married, dtype: int64	

```
In [9]: df['Education'].value_counts()
```

Graduate	489
Not Graduate	134
Name: Education, dtype: int64	

```
In [10]: df['ApplicantIncome'].hist(bins=50)
```

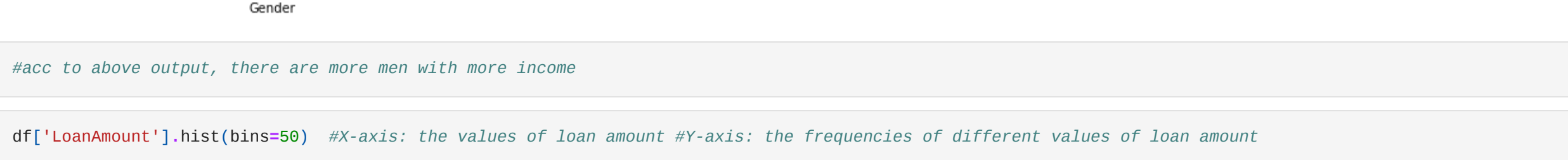


```
In [11]: df.boxplot(column='ApplicantIncome', by='Gender')
```



##acc to above output, there are more men with more income

```
In [13]: df['LoanAmount'].hist(bins=50) ##x-axis: the values of loan amount ##y-axis: the frequencies of different values of loan amount
```



```
In [14]: df['Credit_History'].value_counts()
```

1.0	475
0.0	89
Name: Credit_History, dtype: int64	

##creation of pivot tables using python

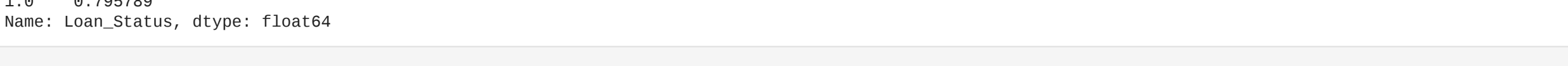
```
pt = df.pivot_table(values='Loan_Status', index=['Credit_History'], aggfunc=lambda x: x.map({'Y':1, 'N':0}), margin=True)
```

Credit_History	0.0	8.978652
1.0	6.785789	
Name: Loan_Status, dtype: float64		

```
fig = plt.figure(figsize=(8,4))
ax1 = fig.add_subplot(121)
```

```
ax1.set_xlabel('Credit_History')
ax1.set_ylabel('Count of Applicants')
ax1.set_title('Applicants by credit_History')
ch.plot(kind='bar')
```

```
ax2 = fig.add_subplot(122)
ax2.set_xlabel('Credit_History')
ax2.set_ylabel('Probability of getting loan')
ax2.set_title('Probability of getting loan by credit history')
pt.plot(kind='bar')
```

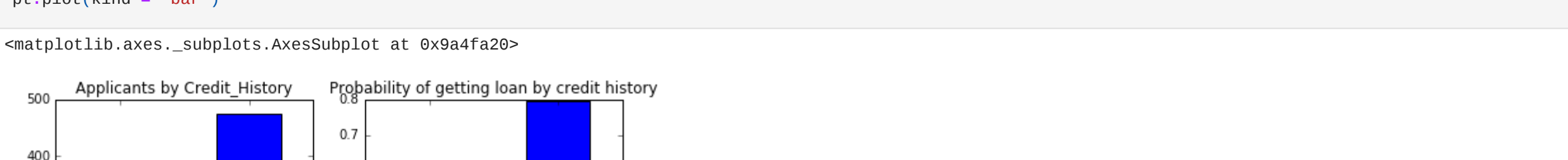


```
In [18]: pv = df.pivot_table(values='Gender', index=['Credit_History'], aggfunc=lambda x: x.map({'Female':1, 'Male':0}), margin=True)
```

Credit_History	0.0	6.197674
1.0	6.185258	
Name: Gender, dtype: float64		

##combining the above 2 plots in a stacked chart

```
stack_charts = pd.crosstab([df['Credit_History'], df['Loan_Status']],
stack_charts.plot(kind='bar', stacked=True, color=['green', 'yellow'], grid=False)
stack_charts = pd.crosstab([df['Gender'], df['Loan_Status']],
stack_charts.plot(kind='bar', stacked=True, color=['green', 'yellow'], grid=False)
```



```
In [20]: df.apply(lambda x: sum(x.isnull()), axis=0) #to tell the missing values in each column
```

Loan_ID	0
Gender	13
Married	3
Dependents	15
Education	0
Self_Employed	32
ApplicantIncome	2
CoapplicantIncome	0
LoanAmount	22
Loan_Amount_Term	14
Credit_History	58
Property_Area	0
Loan_Status	0
dtype: int64	

##to fill the missing values by mean

```
LoanAmount = df.LoanAmount.fillna(df.LoanAmount.mean())
df.describe()
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	614.000000	600.00000	564.000000
mean	5403.459283	1821.245798	146.412162	342.00000	0.842199
std	6109.041673	2926.248369	84.037468	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.00000	0.000000
25%	2877.500000	0.000000	100.250000	360.00000	1.000000
50%	3812.500000	1188.500000	128.000000	360.00000	1.000000
75%	5795.000000	2297.250000	164.750000	360.00000	1.000000
max	81000.000000	41667.000000	700.000000	480.00000	1.000000

```
In [22]: se = df[['Self_Employed']].value_counts(ascending=True)
```

se	head()
Yes	82
No	599
Name: Self_Employed, dtype: int64	

##since 'no' has ~82% so, we will fill the missing values of this column with "no"

```
df['Self_Employed'] = df['Self_Employed'].fillna('No')
df.describe()
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	614.000000	600.00000	564.000000
mean	5403.459283	1821.245798	146.412162	342.00000	0.842199
std	6109.041673	2926.248369	84.037468	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.00000	0.000000
25%	2877.500000	0.000000	100.250000	360.00000	1.000000
50%	3812.500000	1188.500000	128.000000	360.00000	1.000000
75%	5795.000000	2297.250000	164.750000	360.00000	1.000000
max	81000.000000	41667.000000	700.000000	480.00000	1.000000

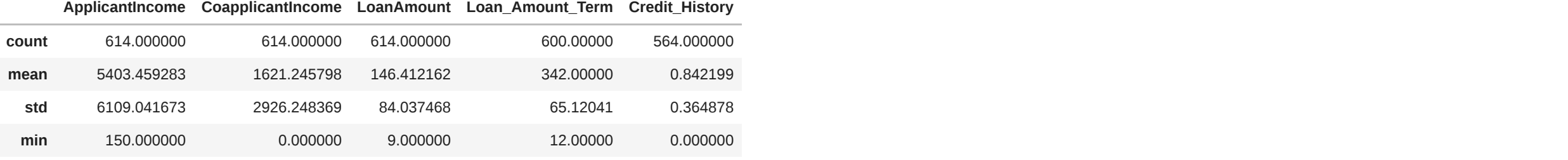
```
In [28]: df['LoanAmount_log'] = np.log(df['LoanAmount']) #log function to reduce the extreme effect in loan amount (which we saw in the previous barplot)
```

```
df['LoanAmount_log'].hist(bins=20)
```



##to decrease the extreme values of ApplicantIncome, we can add CoApplicantIncome to compensate the value

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



```
In [30]: df['Capacity'] = ((df['LoanAmount']/ df['TotalIncome']) * 100).astype(float) #capacity of each applicant of how well he/she is suited to pay back his loan.
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status
0	LP001002	Male	No	0	Graduate	No	5849	0.0	146.412162	360.0	1.0	Urban	Y
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.000000	360.0	1.0	Rural	N
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.000000	360.0	1.0	Urban	Y
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.000000	360.0	1.0	Urban	Y
4	LP001008	Male	No	0	Graduate	No	6000	0.0	141.000000	360.0	1.0	Urban	Y
5	LP001011	Male	Yes	2	Graduate	Yes	5417	4196.0	267.000000	360.0	1.0	Urban	Y
6	LP001013	Male	Yes	0	Not Graduate	No	2333	1516.0	95.000000	360.0	1.0	Urban	Y
7	LP001014	Male	Yes	3+	Graduate	No	3036	2504.0	150.000000	360.0	0.0	Semurban	N
8	LP001018	Male	Yes	2	Graduate	No	4006	1525.0	168.000000	360.0	1.0	Urban	Y
9	LP001021	Male	Yes	1	Graduate	No	12841	10980.0	349.000000	360.0	1.0	Semurban	N
10	LP001024	Male	Yes	2	Graduate	No	3200	700.0	70.000000	360.0	1.0	Urban	Y
11	LP001027	Male	Yes	2	Graduate	No	2500	1840.0	109.000000	360.0	1.0	Urban	Y
12	LP001028	Male	Yes	2	Graduate	No	3073	8106.0	200.000000	360.0	1.0	Urban	Y
13	LP001029	Male	No	0	Graduate	No	1853	2840.0	114.000000	360.0	1.0	Rural	N
14	LP001030	Male	Yes	2	Graduate	No	1299	1086.0	17.000000	120.0	1.0	Urban	Y
15	LP001032	Male	No	0	Graduate	No	4590	0.0	125.000000	360.0	1.0	Urban	Y
16	LP001034	Male	No	1	Not Graduate	No	2696	0.0	100.000000	240.0	NaN	Urban	Y
17	LP001036	Female	No	0	Graduate	No	3510	0.0	76.000000	360.0	0.0	Urban	N
18	LP001038	Male	Yes	0	Not Graduate	No	4887	0.0	133.000000	360.0	1.0	Rural	N
19	LP001041	Male	Yes	0	Graduate	No	2600	3500.0	115.000000	360.0	1.0	Urban	Y

##males are around ~81%, so we'll fill the missing values with "male"

```
df['Gender'] = df['Gender'].fillna('Male')
```

```
gen = df['Gender'].value_counts(ascending=True)
gen.head()
```

Female	112
Male	582
Name: Gender, dtype: int64	

##which gender has more loan amount

```
fig = plt.figure(figsize=(8,4))
ax1 = fig.add_subplot(121)
```

```
ax1.set_xlabel('Gender')
ax1.set_ylabel('LoanAmount')
gen.plot(kind='bar')
```



```
In [34]: md = df[['Married']].value_counts(ascending=True)
```

md	head()
No	213
Yes	398
Name: Married, dtype: int64	

##yes has ~55%

```
##filling the missing values with
df['Married'] = df['Married'].fillna('Yes')
```

```
fig = plt.figure(figsize=(10,5))
ax1 = fig.add_subplot(121)
```

```
ax1.set_xlabel('Married')
ax1.set_ylabel('LoanAmount')
md.plot(kind='bar')
```



```
In [41]: loan = df['Loan_Status'].value_counts(ascending=True)
loan.head()
```

0	192
1	422
Name: Loan_Status, dtype: int64	

```
In [43]: df['LoanAmount_log'] = df.LoanAmount.fillna('1.0')
df.describe()
```

A histogram showing the distribution of 'Capacity%' values. The x-axis ranges from 2 to 7, and the y-axis ranges from 0 to 140. The distribution is unimodal and slightly right-skewed, with a peak frequency of approximately 140 at a capacity of 4.8%.

Figure 1: Histogram of Capacity% showing the distribution of capacity values across the dataset.

The following code calculates the capacity of each applicant based on their loan amount and total income, and displays the first 20 results.

```
df['Capacity%'] = ((df['LoanAmount'] / df['TotalIncome']) * 100).astype(float) #capacity of each applicant of how well he/she is suited to pay back his loan.  
df.head(20)
```