

Task: Determine the eligibility for granting a Home Loan

Objective of this notebook is:

1. To understand the patterns in the data.
2. How to handle the categorical features.
3. How to deal with missing data.
4. Feature Engineering
5. Finding the most important features while taking the decision of granting a loan application.
6. Understanding the Normalization and standardization of the data.

Importing Libraries

```
In [2]: import numpy as np
import pandas as pd
from scipy import stats

import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import ttest_ind, chi2_contingency, shapiro
from statsmodels.graphics.gcfplots import qqplot
from scipy.stats import kstest, norm
```

```
In [3]: df = pd.read_csv('loan.csv')
df.head()
```

```
Out[3]:
```

	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

Log Normal Distribution

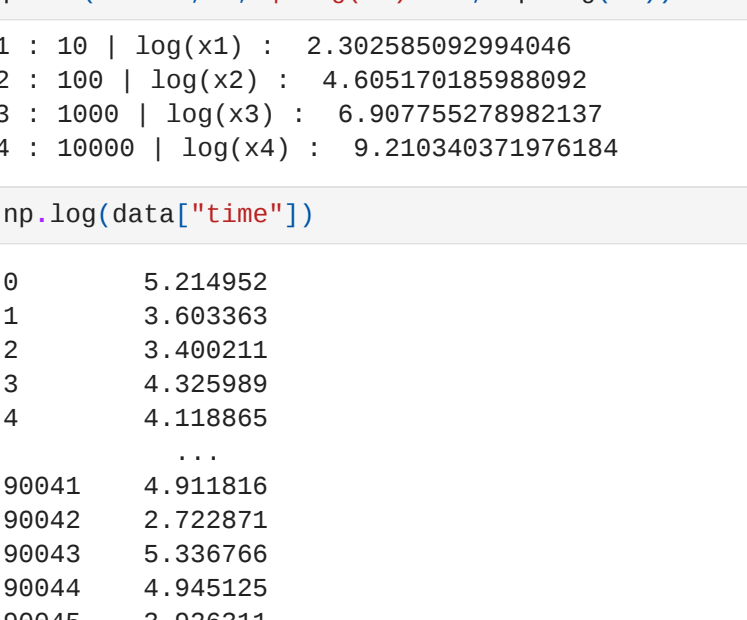
```
In [4]: data=pd.read_csv("wailog_time.csv")
data
```

```
Out[4]:
```

	time
0	184.003076
1	36.721321
2	28.970417
3	75.640285
4	61.489439
...	...
99041	126.889964
99042	16.222970
99043	20.859528
99044	140.488418
99045	50.713544
...	...
99046	706.104005

```
In [5]: sns.histplot(data["time"])
```

```
Out[5]: <AxesSubplot: xlabel='time', ylabel='Count'>
```



```
In [7]: x1=10
x2=100
x3=1000
x4=10000

In [12]: print("x1 :",x1," log(x1) :", np.log(x1))
print("x2 :",x2," log(x2) :", np.log(x2))
print("x3 :",x3," log(x3) :", np.log(x3))
print("x4 :",x4," log(x4) :", np.log(x4))

x1: 10 log(x1): 2.302585092994046
x2: 100 log(x2): 4.605170185988092
x3: 1000 log(x3): 6.907755278982137
x4: 10000 log(x4): 9.210340371976184
```

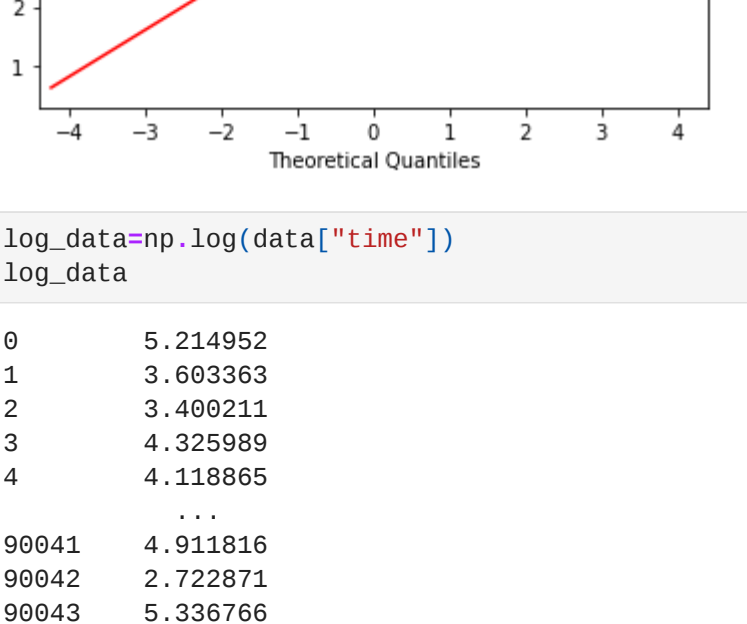
```
In [14]: np.log(data["time"])

Out[14]:
```

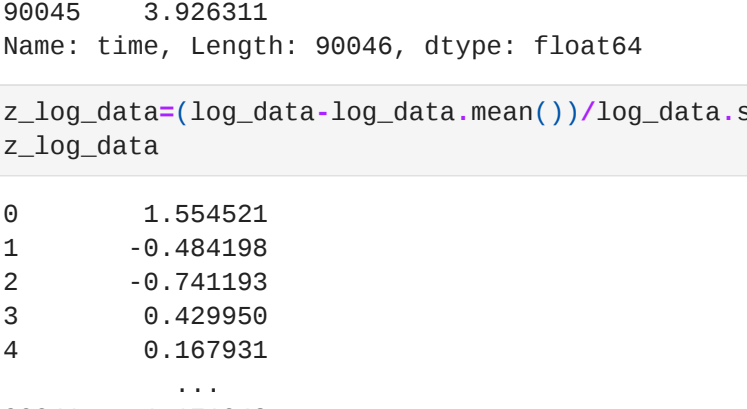
0	5.214952
1	3.603363
2	3.460211
3	4.325989
4	4.118865
...	...
99041	4.911816
99042	2.722871
99043	5.336766
99044	4.945125
99045	3.926311
Name: time, Length: 99046, dtype: float64	

```
In [15]: sns.histplot(np.log(data["time"]))
```

```
Out[15]: <AxesSubplot: xlabel='time', ylabel='Count'>
```



```
In [18]: qqplot(np.log(data["time"]),line='s')
plt.show()
```



```
In [19]: log_data=np.log(data["time"])
log_data

Out[19]:
```

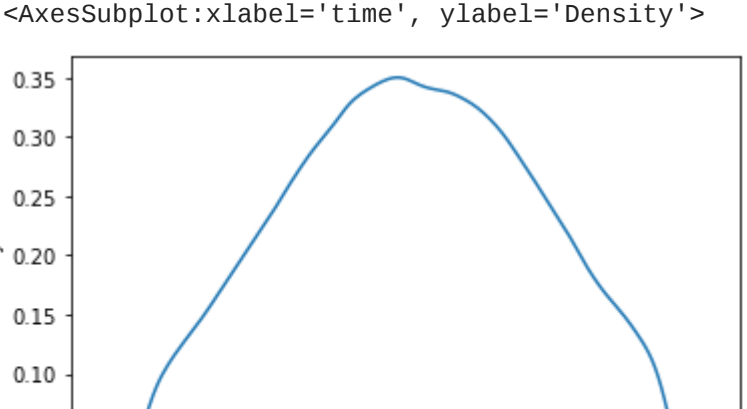
0	5.214952
1	3.603363
2	3.460211
3	4.325989
4	4.118865
...	...
99041	4.911816
99042	2.722871
99043	5.336766
99044	4.945125
99045	3.926311
Name: time, Length: 99046, dtype: float64	

```
In [20]: z_log_data=(log_data-log_data.mean())/log_data.std()
z_log_data

Out[20]:
```

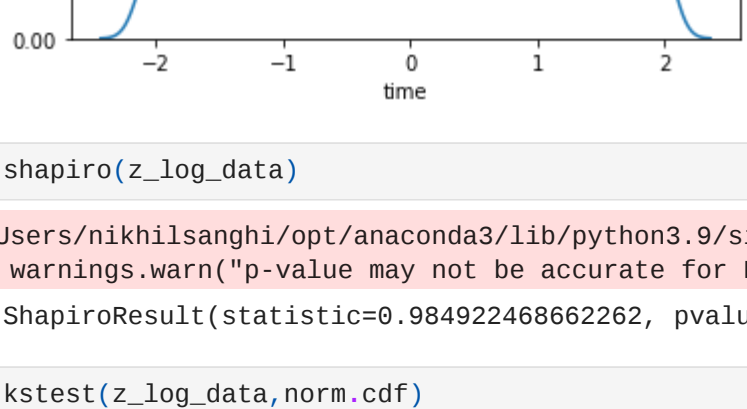
0	1.564521
1	-0.484138
2	-0.741393
3	0.429058
4	0.167931
...	...
99041	1.171943
99042	-2.199852
99043	1.708829
99044	1.121386
99045	-0.755657
Name: time, Length: 99046, dtype: float64	

```
In [22]: qqplot(z_log_data,line='s')
plt.show()
```



```
In [24]: sns.kdeplot(z_log_data)

Out[24]: <AxesSubplot: xlabel='time', ylabel='Density'>
```



```
In [25]: shapiro(z_log_data)

/Users/nikhilshampi/opt/anaconda3/lib/python3.8/site-packages/scipy/stats/_morestats.py:1760: UserWarning: p-value may not be accurate for N > 5000.
warnings.warn("p-value may not be accurate for N > 5000.")

Out[25]: ShapiroResult(statistic=0.984922468662262, pvalue=0.6)
```

```
In [26]: kstest(z_log_data,norm.cdf)

Out[26]: KstestResult(statistic=0.026341647599834476, pvalue=1.034118663995682e-54)
```

Basic Exploration

```
In [27]: df.head()
```

```
Out[27]:
```

	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 [28]: df.info()

Out[28]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 634 entries, 0 to 613
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   Loan_ID               634 non-null   object  
 1   Gender                603 non-null   object  
 2   Married               613 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   float64  
 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

In [29]: df["Dependents"].value_counts()

Out[29]:
0    345
1    182
2    181
3     61
Name: Dependents, dtype: int64

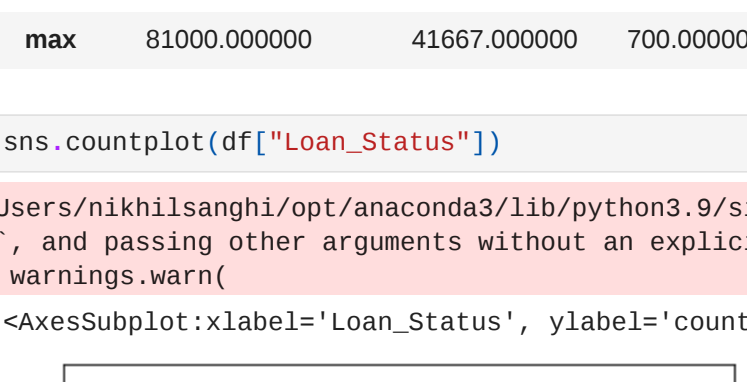
In [30]: df.describe()

Out[30]:
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
count	614.000000	614.000000	592.000000	600.000000	564.000000
mean	5403.462693	1821.245798	146.412162	342.000000	0.848199
std	6109.041673	2926.248369	85.987325	65.12041	0.364878
min	150.000000	0.000000	9.000000	12.000000	0.000000
25%	2877.500000	0.000000	100.000000	360.000000	1.000000
50%	3812.500000	1198.500000	128.000000	360.000000	1.000000
75%	5795.000000	2297.500000	168.000000	360.000000	1.000000
max	81000.000000	41687.000000	700.000000	480.000000	1.000000

```
In [31]: sns.countplot(df["Loan_Status"])

Out[31]: <AxesSubplot: xlabel='Loan_Status', ylabel='count'>
```




Univariate Analysis

```
In [32]: "Applicant_Income"
```

```
In [32]: "Applicant_Income"
```

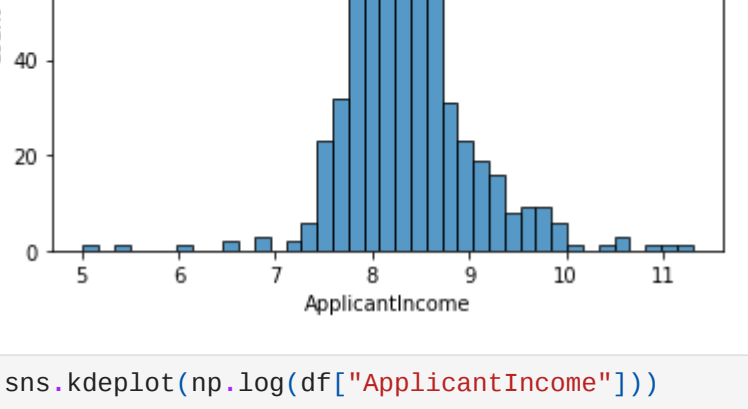
```
In [34]: sns.histplot(df["ApplicantIncome"])

Out[34]: <AxesSubplot: xlabel='ApplicantIncome', ylabel='Count'>
```



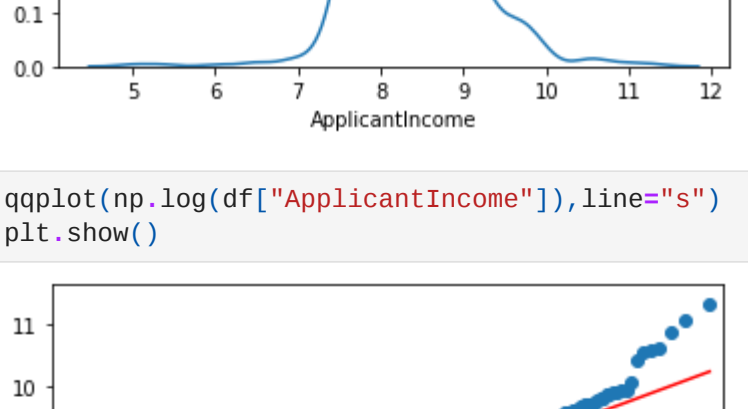
```
In [35]: sns.histplot(np.log(df["ApplicantIncome"]))

Out[35]: <AxesSubplot: xlabel='ApplicantIncome', ylabel='Count'>
```

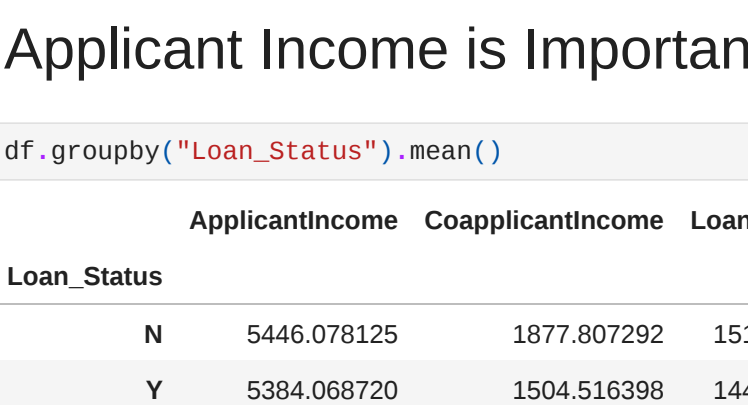


```
In [36]: sns.kdeplot(np.log(df["ApplicantIncome"]))

Out[36]: <AxesSubplot: xlabel='ApplicantIncome', ylabel='Density'>
```



```
In [38]: qqplot(np.log(df["ApplicantIncome"]),line='s')
plt.show()
```



Applicant Income is Important/ Good Predictor?

```
In [42]: df.groupby("Loan_Status").mean()
```

```
Out[42]:
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
Loan_Status					
Y	5446.078125	1877.807292	151.220994	344.064516	0.541899
N	5384.068720	1504.516398	144.294404	341.072664	0.981818

```
In [ ]:
```

```
In [ ]:
```

```
In [40]: df_acc=df.loc[df["Loan_Status"]=="Y"]["ApplicantIncome"]
df_rej=df.loc[df["Loan_Status"]=="N"]["ApplicantIncome"]
```

```
In [50]: df_acc.mean()
```

```
Out[50]: 5384.06872079147
```

```
In [52]: df_rej.mean()
```

```
Out[52]: 5446.078125
```

```
In [53]: ttest_ind(df_acc,df_rej,alternative="less")
```

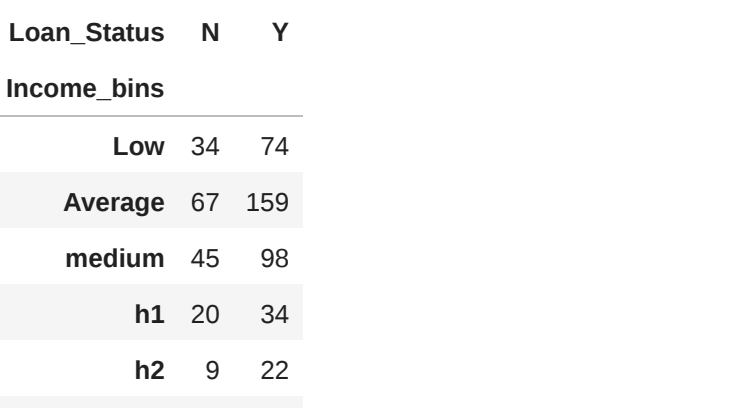
```
Out[53]: Test_IndResult(statistic=-0.11650844828724542, pvalue=0.453643960605259)
```

```
In [55]: ttest_ind(df_acc,df_rej,alternative="two.sided")
```

```
Out[55]: Test_IndResult(statistic=-0.11650844828724542, pvalue=0.99787812139518)
```

```
In [54]: sns.kdeplot(df_acc)
sns.kdeplot(df_rej)
```

```
Out[56]: <AxesSubplot: xlabel='ApplicantIncome', ylabel='Density'>
```



```
In [57]: kstest(df_acc,df_rej)

Out[57]: KstestResult(statistic=0.04393759873617693, pvalue=0.948028334325084)
```

```
In [59]: bins=[0, 2500, 4888, 6000, 8880, 10000, 81000]
labels=["Low", "Average", "Medium", "High", "Very high"]
df["Income_Bins"]=pd.cut(df["ApplicantIncome"],bins=bins,labels=labels)
df
```

```
Out[59]:
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area	Loan_Status	Income_Bins
0	LP001002	Male	No	0	Graduate	No	5849	0.0	NaN	360.0	1.0	Urban	Y	medium
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	1.0	Rural	N	medium
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0	1.0	Urban	Y	Average
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0	1.0	Urban	Y	Average
4	LP001008	Male	No	0	Graduate	No	6000	0.0	141.0	360.0	1.0	Urban	Y	medium
...
609	LP002978	Female	No	0	Graduate	No	2900	0.0	71.0	360.0	1.0	Rural	Y	Average
610	LP002979	Male	Yes	3+	Graduate	No	4106	0.0	40.0	180.0	1.0	Rural	Y	medium
611	LP002983	Male	Yes	1	Graduate	No	8072	240.0	253.0	360.0	1.0	Urban	Y	High
612	LP002984	Male	Yes	2	Graduate	No	7583	0.0	187.0	360.0	1.0	Urban	Y	High
613	LP002990	Female	No	0	Graduate	Yes	4583	0.0	133.0	360.0	0.0	Semurban	N	medium

614 rows x 14 columns

```
In [61]: va=pd.crosstab(df["Income_Bins"],df["Loan_Status"])
va

Out[61]:
```

	Loan_Status	N	Y
Income_Bins			
Low	34	74	
Average	67	159	
medium	45	98	
High	20	34	
High	9	22	
Very high	17	25	

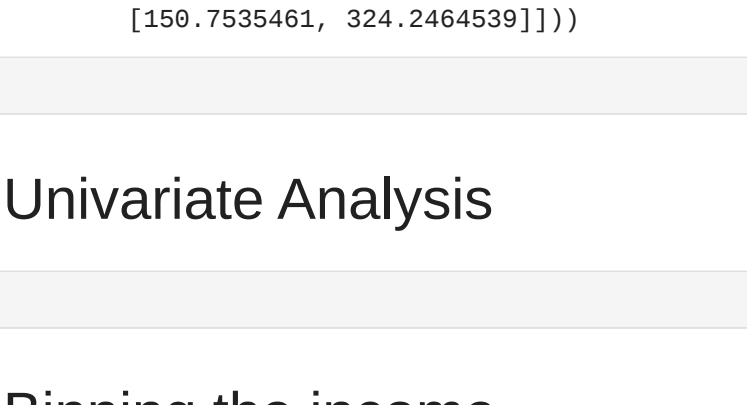
```
In [62]: chi2_contingency(va)

Out[62]: (1.2398175474316956,
0.941079844731327,
5,
```

```
In [64]: Income_bin = pd.crosstab(df["Income_Bins"],df["Loan_Status"])
Income_bin.div(Income_bin.sum(axis=1),axis=0).plot(kind='bar',figsize=(4,4))
labels=["Loan", "Average", "Medium", "High", "Very high"]
plt.xlabel("Percentage")
plt.show()
```

```
In [65]: sns.countplot(x=df["Credit_History"],hue=df["Loan_Status"])

Out[65]: <AxesSubplot: xlabel='Credit_History', ylabel='count'>
```



```
In [66]: va2=pd.crosstab(df["Credit_History"],df["Loan_Status"])
va2

Out[66]:
```

	Loan_Status	N	Y
Credit_History			
0.0	82	7	
1.0	97	378	

```
In [67]: chi2_contingency(va2)

Out[67]: (174.637296842529,
7.18479545791e-09,
4,
```

Univariate Analysis

```
In [ ]:
```

```
In [ ]:
```

```
In [71]: def check_ttest(x,y,alpha):
    _p_value=ttest_ind(x,y)
    if _p_value <= alpha:
        print("Feature x is a good Predictor")
    else:
        print("Feature x is a bad Predictor")

In [72]: check_ttest(df_acc,df_rej,0.05)
```

```
Feature x is a bad Predictor
```

```
In [ ]:
```

```
In [70]: cat_cols=seriesndf dtypes="object"
cat_columns=list(cat_cols.series[cat_cols.series.index])
cat_columns

Out[70]:
```

"Loan_ID",
"Gender",
"Married",
"Dependents",
"Education",
"Self_Employed",
"Property_Area",
"Loan_Status"

```
In [80]: cat_columns.remove("Loan_ID")

In [86]: # cat_columns.remove("Dependents")
# cat_columns.remove("Loan_Status")
```

```
In [106]: def check_chi2contingency(x,y,alpha):
    va=pd.crosstab(x,y)
    _p_value=chi2_contingency(va)
    if _p_value <= alpha:
        print("Feature x is a good Predictor for target variable Loan_Status")
    else:
        print("Feature x is a bad Predictor for target variable Loan_Status")

In [107]: cat_columns
```

```
In [107]: ["Gender", "Married", "Education", "Self_Employed", "Property_Area"]
```

```
In [108]: for col in cat_columns:
    check_chi2contingency(df[col],df["Loan_Status"],0.05)

Feature Married is a good Predictor for target variable Loan_Status
Feature Education is a good Predictor for target variable Loan_Status
Feature Self_Employed is a bad Predictor for target variable Loan_Status
Feature Property_Area is a good Predictor for target variable Loan_Status
```

