

# Credit Card Clients in Taiwan in 2005-Box Plots, Count Plots & Histograms REPORT



**SAMUEL**

SKILL LYNC TRAINEE EMPLOYEE

06.05.202

## **Index:**

- 1. Introduction**
- 2. Problem statement**
- 3. Solution approach**
- 4. Data overview**
  - a) Data extraction from a dataset**
  - b) Plots of dataset**
  - c) Analysing the data and making the more understandable**
- 5. Advantages of Credit card**
- 6. Technologies used for data visualisation and analysis**
- 7. Conclusion**

**Introduction:** Based on the financial capability of a client, they get a credit limit, i.e., the maximum amount they can spend in a month through a credit card.

Credit card companies maintain comprehensive data about each of their clients. By analysing the data, they can know what would be the maximum amount they won't be able to recover from their clients yet able to make a significant profit in a financial year to run a sustainable business.

**Problem statement:** The credit card clients dataset is full of irregularities and incorrect values. You need to replace them with the right values. Additionally, you have to create box plots, count plots and histograms to find a specific trend (if there exists) in the dataset.

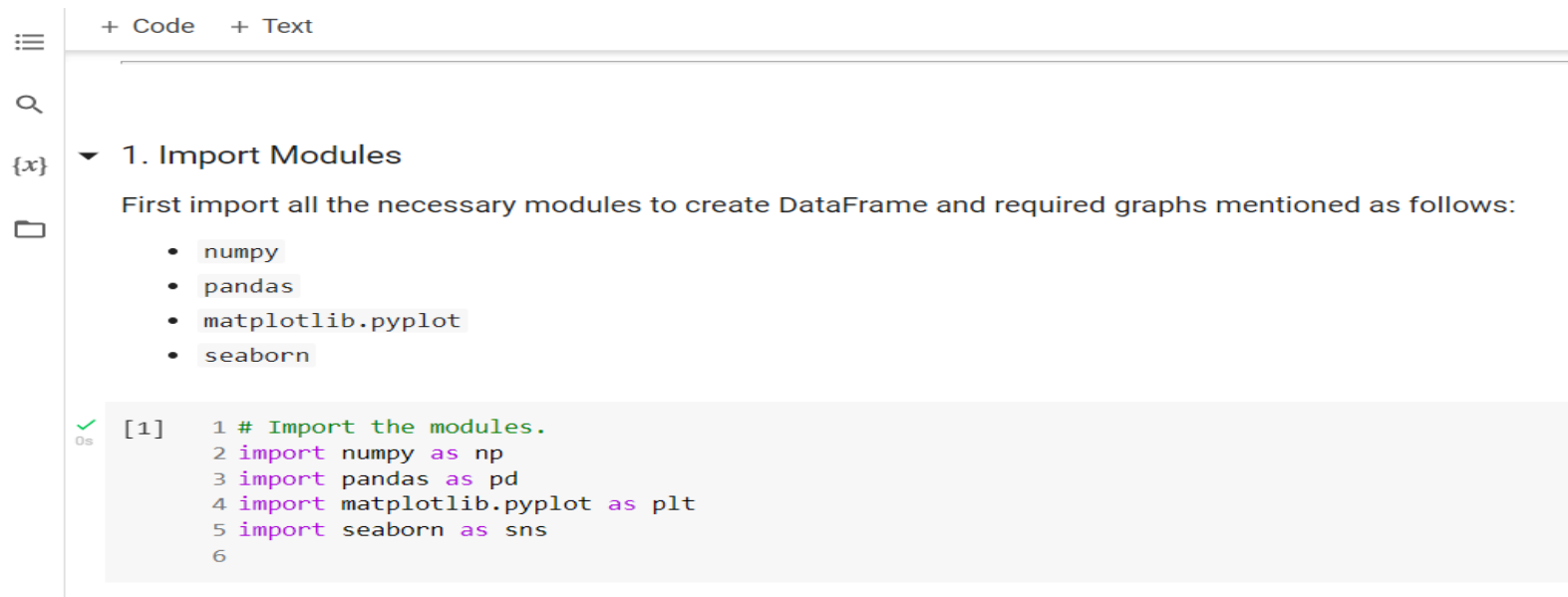
**Solution Approach:** Here we use python and machine learning tools to make the values right from given irregularities values in the given dataset of credit cards, Additionally we create box plots, count plots and histograms to find the specific trend in the dataset, from the plots we analyse the client's credit card data and make models for better credit card service

## Data overview

### Data:

Here we are abstracting the required data from credit card resources by using **Python Libraries**.

Here are the screenshots of data abstraction from the source :



The screenshot shows a Jupyter Notebook interface. On the left, there is a sidebar with icons for a menu, search, and file explorer. The main area has a tab labeled '+ Code' and '+ Text'. Below the tab, there is a section titled '1. Import Modules' with a dropdown arrow. Under this section, there is a text prompt: 'First import all the necessary modules to create DataFrame and required graphs mentioned as follows:'. Below this prompt, there is a bulleted list of libraries: numpy, pandas, matplotlib.pyplot, and seaborn. At the bottom, there is a code cell with a green checkmark and '0s' indicating successful execution. The code cell contains the following Python code:

```
[1] 1 # Import the modules.  
2 import numpy as np  
3 import pandas as pd  
4 import matplotlib.pyplot as plt  
5 import seaborn as sns  
6
```

Figure 1

Source link: "[https://raw.githubusercontent.com/m-narayanan22/datasets/main/UCI\\_Credit\\_Card.csv](https://raw.githubusercontent.com/m-narayanan22/datasets/main/UCI_Credit_Card.csv)"

From the source link, we got the dataset values of credit card client details like their Age, Education, Limit balance and etc.,

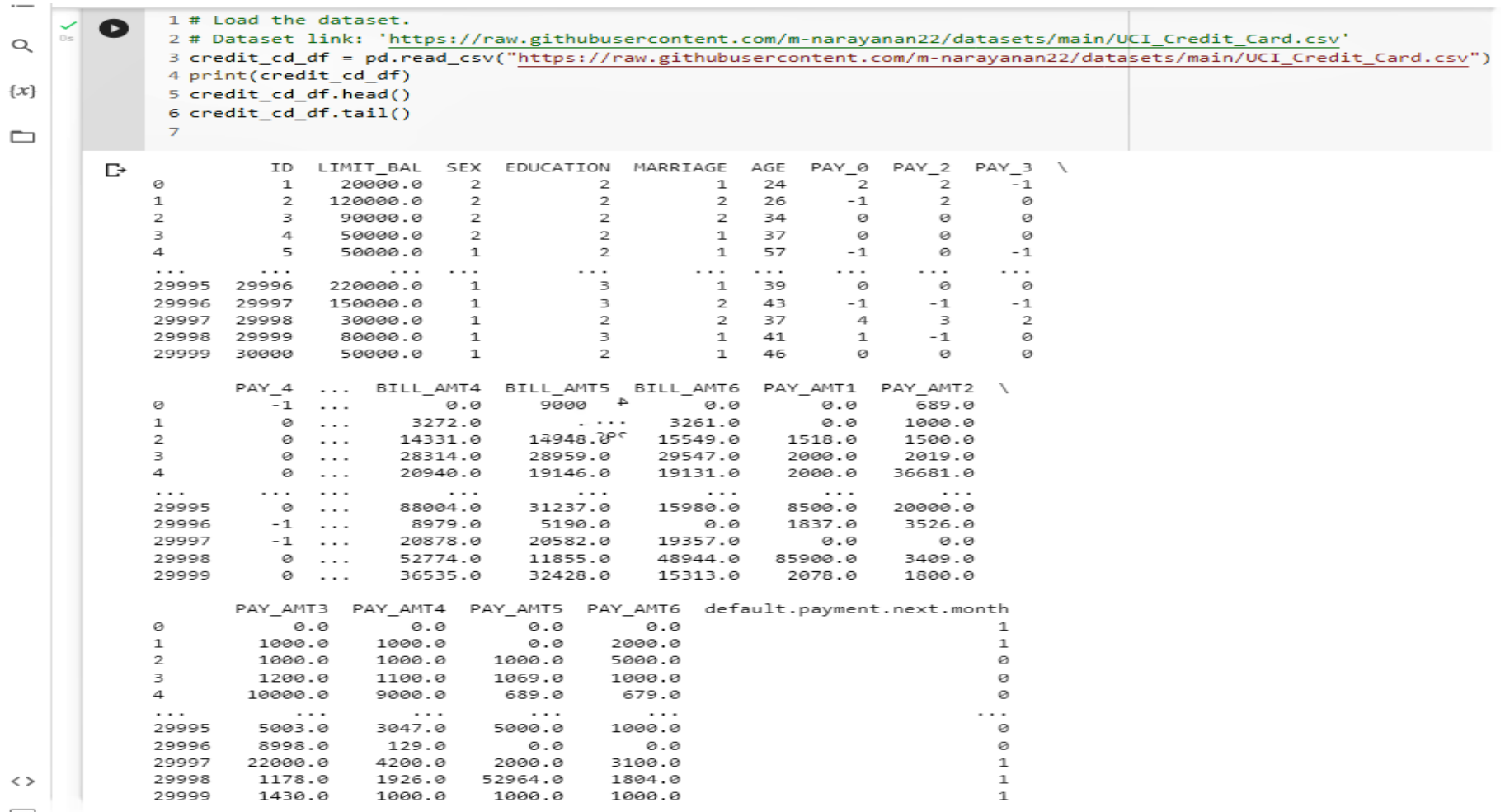
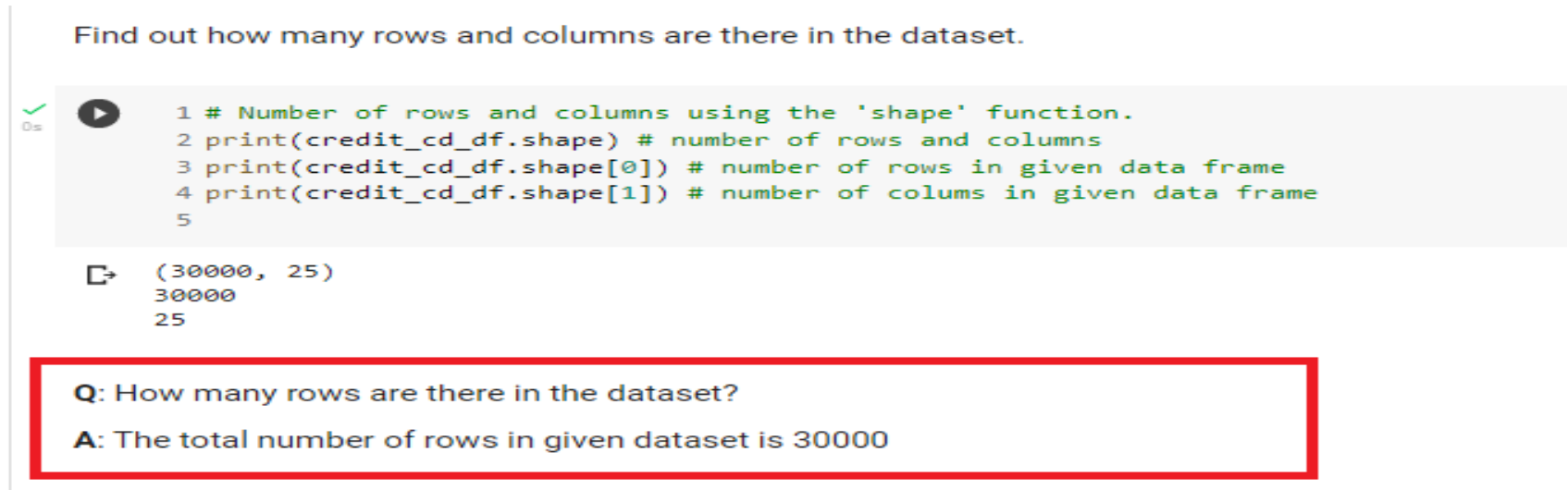


Figure 2

From above Figure 2, we got know that the clients' credit dataset is huge, so for to analyse the given data use

python libraries, and we will plot the data, and we will see the plot in the upcoming Figures

Here we are finding the rows and columns are there in the dataset.



Find out how many rows and columns are there in the dataset.

```
1 # Number of rows and columns using the 'shape' function.
2 print(credit_cd_df.shape) # number of rows and columns
3 print(credit_cd_df.shape[0]) # number of rows in given data frame
4 print(credit_cd_df.shape[1]) # number of columns in given data frame
5
```

```
(30000, 25)
30000
25
```

**Q:** How many rows are there in the dataset?

**A:** The total number of rows in given dataset is 30000

Figure 3

From the above Figure 3, the number of rows and columns is (3000, 25).

The red marked box represents the Question and answers for the specified task.

Checking For The Missing Values

Now, check whether the dataset contains any NaN or null or missing values

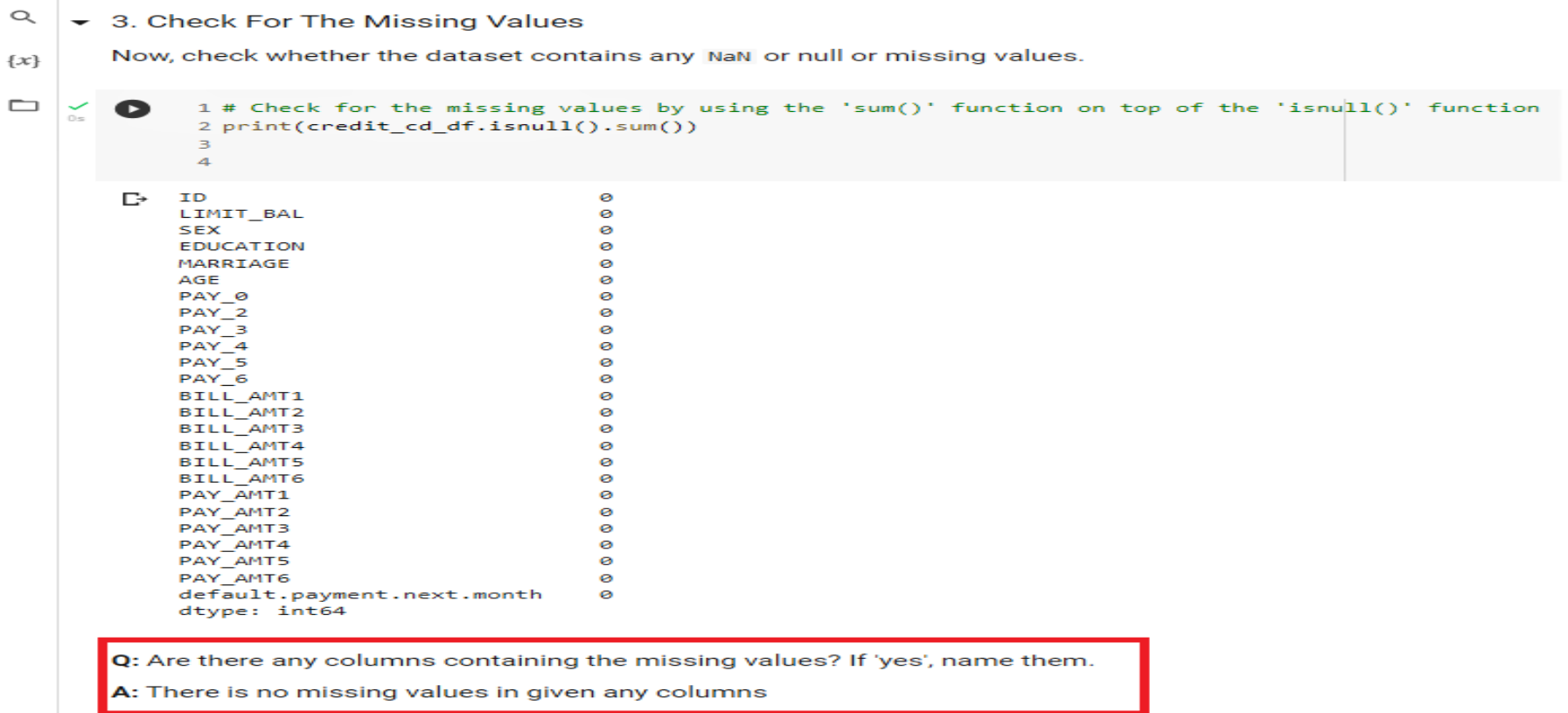


Figure 4

From the above Figure 4, we can observe that there are no missing values in the given dataset.

Now, We are going to find the education column's value, count values and plotting of the education columns.

**Q:** What value(s) is/are contained in the `EDUCATION` column apart from the values 1 to 5? And what will you do with them?

**A:** 0 and 6 is the values apart from the 1 to 5 and We have to replace 0 and 6 value with 5

```
[9] 1 # Replace the unwanted values in the 'EDUCATION' column. Ignore if there are none.
2 replace_indices1 = credit_cd_df[credit_cd_df['EDUCATION'] == 6].index
3 credit_cd_df.loc[replace_indices1, ['EDUCATION']] = 5
4
5 replace_indices2 = credit_cd_df[credit_cd_df['EDUCATION'] == 0].index
6 credit_cd_df.loc[replace_indices2, ['EDUCATION']] = 5
7 replace_indices2
8
9
Int64Index([ 3769,  5945,  6876, 14631, 15107, 16881, 16896, 17414, 19920,
            20030, 23234, 24137, 27155, 27270],
            dtype='int64')
```

**Hint:** You can replace the rows with 0 and 6 in the `Education` column by using the `loc[]` function.

**Syntax:** `df.loc[df['column_name'] == old_value, 'EDUCATION'] = new_value`

Calculate the percentage of each value in the `EDUCATION` column.

```
1 # Percentage of each value in the 'EDUCATION' column.
2 per_1 = credit_cd_df["EDUCATION"].value_counts() * 100 / credit_cd_df.shape[0]
3 per_1
4
5
2    46.766667
1    35.283333
3    16.390000
5     1.150000
4     0.410000
Name: EDUCATION, dtype: float64
```

**Hint:** You can get the total number of counts of each value in the column by using the `value_counts()` function. Then you can calculate the percentage of each value by multiplying the total number of counts of each value with 100 and dividing the resultant value by the total number of rows in the DataFrame (`df.shape[0]`).

**Q:** What percent of clients were university graduates?

**A:** From above data we got..46.7% percent of cilents were university graduate (round off for 46.7666 is 47)

**Figure 5**

From above Figure 5, we can understand that the percentage of clients who were university graduates is 47%

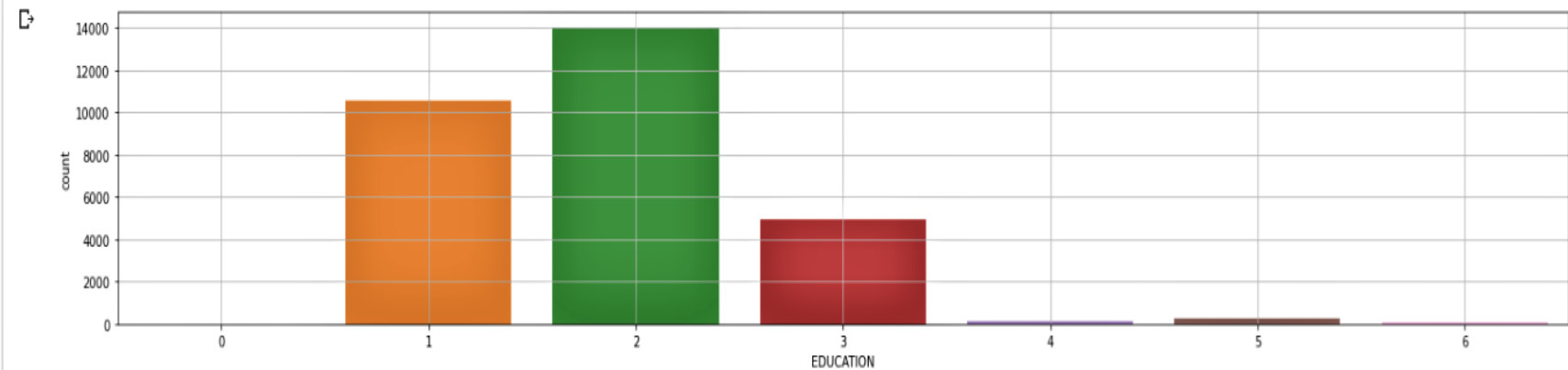


So we can observe that half of the clients are Educated.

### Creating a count plot for the 'EDUCATION' column:

Create a count plot for the 'EDUCATION' column.

```
1 # Count plot for the 'EDUCATION' column using the 'countplot()' function of 'seaborn' module.  
2 plt.figure(figsize=(24,4))  
3 sns.countplot(x="EDUCATION", data =credit_cd_df)  
4 plt.grid()  
5  
6
```



**Figure 6**

From above Figure 6, we can say that there are 6 bars with different count values, the green having the highest count of 14000.

## Checking the data types of all the columns using the 'info()' function.

You may require to check the data-type of every column. So, instead of applying the `dtype` keyword one-by-one for each column, you can use the `info()` function to check the data-types of all the columns at once. It also tells you the total number of rows and columns in a DataFrame. Here's the syntax:

**Syntax:** `data_frame.info()`

where `data_frame` is a variable storing some Pandas DataFrame.

**Note:** This function is applicable only to Pandas DataFrame.

```
1 # Check the data-types of all the columns using the 'info()' function.
2 credit_cd_df.info()
3
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 25 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0    ID                                         30000 non-null   int64
1    LIMIT_BAL                                30000 non-null   float64
2    SEX                                       30000 non-null   int64
3    EDUCATION                                30000 non-null   int64
4    MARRIAGE                                 30000 non-null   int64
5    AGE                                       30000 non-null   int64
6    PAY_0                                    30000 non-null   int64
7    PAY_2                                    30000 non-null   int64
8    PAY_3                                    30000 non-null   int64
9    PAY_4                                    30000 non-null   int64
10   PAY_5                                    30000 non-null   int64
11   PAY_6                                    30000 non-null   int64
12   BILL_AMT1                                30000 non-null   float64
13   BILL_AMT2                                30000 non-null   float64
14   BILL_AMT3                                30000 non-null   float64
15   BILL_AMT4                                30000 non-null   float64
16   BILL_AMT5                                30000 non-null   float64
17   BILL_AMT6                                30000 non-null   float64
18   PAY_AMT1                                30000 non-null   float64
19   PAY_AMT2                                30000 non-null   float64
20   PAY_AMT3                                30000 non-null   float64
21   PAY_AMT4                                30000 non-null   float64
22   PAY_AMT5                                30000 non-null   float64
23   PAY_AMT6                                30000 non-null   float64
24   default.payment.next.month              30000 non-null   int64
dtypes: float64(13), int64(12)
memory usage: 5.7 MB
```

Figure 7

## Marital Status of Clients:

The below figure represents the client's marital status which means who is married, single or divorced.

### 4.3 Marital Status of Clients

The `MARRIAGE` column contains the following three different types of values:

- 1 denotes that a client is married
- 2 denotes that a client is single
- 3 denotes all other possible marital statuses such as divorced, widowed etc.

If there are any other values, then they should be replaced with 3 because it covers all the other possible cases of marital status of a client.

Calculate the counts of each value in the `MARRIAGE` column.

```
1 # Counts of each value in the 'MARRIAGE' column.
2 credit_cd_df['MARRIAGE'].value_counts()
3
```

```
2    15964
1    13659
3     323
0        54
Name: MARRIAGE, dtype: int64
```

**Q:** What value(s) is/are contained in the `MARRIAGE` column apart from the values 1, 2 and 3? What are their counts?

**A:** 0 is value is apart from 1,2,3 and the count value is 54

```
[15] 1 # Replace the unwanted values ('0') in the 'MARRIAGE' column with '3'. Ignore if there are none.
2 credit_cd_df.loc[credit_cd_df['MARRIAGE'] == 0, 'MARRIAGE'] = 3
3 Count_meg = credit_cd_df['MARRIAGE'].value_counts()
4 Count_meg
5
6
```

```
2    15964
1    13659
3     377
Name: MARRIAGE, dtype: int64
```

Figure 8

From the above Figure 8, we can observe that the number of married, single or divorced and 1,2 and 3 represent the

marital status of clients.

Calculating the percentage of the values in the MARRIAGE column in Figure 9 below:

Calculate the percentage of the values in the MARRIAGE column.

```
✓ [16] 1 # Percentage of the values in the 'MARRIAGE' column.\
0s    2 Count_meg_per = Count_meg*100/credit_cd_df.shape[0]
      3 Count_meg_per
      4
      5
      6
      7
      8

      2    53.213333
      1    45.530000
      3     1.256667
      Name: MARRIAGE, dtype: float64
```

**Q:** What of clients were married?

**A:** The percentage of people who got married are 45.53

Figure 9

From above Figure 9, we can say that the percentage of married clients is 46 (Round off value).

Creating a count plot for the MARRIAGE column so that we can analyse the data:

```
1 # Count plot for the 'MARRIAGE' column.
2 plt.figure(figsize =(20,4))
3 sns.countplot(x='MARRIAGE', data =credit_cd_df, edgecolor ='Black')
4 plt.grid()
5
6
7
```

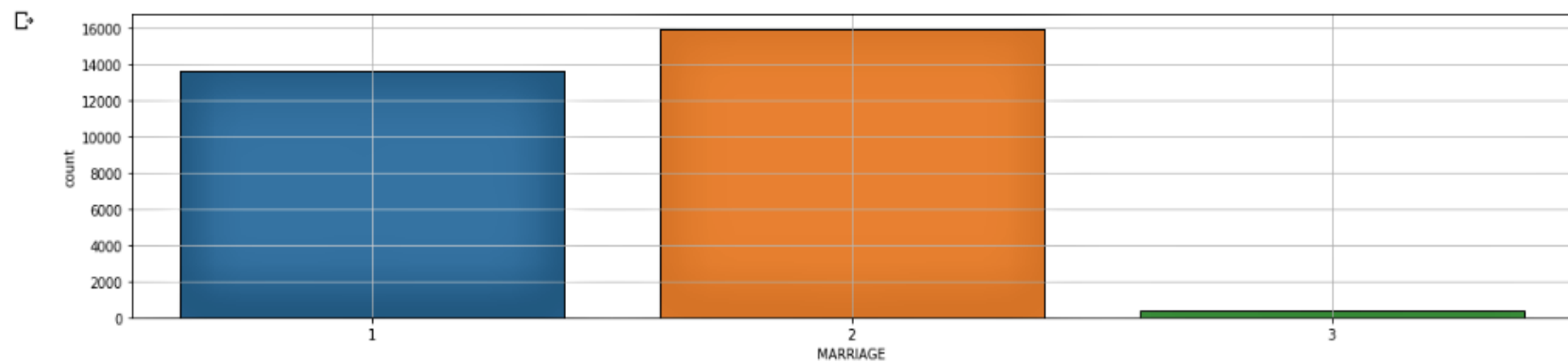


Figure 10

**Creating box plots and histograms for the columns containing continuous numeric values.**

Box Plot & Histogram For The AGE Column

All the histograms are in grids.

**Creating a box plot for the AGE column by using the given dataset:**

#### 5. Box Plots & Histograms

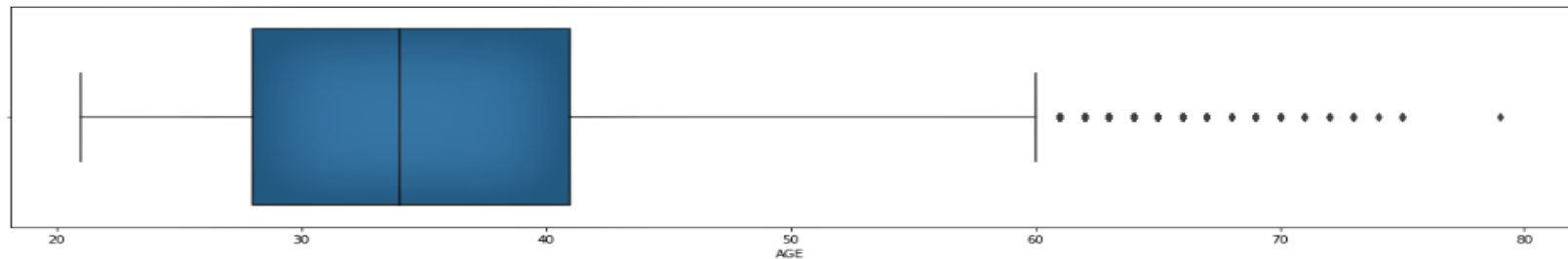
The final task is to create box plots and histograms for the columns containing continuous numeric values.

**Note:** All the histograms must have grids.

##### 5.1 Box Plot & Histogram For The AGE Column

Create a box plot for the AGE column.

```
1 # Box plot for the 'AGE' column using the 'boxplot()' function.  
2 plt.figure(figsize=(20,4))  
3 sns.boxplot(x=credit_cd_df['AGE'])  
4 plt.show()  
5  
6
```



**Q:** From the box plot for the AGE column, what is the approx median age of a credit card holder?

**A:** The approx. median age of a credit card holder is 35

**Figure 11**

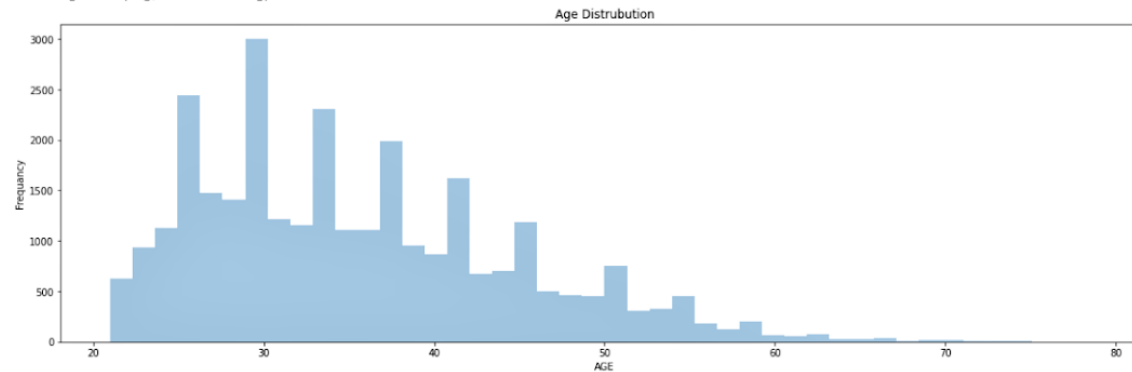
**Here,** From the above Figure,11 10we can observe the boxplot of age in different distributions at different intervals, also we can say that the Median age of the client is 35, and the points dots are very negligible values in the whole distribution in a boxplot.

### Creating Histogram distplot for the column of Age:

Create a histogram for the AGE column.

```
1 # Histogram for the 'AGE' column using 'distplot()' function from the 'seaborn' module.
2 plt.figure(figsize=(20,6))
3 sns.distplot(credit_cd_df['AGE'], bins = 44, kde = False)
4 plt.ticklabel_format(style = 'plain') #The values of x coordinates are high it will represents in 'e' notations by default, so we use ticklabel_format(style = 'plain') to remove 'e'.
5 plt.xlabel('AGE')
6 plt.ylabel('Frequency')
7 plt.title('Age Distrubution')
8 plt.show()
9
10
```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level



Q: Is there some peculiar pattern in the AGE histogram?

A: yes, (peculiar pattern is defined as when the pattern are unusual in shape)

**Figure 12**

**Here,** From above 12 we can say that the histogram of the Age column is in a peculiar pattern, which means the Age distribution is not uniform.

## Creating a box plot & histogram For The LIMIT\_BAL column:

### ▼ 5.2 Box Plot & Histogram For The LIMIT\_BAL Column

Create a box plot for the LIMIT\_BAL column.

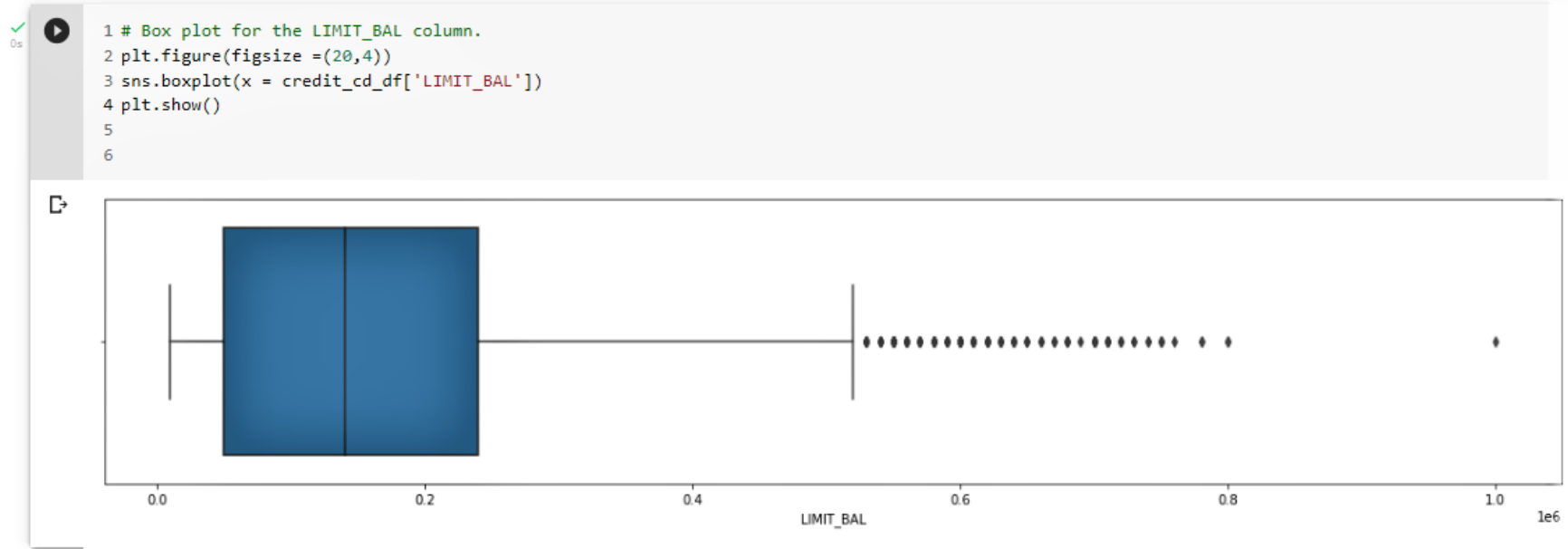


Figure 13

**Here**, in Figure 13 we can say that the Median of the LiMIT\_BAL column is Median is 140000 and max value is 1000000, and the points which are dotted in shape are negligible values

## Creating a dist plot & histogram For The LIMIT\_BAL column:

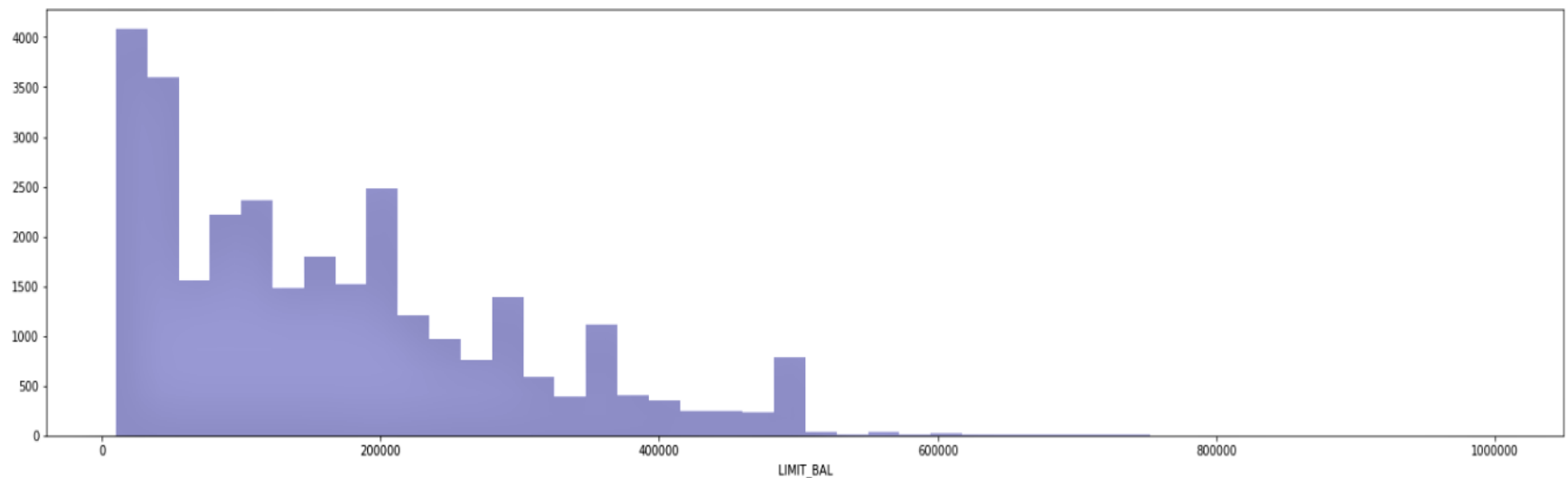


```

1 # Histogram for the 'LIMIT_BAL' column using 'distplot()' function from the 'seaborn' module.
2 plt.figure(figsize =(25,6))
3 sns.distplot(credit_cd_df['LIMIT_BAL'], bins = 44, kde = False, color ='darkblue')
4 plt.ticklabel_format(style = 'plain')
5 plt.show()
6
7
8

```

⚠ /usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to us  
warnings.warn(msg, FutureWarning)



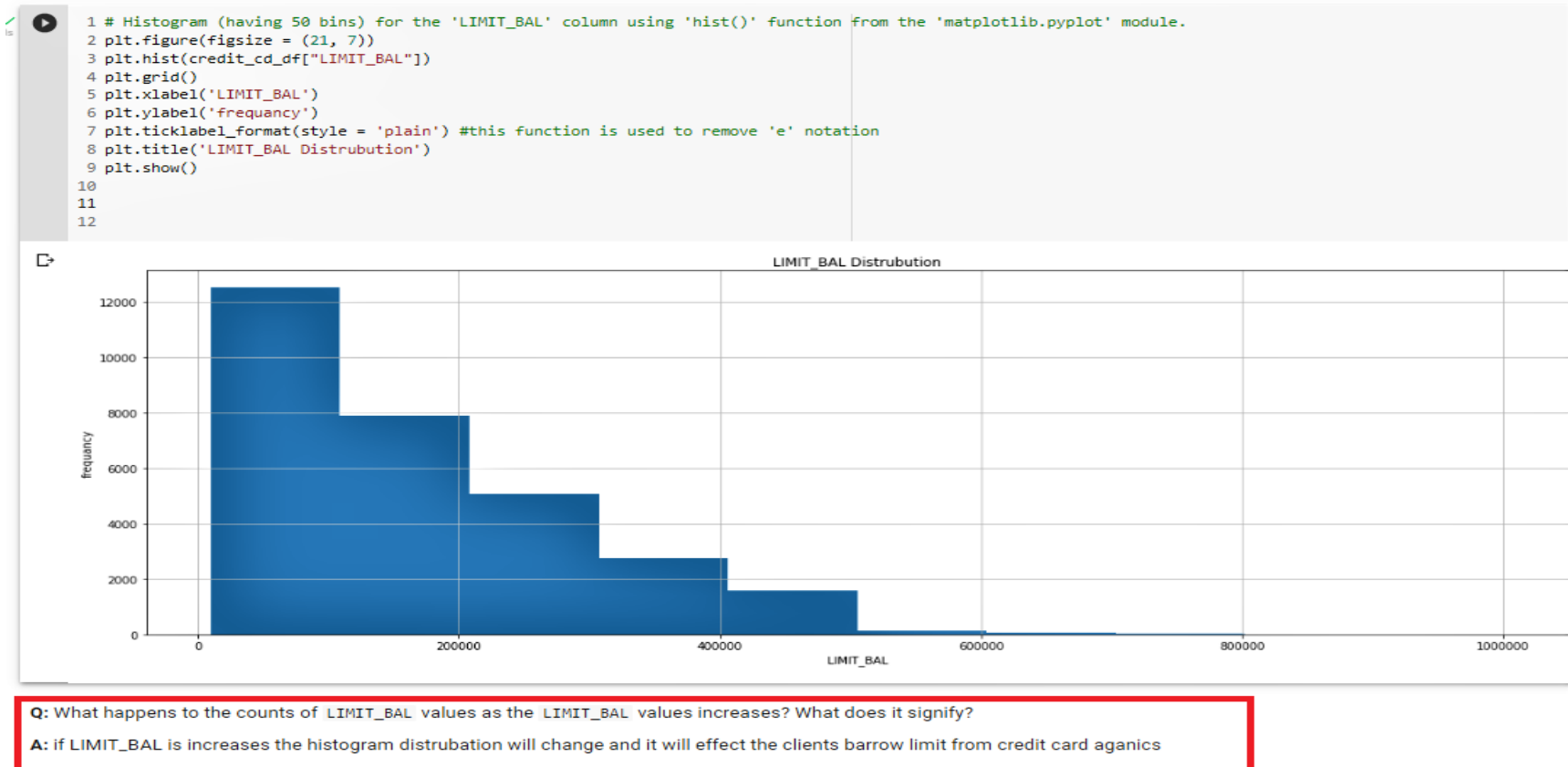
**Figure 14**

Here, in Figure 14 we observe that the LIMIT\_BAL Column count values are in the high range, and the distribution is not in the distribution that's why the pattern of the histogram is a peculiar pattern.

**Creating Histogram (having 50 bins) for the 'LIMIT\_BAL' column :**

Here is down below the histogram is

created with the required conditions :



**Figure 15**

Here, in Figure 15 we can say that the highest frequency is above 12000 and the distribution is decreasing with the value of LIMIT\_BAL and the distribution is in decreasing pattern and if every client LIMIT\_BAL is increasing the x coordinate only will affect. And frequency will not be

affected but the client will spend on things it will affect his lifestyle but the client spends more he/she will be debit in future.

## Advantages of a Credit Card

### 1. One-Time Bonuses

There's nothing like an initial bonus opportunity when [getting a new credit card](#). Oftentimes, applicants with [good credit](#) or excellent credit can get approved for credit cards that offer bonuses worth \$150 or more (sometimes much more) in exchange for spending a certain amount (anywhere from \$500 to several thousands of dollars) in the first several months the account is open.

### 2. Cash Back

The cash-back credit card was first popularized in the United States by Discover, and the idea was simple: Use the card and get 1% of your purchases rebated in the form of cashback. Today, the concept has grown and matured. Now, some cards now offer 2%, 3% or even as much as 6% cashback on selected purchases, though such lucrative offers involve quarterly or annual spending caps. The [best cash-back cards](#) are those that charge minimal fees and interest while offering a high rewards rate.

### 3. Rewards Points

Credit cards are set up to allow cardholders to earn one or more points per dollar in spending. Many [reward credit cards](#) provide bonus points for certain categories of spending like restaurants, groceries or [gasoline](#). When certain earnings thresholds are reached, points can be redeemed for travel, [gift cards](#) from retailers

and restaurants, or merchandise items through the credit card company's online rewards portal.

## 5. Safety

Paying with a credit card makes it easier to avoid losses from fraud. When your debit card is used by a thief, the money is missing from your account instantly. Legitimate expenses for which you've scheduled online payments or mailed checks may bounce, triggering insufficient funds fees and affecting your credit. Even if not your fault, these late or missed payments can lower your [credit score](#). It can take time for fraudulent transactions to be reversed and money restored to your account while the bank investigates.

## Technologies used for the project:

**We used Python tools to analyse the given data and make a plot of the different columns like Age, Education, Marriage etc., with different conditions for a better understanding of the dataset, and we find every possible condition in the dataset like Max, Median, Min etc., this makes to analyse given data and makes easy to credit card companies for clients requirements.**

## **Conclusion:**

**From the given dataset of Credit Card Clients in Taiwan in 2005, we concluded that the age and sex is not a barrier to the user not having a credit card in the modern era, the credit card is becoming a part of human life and having a credit card makes users lead a better life and the user can make any transaction by simply using a credit card like online travel booking, making online purchase etc,**

**.....END.....**