

EDA Software Analytics Coursework

Candidate Number: 418207

In [118]:

```
# Importing the required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix
```

In [119]:

```
# Loading the CSV file into the dataframe and displaying the first 10 entries
df_1 = pd.read_csv('loanapp.csv')
df_1.head(10)
```

	married	race	loan_decision	occupancy	loan_amount	applicant_income	num_units	num_dependants	self_employed	monthly_income	purchase_price	liquid_assets	mortgage_payment_history
0	True	white	reject	1	128	74	1.0	1.0	False	4583			
1	False	white	approve	1	128	84	1.0	0.0	False	2666			
2	True	white	approve	1	66	36	1.0	0.0	True	3000			
3	True	white	approve	1	120	59	1.0	0.0	False	2583			
4	False	white	approve	1	111	63	1.0	0.0	False	2008			
5	False	white	approve	1	141	72	1.0	0.0	False	6200			
6	False	white	approve	1	276	90	1.0	0.0	False	7500			
7	True	white	approve	1	100	72	1.0	1.0	False	3234			
8	True	white	approve	1	267	144	1.0	2.0	True	5417			
9	True	white	approve	1	175	97	1.0	0.0	False	6292			

In [120]:

```
# Displaying the column names and the datatypes in the dataframe
df_1.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1988 entries, 0 to 1987
Data columns (total 17 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   married               1988 non-null  object  
 1   race                  1988 non-null  object  
 2   loan_decision         1988 non-null  object  
 3   occupancy             1988 non-null  int64   
 4   loan_amount           1988 non-null  int64   
 5   applicant_income      1988 non-null  int64   
 6   num_units             1984 non-null  float64  
 7   num_dependants        1985 non-null  float64  
 8   self_employed         1988 non-null  bool     
 9   monthly_income        1988 non-null  int64   
10  purchase_price        1988 non-null  float64  
11  liquid_assets          1988 non-null  float64  
12  mortgage_payment_history 1988 non-null  int64   
13  consumer_credit_history 1988 non-null  int64   
14  filed_bankruptcy       1988 non-null  bool     
15  property_type          1988 non-null  int64   
16  gender                1974 non-null  object  
dtypes: bool(2), float64(4), int64(7), object(4)
memory usage: 237.0+ KB
```

In [121]:

```
# Displaying the number of rows and columns in the dataframe
df_1.shape
```

Out[121]:

```
(1988, 17)
```

1. Display descriptive statistics on the dataset.

For displaying the descriptive statistics of the dataset, we can call the describe function in pandas. It will display the mean, min/max values, standard deviation etc, for each column.

In [122]:

```
# Calling the describe function on the dataframe
df_1.describe()
```

	occupancy	loan_amount	applicant_income	num_units	num_dependants	monthly_income	purchase_price	liquid_assets	mortgage_payment_history
count	1988.000000	1988.000000	1988.000000	1984.000000	1985.000000	1988.000000	1988.000000	1988.000000	
mean	1.031698	143.272636	84.684105	1.122480	0.771285	5195.220825	196.304088	4620.333873	
std	0.191678	80.531477	87.079777	0.437315	1.104464	5270.360946	128.136030	67142.936043	
min	1.000000	2.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	100.000000	46.000000	1.000000	0.000000	2875.750000	129.000000	20.000000	
50%	1.000000	126.000000	64.000000	1.000000	0.000000	3812.500000	163.000000	38.000000	
75%	1.000000	165.000000	88.000000	1.000000	1.000000	5594.500000	225.000000	83.000000	
max	3.000000	980.000000	972.000000	4.000000	8.000000	81000.000000	1535.000000	1000000.000000	

1. Check if any records in the data have any missing values; handle the missing data as appropriate.

For checking if a column is null or not we can use isnull() which will return either True (for Null values) or False. To find the sum of all null values we can use the sum() over the isnull(). After finding the number of columns with null values, we can either drop those rows or interpolate values. Using dropna() and giving the columns with Null entries as the parameter we can delete rows with Null values. After dropping those rows, using info() we can check the new shape of the dataframe. The number of rows has decreased from 1988 to 1969 after cleaning the data.

In [123]:

```
# Finding the count of null values using isnull()
df_1.isnull().sum()
```

married	0
race	0
loan_decision	0
occupancy	0
loan_amount	0
applicant_income	0
num_units	4
num_dependants	3
self_employed	0
monthly_income	0
purchase_price	0
liquid_assets	0
mortgage_payment_history	0
consumer_credit_history	0
filed_bankruptcy	0
property_type	0
gender	14
dtype:	int64

In [124]:

```
# dropping the rows with Null values in the below columns
df_1 = df_1.dropna(subset=['num_units','num_dependants','married','gender'])
```

In [125]:

```
# Finding the count of null values using isnull()
df_1.isnull().sum()
```

married	0
race	0
loan_decision	0
occupancy	0
loan_amount	0
applicant_income	0
num_units	0
num_dependants	0
self_employed	0
monthly_income	0
purchase_price	0
liquid_assets	0
mortgage_payment_history	0
consumer_credit_history	0
filed_bankruptcy	0
property_type	0
gender	0
dtype:	int64

In [126]:

```
df_1.shape
```

Out[126]:

```
(1969, 17)
```

In [127]:

```
df_1.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 1969 entries, 0 to 1987
Data columns (total 17 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   married               1969 non-null  object  
 1   race                  1969 non-null  object  
 2   loan_decision         1969 non-null  object  
 3   occupancy             1969 non-null  int64   
 4   loan_amount           1969 non-null  int64   
 5   applicant_income      1969 non-null  int64   
 6   num_units             1969 non-null  float64  
 7   num_dependants        1969 non-null  float64  
 8   self_employed         1969 non-null  bool     
 9   monthly_income        1969 non-null  int64   
10  purchase_price        1969 non-null  float64  
11  liquid_assets          1969 non-null  float64  
12  mortgage_payment_history 1969 non-null  int64   
13  consumer_credit_history 1969 non-null  int64   
14  filed_bankruptcy       1969 non-null  bool     
15  property_type          1969 non-null  int64   
16  gender                1969 non-null  object  
dtypes: bool(2), float64(4), int64(7), object(4)
memory usage: 250.0+ KB
```

In [128]:

```
df_1.describe()
```

	occupancy	loan_amount	applicant_income	num_units	num_dependants	monthly_income	purchase_price	liquid_assets	mortgage_payment_history
count	1969.000000	1969.000000	1969.000000	1969.000000	1969.000000	1969.000000	1969.000000	1969.000000	
mean	1.031996	143.505333	84.908075	1.122905	0.770442	5204.838497	196.748556	4664.390041	
std	0.192576	80.802291	87.439115	0.438410	1.103450	5290.834720	126.785218	67464.766307	
min	1.000000	2.000000	0.000000	1.000000	0.000000	0.000000	25.000000	0.000000	
25%	1.000000	100.000000	48.000000	1.000000	0.000000	2876.000000	130.000000	20.000000	
50%	1.000000	126.000000	64.000000	1.000000	0.000000	3812.000000	163.000000	38.000000	
75%	1.000000	165.000000	88.000000	1.000000	1.000000	5600.000000	225.000000	83.000000	
max	3.000000	980.000000	972.000000	4.000000	8.000000	81000.000000	1535.000000	1000000.000000	

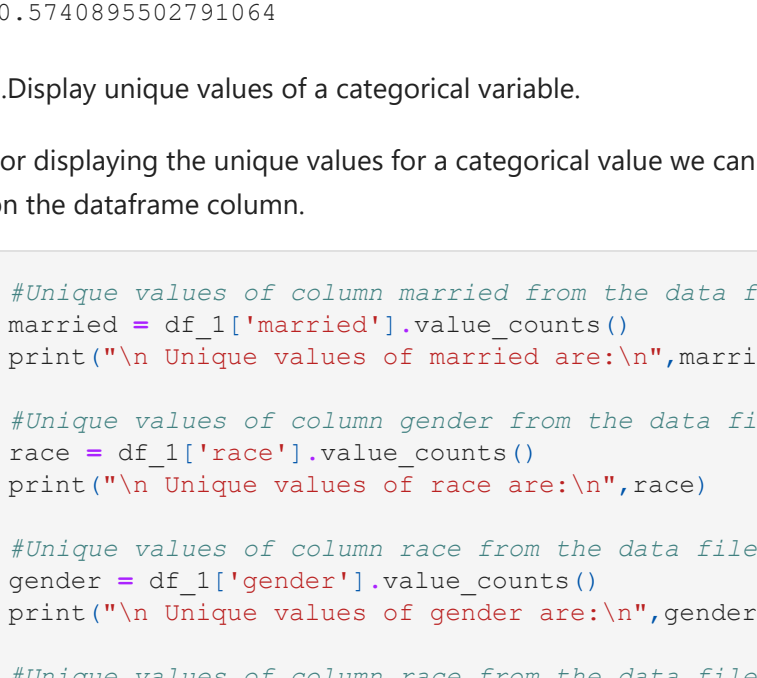
1. Build a graph visualizing the distribution of one or more individual continuous variables of the dataset.

To plot graphs in Python, we can use plot over the dataframe variable. To plot a histogram we will call plot over the dataframe df_1 variable 'loan_amount' with hist() for histogram). The bin size can be mentioned as a parameter inside the hist function. From the graph we can see that the maximum loan amount that has been given out is in the 0-200 range. Very less number of people have been given loans greater than 200.

In [129]:

```
# Plotting a histogram based on loan amount
df_1[['loan_amount']].plot.hist(alpha=0.7, bins=25)
```

Out[129]:



4. Build a graph visualizing the relationship in a pair of continuous variables. Determine the correlation between them.

For plotting a scatter plot, we can use the plot function again and then mention the plot type as scatter in the parameter kind. We also need to mention the variables to be used for the x and y axis. We can also specify the color for the scatter plot. From the graph we can interpret that there is no proper relationship between loan amount and applicant income. For calculating the correlation between a pair of continuous variables, we are using Pearson's correlation. We will call the corr() and mention the columns with which we want to check the correlation and then specify the method as 'pearson' as the parameter for corr().

In [130]:

```
# Plotting a scatter plot based on loan amount and applicant income
df_1.plot(kind='scatter', x='loan_amount', y='applicant_income', color='blue')
plt.show()
```



In [131]:

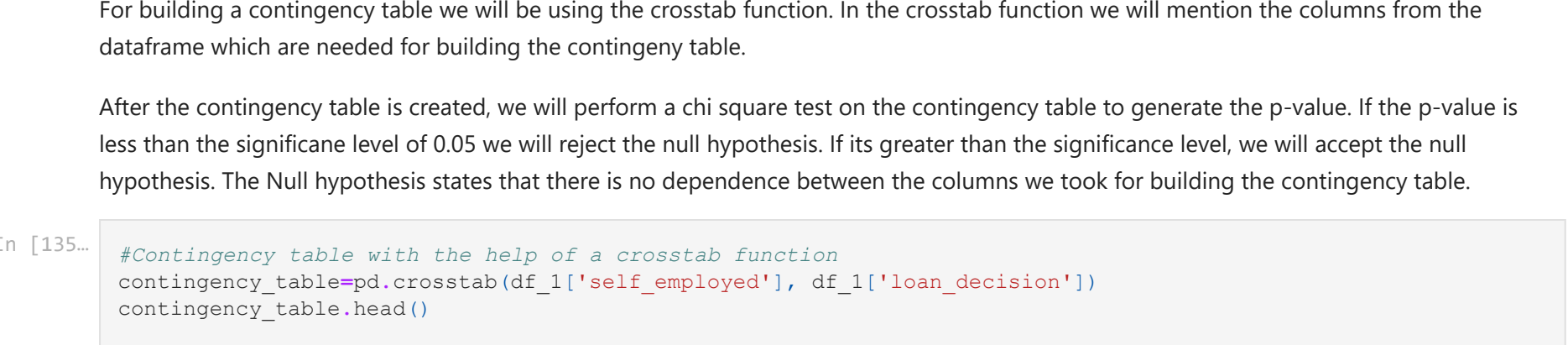
```
# Pearson's Correlation using correlation function
df_1[['loan_amount']].corr(df_1[['applicant_income']], method='pearson')
```

Out[131]:

```
0.43271137519859754
```

In [132]:

```
# Plotting a scatter plot based on loan amount and monthly income
df_1.plot(kind='scatter', x='loan_amount', y='monthly_income', color='blue')
plt.show()
```



In [133]:

```
# Pearson's Correlation using correlation function
df_1[['loan_amount']].corr(df_1[['monthly_income']], method='pearson')
```

Out[133]:

```
0.5740895502791064
```

5. Display unique values of a categorical variable.

For displaying the unique values of a categorical variable we can either use value_counts() on the columns of a dataframe or use the unique() on the dataframe column.

In [134]:

```
# Unique values of column married from the data file
married = df_1['married'].value_counts()
print("\n Unique values of married are:\n", married)

# Unique values of column gender from the data file
gender = df_1['gender'].value_counts()
print("\n Unique values of gender are:\n", gender)

# Unique values of column race from the data file
loan_decision = df_1['loan_decision'].value_counts()
print("\n Unique values of loan_decision are:\n", loan_decision)

# Unique values of column race from the data file
self_employed = df_1['self_employed'].value_counts()
print("\n Unique values of self_employed are:\n", self_employed)

# Unique values of column race from the data file
filed_bankruptcy = df_1['filed_bankruptcy'].value_counts()
print("\n Unique values of filed_bankruptcy are:\n", filed_bankruptcy)

Unique values of married are:
True    1298
False    671
Name: married, dtype: int64

Unique values of race are:
white    1666
black    195
hispan    108
Name: race, dtype: int64

Unique values of gender are:
male    1601
female   368
Name: gender, dtype: int64

Unique values of loan_decision are:
approve   1727
reject    242
Name: loan_decision, dtype: int64

Unique values of self_employed are:
False    1712
True      257
Name: self_employed, dtype: int64

Unique values of filed_bankruptcy are:
False    1834
True     135
Name: filed_bankruptcy, dtype: int64
```

6. Build a contingency table of two potentially related categorical variables. Conduct a statistical test of the independence between them.

For building a contingency table we will be using the crosstab function. In the crosstab function we will mention the columns from the dataframe which are needed for building the contingency table.

After the contingency table is created, we will perform a chi-square test on the contingency table to generate the p-value. If the p-value is less than the significance level of 0.05, we will reject the null hypothesis. If it is greater than the significance level, we will accept the null hypothesis. The Null hypothesis states that there is no dependence between the columns we took for building the contingency table.

In [135]:

```
# Contingency table with the help of a crosstab function
contingency_table = pd.crosstab(df_1['self_employed'], df_1['loan_decision'])
contingency_table.head()
```

Out[135]:

```
self_employed
```

```
loan_decision
```

```
approve
```

```
reject
```

```
False
```

```
1510
```

```
202
```

```
True
```

```
217
```

```
40
```

In [136]:

```
# Running a chi-square test and calculating the p-value
chi2, p_value, dof, expected = stats.chi2_contingency(contingency_table)
print(f"p-value: {p_value}")
```

Out[136]:

```
p-value: 0.1068853896762208
```

In [137]:

```
# Contingency table with the help of a crosstab function
contingency_table_2 = pd.crosstab(df_1['filed_bankruptcy'], df_1['loan_decision'])
contingency_table_2.head()
```

Out[137]:

```
loan_decision
```

```
approve
```

```
reject
```

```
filed_bankruptcy
```

```
False
```

```
1653
```

```
181
```

```
True
```

```
74
```

```
61
```

In [138]:

```
# Running a chi-square test and calculating the p-value
chi2, p_value, dof, expected = stats.chi2_contingency(contingency_table_2)
print(f"p-value: {p_value}")
```

Out[138]:

```
p-value: 8.677306949234421e-33
```

7. Retrieve one or more subset of data based on two or more criteria and present descriptive statistics on the subsets).

For retrieving a subset of the data, we will specify the conditions using logical operators or by specifying the condition with which we need the data. Here we have taken 2 subsets of the dataframe based on gender, approve as the loan decision and where the number of dependants are less than 2. For generating the descriptive statistical report of the subset, we can use the describe(). From the data we can see that the mean loan amount for male is greater than that of females.

In [139]:

```
# Creating a subset of the dataset based on gender, loan_decision and number of dependants.
df_1_subset_male = df_1[(df_1['gender'] == 'male') & (df_1['loan_decision'] == 'approve') & (df_1['num_dependants'] < 2)]
df_1_subset_female = df_1[(df_1['gender'] == 'female') & (df_1['loan_decision'] == 'approve') & (df_1['num_dependants'] < 2)]
```

Out[139]:

	occupancy	loan_amount	applicant_income	num_units	num_dependants	monthly_income	purchase_price	liquid_assets	mortgage_payment_history
count	1013.000000	144.150049	83.465943	1.107601	0.228036	5007.127345	196.387868	6018.981464	
mean	1.083752	80.732254	74.259001	0.416165	0.419773	4945.086190	127.778497	67613.678020	
min	1.000000	2.000000	5.000000	1.000000	0.000000	0.000000	25.000000	0.000000	
25%	1.000000	104.000000	51.000000	1.000000	0.000000	2773.000000	132.000000	21.000000	
50%	1.000000	129.000000	66.000000	1.000000	0.000000	3727.000000	165.000000	40.450000	
75%	1.000000	166.000000	89.000000	1.000000	0.000000	5500.000000	225.000000	87.100000	
max	3.000000	980.000000	796.000000	4.000000	1.000000	67000.000000	1535.000000	1000000.000000	

In [140]:

```
# Describing the subset based on the loan_amount
df_1_subset_male['loan_amount'].describe()
```

Out[140]:

```
count    1013.000000
mean     144.150049
std       80.732254
min        2.000000
25%      104.000000
50%      129.000000
75%      166.000000
max       980.000000
Name: loan_amount, dtype: float64
```

In [141]:

```
# Creating a subset of the dataset based on gender, loan_decision and number of dependants.
df_1_subset_female = df_1[(df_1['gender'] == 'female') & (df_1['loan_decision'] == 'approve') & (df_1['num_dependants'] < 2)]
df_1_subset_female.describe()
```

Out[141]:

OLS Regression Results				
Dep. Variable:	loan_amount	R-squared:	0.700	
Model:	OLS	Adj. R-squared:	0.700	
Method:	Least Squares	F-statistic:	2297.	
Date:	Thu, 24 Mar 2022	Prob (F-statistic):	0.00	
Time:	23:09:22	Log-Likelihood:	-10255.	
No. Observations:	1969	AIC:	2.052e+04	
Df Residuals:	1966	BIC:	2.053e+04	
Df Model:	2			
Covariance Type:	nonrobust			
	coef	std err	t	P> t [0.025 0.975]
Intercept	38.7057	1.851	20.916	0.000 35.077 42.336