

CHANGING WORKING DIRECTORY

In [1]:

```
1 import os  
2 os.chdir(r'D:\MCT\Python Project\Project_Final')  
3 os.getcwd()
```

Out[1]: 'D:\\MCT\\Python Project\\Project_Final'

In [2]:

```
1 import warnings  
2 warnings.filterwarnings('ignore')
```

IMPORT LIBRARY

In [3]:

```
1 import pandas as pd  
2 import numpy as np  
3 import pandas_profiling as pp  
4 import seaborn as sns  
5 from datetime import datetime, timedelta, date  
6 import matplotlib.pyplot as plt  
7 %matplotlib inline  
8 import sklearn  
9 from sklearn.preprocessing import LabelEncoder, OneHotEncoder  
10 from sklearn.preprocessing import StandardScaler  
11 from sklearn.model_selection import train_test_split
```

In [13]:

```
1 app_df=pd.read_csv('D:\\MCT\\Python Project\\Project Credit Card_Pavithra'
2 credit_df=pd.read_csv('D:\\MCT\\Python Project\\Project Credit Card_Pavi
3 app_df
```

Out[13]:

	ID	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AI
0	5008804	M	Y	Y	0	
1	5008805	M	Y	Y	0	
2	5008806	M	Y	Y	0	
3	5008808	F	N	Y	0	
4	5008809	F	N	Y	0	
...
438552	6840104	M	N	Y	0	
438553	6840222	F	N	N	0	
438554	6841878	F	N	N	0	
438555	6842765	F	N	Y	0	
438556	6842885	F	N	Y	0	

438557 rows × 18 columns

**Days_Birth and Days Employed is negative -- need to check**

In [14]: 1 credit_df

Out[14]:

	ID	MONTHS_BALANCE	STATUS
0	5001711	0	X
1	5001711	-1	0
2	5001711	-2	0
3	5001711	-3	0
4	5001712	0	C
...
1048570	5150487	-25	C
1048571	5150487	-26	C
1048572	5150487	-27	C
1048573	5150487	-28	C
1048574	5150487	-29	C

1048575 rows × 3 columns

Month_balance is in negative need to do further analysis

MERGING DATASET

we merge both dataset using inner join, which means use intersection of keys from both frames,preserve the order of the left keys.

In [15]: 1 df = app_df.merge(credit_df, how='inner', on=['ID'])

SHAPE OF DATAFRAME

In [16]: 1 print("Shape of Application Dataframe:",app_df.shape)
2 print("Shape of Credit Dataframe:",credit_df.shape)
3 print("Shape of Merged Dataframe:",df.shape)

Shape of Application Dataframe: (438557, 18)
Shape of Credit Dataframe: (1048575, 3)
Shape of Merged Dataframe: (777715, 20)

Application Dataframe has 438557 rows and 18 columns and Credit dataframe has 1048575 rows and 3 columns. After merging using inner joining, dataframe has 777715 rows and 20 columns.

DATA PROFILING

```
In [ ]: 1 Credit_Report = pp.ProfileReport(df,title="Credit_Report")
2 Credit_Report.to_file("Credit_Report.html")
```

UNIQUE VALUE

```
In [17]: 1 app_df.apply(lambda x: len(x.unique()))
```

```
Out[17]: ID                438510
CODE_GENDER                 2
FLAG_OWN_CAR                 2
FLAG_OWN_REALTY                2
CNT_CHILDREN                  12
AMT_INCOME_TOTAL                866
NAME_INCOME_TYPE                  5
NAME_EDUCATION_TYPE                5
NAME_FAMILY_STATUS                  5
NAME_HOUSING_TYPE                  6
DAYS_BIRTH                     16379
DAYS_EMPLOYED                   9406
FLAG_MOBIL                      1
FLAG_WORK_PHONE                   2
FLAG_PHONE                      2
FLAG_EMAIL                      2
OCCUPATION_TYPE                  19
CNT_FAM_MEMBERS                  13
dtype: int64
```

```
In [18]: 1 credit_df.apply(lambda x: len(x.unique()))
```

```
Out[18]: ID                45985
MONTHS_BALANCE                  61
STATUS                          8
dtype: int64
```

In [19]: 1 df.nunique()

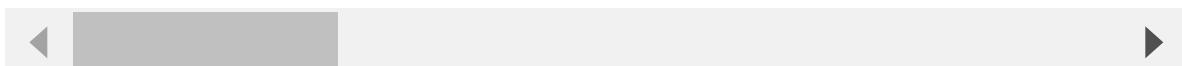
```
Out[19]: ID           36457
CODE_GENDER      2
FLAG_OWN_CAR     2
FLAG_OWN_REALTY  2
CNT_CHILDREN     9
AMT_INCOME_TOTAL 265
NAME_INCOME_TYPE 5
NAME_EDUCATION_TYPE 5
NAME_FAMILY_STATUS 5
NAME_HOUSING_TYPE 6
DAYS_BIRTH        7183
DAYS_EMPLOYED     3640
FLAG_MOBIL        1
FLAG_WORK_PHONE   2
FLAG_PHONE         2
FLAG_EMAIL         2
OCCUPATION_TYPE   18
CNT_FAM_MEMBERS   10
MONTHS_BALANCE    61
STATUS             8
dtype: int64
```

EDA (Exploratory Data Analysis)

In [20]: 1 #checking head of dataset
2 df.head()

Out[20]:

	ID	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_IN
0	5008804	M		Y		0
1	5008804	M		Y		0
2	5008804	M		Y		0
3	5008804	M		Y		0
4	5008804	M		Y		0

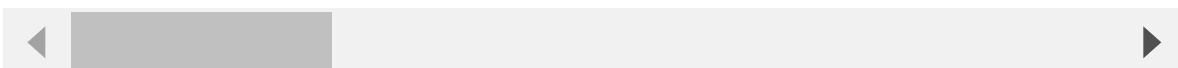


In [21]:

```
1 #checking tail od dataset
2 df.tail()
```

Out[21]:

ID	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AI
777710	5150337	M	N	Y	0
777711	5150337	M	N	Y	0
777712	5150337	M	N	Y	0
777713	5150337	M	N	Y	0
777714	5150337	M	N	Y	0



In [22]:

```
1 df.dtypes
```

Out[22]:

ID	int64
CODE_GENDER	object
FLAG_OWN_CAR	object
FLAG_OWN_REALTY	object
CNT_CHILDREN	int64
AMT_INCOME_TOTAL	float64
NAME_INCOME_TYPE	object
NAME_EDUCATION_TYPE	object
NAME_FAMILY_STATUS	object
NAME_HOUSING_TYPE	object
DAYS_BIRTH	int64
DAYS_EMPLOYED	int64
FLAG_MOBIL	int64
FLAG_WORK_PHONE	int64
FLAG_PHONE	int64
FLAG_EMAIL	int64
OCCUPATION_TYPE	object
CNT_FAM_MEMBERS	int64
MONTHS_BALANCE	int64
STATUS	object
dtype:	object

There are 9 categorical variable and 11 numerical variable

In [23]:

```
1 #checking Column names
2 df.columns
```

Out[23]:

```
Index(['ID', 'CODE_GENDER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'CNT_CHILDREN',
       'AMT_INCOME_TOTAL', 'NAME_INCOME_TYPE', 'NAME_EDUCATION_TYPE',
       'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE', 'DAYS_BIRTH',
       'DAYS_EMPLOYED', 'FLAG_MOBIL', 'FLAG_WORK_PHONE', 'FLAG_PHONE',
       'FLAG_EMAIL', 'OCCUPATION_TYPE', 'CNT_FAM_MEMBERS', 'MONTHS_BALANCE',
       'STATUS'],
      dtype='object')
```

In [24]:

```
1 #Descriptive Statistics
2 pd.options.display.float_format = "{:.2f}".format
3 df.describe()
```

Out[24]:

	ID	CNT_CHILDREN	AMT_INCOME_TOTAL	DAYS_BIRTH	DAYS_EMPLOYED	FI
count	777715.00	777715.00	777715.00	777715.00	777715.00	
mean	5078742.95	0.43	188534.80	-16124.94	57775.83	
std	41804.42	0.75	101622.45	4104.30	136471.74	
min	5008804.00	0.00	27000.00	-25152.00	-15713.00	
25%	5044568.50	0.00	121500.00	-19453.00	-3292.00	
50%	5069530.00	0.00	162000.00	-15760.00	-1682.00	
75%	5115551.00	1.00	225000.00	-12716.00	-431.00	
max	5150487.00	19.00	1575000.00	-7489.00	365243.00	



In [25]:

```

1 #Data Description
2 df.describe(include='all').T

```

Out[25]:

		count	unique	top	freq	mean	std	
	ID	777715.00	NaN	NaN	NaN	5078742.95	41804.42	50088
	CODE_GENDER	777715	2	F	518851	NaN	NaN	
	FLAG_OWN_CAR	777715	2	N	473355	NaN	NaN	
	FLAG_OWN_REALTY	777715	2	Y	512948	NaN	NaN	
	CNT_CHILDREN	777715.00	NaN	NaN	NaN	0.43	0.75	
	AMT_INCOME_TOTAL	777715.00	NaN	NaN	NaN	188534.80	101622.45	270
	NAME_INCOME_TYPE	777715	5	Working	400164	NaN	NaN	
	NAME_EDUCATION_TYPE	777715	5	Secondary / secondary special	524261	NaN	NaN	
	NAME_FAMILY_STATUS	777715	5	Married	546619	NaN	NaN	
	NAME_HOUSING_TYPE	777715	6	House / apartment	697151	NaN	NaN	
	DAYS_BIRTH	777715.00	NaN	NaN	NaN	-16124.94	4104.30	-251
	DAYS_EMPLOYED	777715.00	NaN	NaN	NaN	57775.83	136471.74	-157
	FLAG_MOBIL	777715.00	NaN	NaN	NaN	1.00	0.00	
	FLAG_WORK_PHONE	777715.00	NaN	NaN	NaN	0.23	0.42	
	FLAG_PHONE	777715.00	NaN	NaN	NaN	0.30	0.46	
	FLAG_EMAIL	777715.00	NaN	NaN	NaN	0.09	0.29	
	OCCUPATION_TYPE	537667	18	Laborers	131572	NaN	NaN	
	CNT_FAM_MEMBERS	777715.00	NaN	NaN	NaN	2.21	0.91	
	MONTHS_BALANCE	777715.00	NaN	NaN	NaN	-19.37	14.08	-
	STATUS	777715	8	C	329536	NaN	NaN	



CORRELATION MAP

In [26]:

```
1 fig, ax = plt.subplots(figsize=(25,10))
2 _ = sns.heatmap(df.corr(), annot=True, cmap=sns.cubehelix_palette(as_cmap=True))
```



INFERENCE:

- 1) FLAG_MOBIL Column has no relationship with any variables.
- 2) CNT_FAM_MEMBERS and CNT_CHILDREN have positive relationship with each other.
- 3) FLAG_WORK_PHONE and FLAG_PHONE have relationship.
- 4) CNT_FAM_MEMBERS and DAYS_BIRTH have relationship with each other.

PAIRPLOT

In []:

```
1 sns.pairplot(data=df)
```

CHECKING DUPLICATE VALUES

In [27]:

```
1 # Check Duplicate data
2 def check_duplicate(df):
3     if df.duplicated().all():
4         return 'There are duplicate data in Dataframe Need to be remove'
5     else :
6         return 'Data is clean,No duplicate data found.'
```

```
In [28]: 1 check_duplicate(df)
```

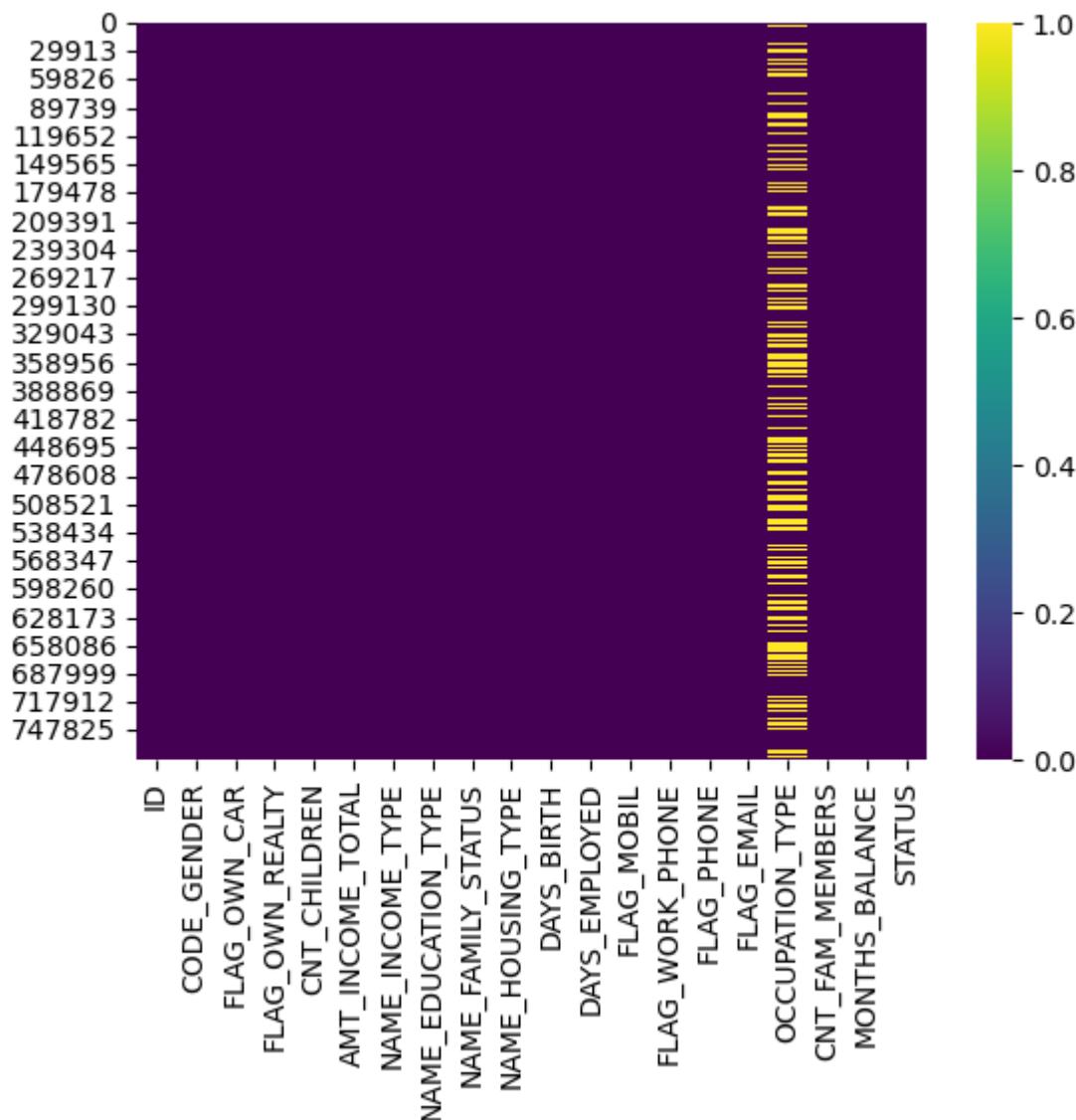
```
Out[28]: 'Data is clean, No duplicate data found.'
```

CHECKING MISSING VALUES

```
In [29]: 1 df.isnull().sum().sort_values(ascending=False)
```

```
Out[29]: OCCUPATION_TYPE      240048
ID                      0
CODE_GENDER                0
MONTHS_BALANCE              0
CNT_FAM_MEMBERS              0
FLAG_EMAIL                  0
FLAG_PHONE                  0
FLAG_WORK_PHONE              0
FLAG_MOBIL                  0
DAYS_EMPLOYED                 0
DAYS_BIRTH                  0
NAME_HOUSING_TYPE              0
NAME_FAMILY_STATUS              0
NAME_EDUCATION_TYPE              0
NAME_INCOME_TYPE                 0
AMT_INCOME_TOTAL                 0
CNT_CHILDREN                  0
FLAG_OWN_REALTY                 0
FLAG_OWN_CAR                  0
STATUS                      0
dtype: int64
```

```
In [30]: 1 sns.heatmap(df[df.columns].isnull(), cmap='viridis')
2 plt.show()
```



```
In [31]: 1 def percentage_of_miss():
2     df1=df[df.columns[df.isnull().sum()>=1]]# I get a subset of data that
3     total_miss = df1.isnull().sum().sort_values(ascending=False)
4     percent_miss = ((df1.isnull().sum()/df1.isnull().count())*100).sort_va
5     missing_data = pd.concat([total_miss, percent_miss], axis=1, keys=[ 'Nu
6     return(missing_data)
```

```
In [32]: 1 percentage_of_miss()
```

Out[32]:

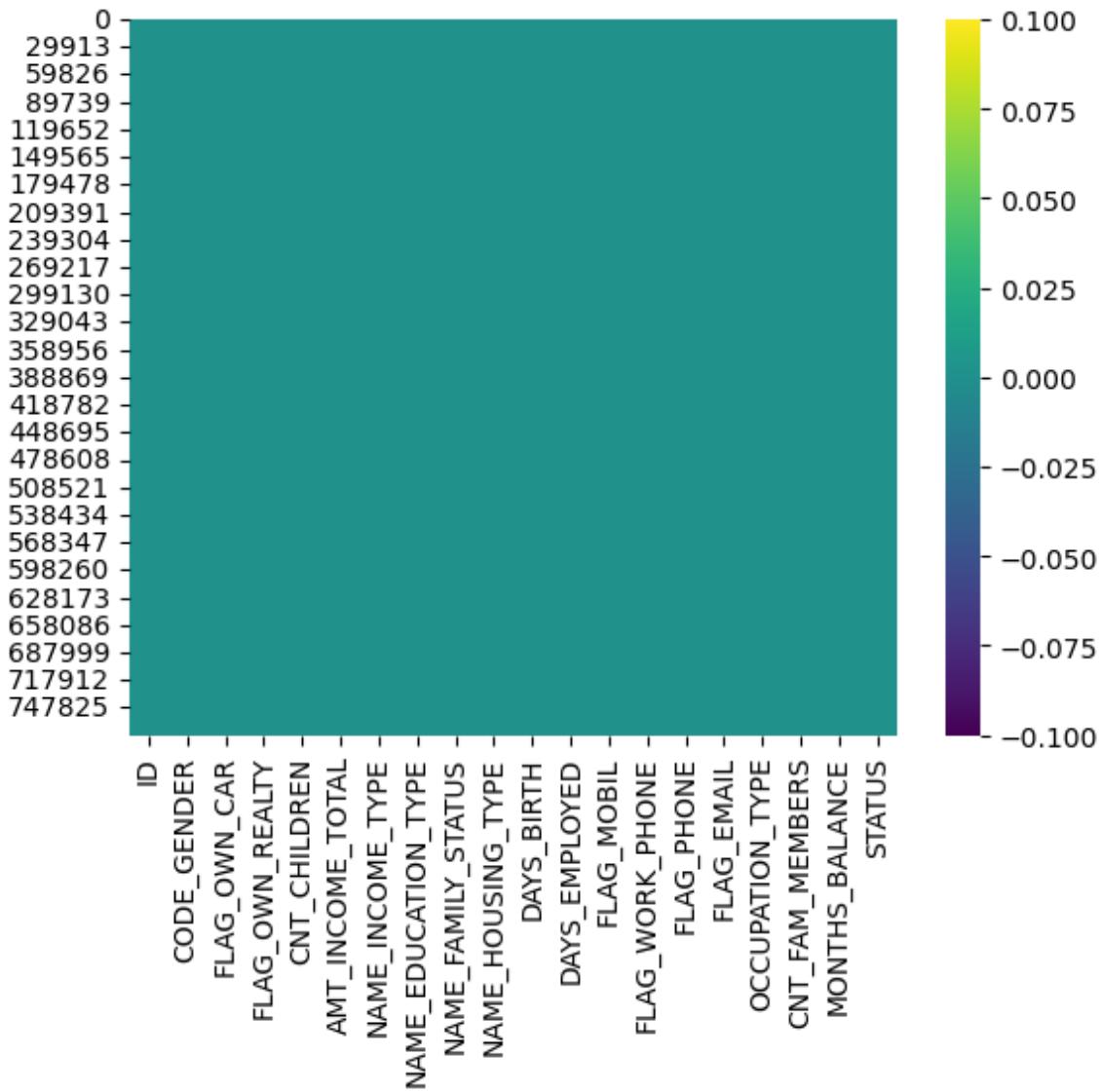
	Number of Missing	Percentage
OCCUPATION_TYPE	240048	30.87

Occupation type is having 240048 data missing - now ,we will replace NAN with others

```
In [33]: 1 #filling Null Values
2 df['OCCUPATION_TYPE'] =df['OCCUPATION_TYPE'].replace(np.nan, 'others')
```

```
In [34]: 1 #cols =df.columns
2 sns.heatmap(df[df.columns].isnull(), cmap='viridis')
```

Out[34]: <AxesSubplot:>



```
In [35]: 1 df.OCCUPATION_TYPE.value_counts()
```

```
Out[35]: others          240048
Laborers           131572
Core staff          77112
Sales staff          70362
Managers            67738
Drivers              47678
High skill tech staff 31768
Accountants         27223
Medicine staff       26691
Cooking staff        13416
Security staff       12400
Cleaning staff       11399
Private service staff 6714
Low-skill Laborers    3623
Secretaries          3149
Waiters/barmen staff 2557
HR staff              1686
IT staff              1319
Realty agents         1260
Name: OCCUPATION_TYPE, dtype: int64
```

NUMERICAL COLUMNS IN DATASET

```
In [36]: 1 # Get numerical variables
2 df_num=df.select_dtypes(exclude='object')# Just get me the numeric vari
3 df_num.columns
```

```
Out[36]: Index(['ID', 'CNT_CHILDREN', 'AMT_INCOME_TOTAL', 'DAYS_BIRTH', 'DAYS_EMPLOYED',
   'FLAG_MOBIL', 'FLAG_WORK_PHONE', 'FLAG_PHONE', 'FLAG_EMAIL',
   'CNT_FAM_MEMBERS', 'MONTHS_BALANCE'],
  dtype='object')
```

CATEGORICAL COLUMN IN DATASET

```
In [37]: 1 # Get categorical variables
2 df_cat=df.select_dtypes(include='object')# Just get me the categorical v
3 df_cat.columns
```

```
Out[37]: Index(['CODE_GENDER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'NAME_INCOME_TYPE',
   'NAME_EDUCATION_TYPE', 'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE',
   'OCCUPATION_TYPE', 'STATUS'],
  dtype='object')
```

UNIVARIATE ANALYSIS

LETS START CHECKING EACH COLUMN ONE BY ONE

DATA CLEANING - FIXING INVALID VALUES

Cleaning column ID

```
In [38]: 1 df.ID.value_counts()
```

```
Out[38]: 5090630    61
      5148524    61
      5066707    61
      5061848    61
      5118380    61
      ..
      5024557     1
      5062311     1
      5024365     1
      5024364     1
      5041568     1
Name: ID, Length: 36457, dtype: int64
```

```
In [ ]: 1 # Data is clean and has unique value for each customer.
```

Cleaning Column CODE_GENDER

```
In [39]: 1 df.CODE_GENDER.value_counts()
```

```
Out[39]: F    518851
      M    258864
Name: CODE_GENDER, dtype: int64
```

```
In [ ]: 1 # Data is clean Female is more compared to Male.
```

Cleaning column FLAG_own_car

```
In [40]: 1 df.FLAG_own_car.value_counts()
```

```
Out[40]: N    473355
      Y    304360
Name: FLAG_own_car, dtype: int64
```

```
In [ ]: 1 # Data is clean. People who own car is less compared to who doesnt own a
```

Cleaning column FLAG_own_REALTY

```
In [41]: 1 df.FLAG_own_REALTY.value_counts()
```

```
Out[41]: Y      512948  
N      264767  
Name: FLAG_own_REALTY, dtype: int64
```

```
In [ ]: 1 # Data is clean - Customers having own property apply for credit card is
```

Cleaning Column CNT_CHILDREN

```
In [42]: 1 df.CNT_CHILDREN.value_counts()
```

```
Out[42]: 0      540639  
1      155638  
2      70399  
3      9328  
4      1224  
5      324  
14     111  
7      46  
19     6  
Name: CNT_CHILDREN, dtype: int64
```

```
In [ ]: 1 # Data is Clean  
2 # But there are outliers which we will work on.
```

Cleaning Column AMT_INCOME_TOTAL

```
In [43]: 1 df.AMT_INCOME_TOTAL.value_counts()
```

```
Out[43]: 135000.00    90217  
180000.00    68579  
157500.00    62686  
112500.00    61622  
225000.00    61399  
...  
51750.00      7  
89550.00      5  
594000.00      4  
151200.00      2  
134995.50      2  
Name: AMT_INCOME_TOTAL, Length: 265, dtype: int64
```

```
In [ ]: 1 # Data is clean - This column gives the annual income of customers.
```

Cleaning Column NAME_INCOME_TYPE

```
In [44]: 1 df.NAME_INCOME_TYPE.value_counts()
```

```
Out[44]: Working          400164  
Commercial associate    183385  
Pensioner              128392  
State servant           65437  
Student                 337  
Name: NAME_INCOME_TYPE, dtype: int64
```

```
In [ ]: 1 # Data is Clean -- Working people applying for credit card is high.
```

Cleaning Column NAME_EDUCATION_TYPE

```
In [45]: 1 df.NAME_EDUCATION_TYPE.value_counts()
```

```
Out[45]: Secondary / secondary special    524261  
Higher education           213633  
Incomplete higher          30329  
Lower secondary             8655  
Academic degree            837  
Name: NAME_EDUCATION_TYPE, dtype: int64
```

In []: 1 # In the column NAME_EDUCATION_TYPE there is a value Secondary / seconda
2 # be value as Secondary / secondary special -- Lets clean it and keep a

In [46]: 1 def education_type(x):
2 if x == 'Secondary / secondary special' :
3 x= x.split(' /')[0]
4 return x

In [47]: 1 df['NAME_EDUCATION_TYPE'] = df['NAME_EDUCATION_TYPE'].apply(education_ty

In [48]: 1 df.NAME_EDUCATION_TYPE.value_counts()

Out[48]: Secondary 524261
Higher education 213633
Incomplete higher 30329
Lower secondary 8655
Academic degree 837
Name: NAME_EDUCATION_TYPE, dtype: int64

Cleaning Column NAME_FAMILY_STATUS

In [49]: 1 df.NAME_FAMILY_STATUS.value_counts()

Out[49]: Married 546619
Single / not married 94335
Civil marriage 60342
Separated 45255
Widow 31164
Name: NAME_FAMILY_STATUS, dtype: int64

In []: 1 #Data is clean

Cleaning Column NAME_HOUSING_TYPE

In [50]: 1 df.NAME_HOUSING_TYPE.value_counts()

Out[50]: House / apartment 697151
With parents 35735
Municipal apartment 24640
Rented apartment 10898
Office apartment 5636
Co-op apartment 3655
Name: NAME_HOUSING_TYPE, dtype: int64

```
In [ ]: 1 # In the column Name_housing_type there is a value House/apartment which  
2 # be House or Apartment -- Lets clean it and keep all the values as Hou  
3 # apartment but not a house option so I choose House.
```

```
In [51]: 1 def get_apartment(x):  
2     if x == 'House / apartment' :  
3         x=x.split(' /')[0]  
4     return x
```

```
In [52]: 1 df['NAME_HOUSING_TYPE'] = df['NAME_HOUSING_TYPE'].apply(get_apartment)
```

```
In [ ]: 1 df.NAME_HOUSING_TYPE.value_counts()
```

Cleaning Column DAYS_BIRTH

```
In [ ]: 1 #Getting DOB -- Changing negative values to readable date format.
```

```
In [ ]: 1 # Days_birth column name was transformed to Birth_day for better underst  
2 # date using timedelta.
```

```
In [53]: 1 df.DAYS_BIRTH.value_counts()
```

```
Out[53]: -14667    1018  
-15140     928  
-15675     835  
-15519     799  
-16995     799  
...  
-18983      1  
-10119      1  
-19526      1  
-7489       1  
-9177       1  
Name: DAYS_BIRTH, Length: 7183, dtype: int64
```

```
In [54]: 1 def date_of_birth (day_num):  
2     today = date.today()  
3     birthDay = (today + timedelta(days=day_num)).strftime('%Y-%m-%d')  
4     return birthDay
```

```
In [55]: 1 df['DAYS_BIRTH']=df['DAYS_BIRTH'].apply(date_of_birth)
```

Cleaning Column DAYS_EMPLOYED

```
In [ ]: 1 # Getting Employed days from days_employed and transform the column as E
2 # The negative values are transformed to readable date format.
```

```
In [56]: 1 df.DAYS_EMPLOYED.value_counts()
```

```
Out[56]: 365243    127972
-1751      1601
-1539      1545
-401       1498
-2531      1319
...
-3294       1
-3891       1
-7049       1
-7765       1
-2848       1
Name: DAYS_EMPLOYED, Length: 3640, dtype: int64
```

```
In [57]: 1 def day_employed(day_num):
2     today = date.today()
3     employedDay = (today + timedelta(days=day_num)).strftime('%Y-%m-%d')
4     result = 0
5     #if employedDay > date.today().strftime('%Y-%m-%d') :
6     #result = 0
7     #else:
8     #result = employedDay
9     #return result
10    return employedDay
```

```
In [58]: 1 df['DAYS_EMPLOYED'] = df['DAYS_EMPLOYED'].apply(day_employed)
```

In [59]: 1 df.DAYS_EMPLOYED.value_counts()

```
Out[59]: 3023-03-14    127972
          2018-05-27    1601
          2018-12-25    1545
          2022-02-05    1498
          2016-04-07    1319
          ...
          2014-03-06     1
          2012-07-17     1
          2003-11-24     1
          2001-12-08     1
          2015-05-26     1
Name: DAYS_EMPLOYED, Length: 3640, dtype: int64
```

In []: 1 # The Negative values are changed to readable date format.

Cleaning Column FLAG_MOBIL

In [60]: 1 df.FLAG_MOBIL.value_counts()

```
Out[60]: 1    777715
Name: FLAG_MOBIL, dtype: int64
```

In []: 1 # Data is clean -- Looks Like every one is having the phone.

Cleaning Column FLAG_WORK_PHONE

In [61]: 1 df.FLAG_WORK_PHONE.value_counts()

```
Out[61]: 0    597427
          1    180288
Name: FLAG_WORK_PHONE, dtype: int64
```

In []: 1 # Data is clean, People not having seperate work phone is high

Clean Column FLAG_PHONE

In [62]: 1 df.FLAG_PHONE.value_counts()

```
Out[62]: 0    543650
          1    234065
Name: FLAG_PHONE, dtype: int64
```

```
In [ ]: 1 # Data is clean Phone users are less.
```

Clean column FLAG_EMAIL

```
In [63]: 1 df.FLAG_EMAIL.value_counts()
```

```
Out[63]: 0    706418  
1     71297  
Name: FLAG_EMAIL, dtype: int64
```

```
In [ ]: 1 # Data is clean. Customer having Email is less as per the info collected
```

Cleaning Column OCCUPATION_TYPE

```
In [64]: 1 df.OCCUPATION_TYPE.value_counts()
```

```
Out[64]: others                240048  
Laborers               131572  
Core staff              77112  
Sales staff              70362  
Managers                 67738  
Drivers                  47678  
High skill tech staff   31768  
Accountants              27223  
Medicine staff            26691  
Cooking staff              13416  
Security staff             12400  
Cleaning staff              11399  
Private service staff    6714  
Low-skill Laborers        3623  
Secretaries                3149  
Waiters/barmen staff     2557  
HR staff                  1686  
IT staff                  1319  
Realty agents              1260  
Name: OCCUPATION_TYPE, dtype: int64
```

```
In [ ]: 1 # Data is clean. we replaced missing column with Others
```

Cleaning Column CNT_FAM_MEMBERS

```
In [65]: 1 df.CNT_FAM_MEMBERS.value_counts()
```

```
Out[65]: 2    423723
1    141477
3    134894
4    66990
5    8999
6    1196
7    273
15   111
9    46
20   6
Name: CNT_FAM_MEMBERS, dtype: int64
```

```
In [ ]: 1 # Data is clean. No missing values -- but it has outliers.
```

Cleaning Column MONTHS_BALANCE

```
In [ ]: 1 # The month of the extracted data is the starting point with 0 is the current month
2 # -1 is the previous month, and so on
```

```
In [66]: 1 df.MONTHS_BALANCE.value_counts()
```

```
Out[66]: -1    24963
-2    24871
0     24672
-3    24644
-4    24274
...
-56   1588
-57   1253
-58   955
-59   627
-60   321
Name: MONTHS_BALANCE, Length: 61, dtype: int64
```

Cleaning column STATUS

```
In [ ]: 1 ...
2 0: 1-29 days past due
3 1: 30-59 days past due
4 2: 60-89 days overdue
5 3: 90-119 days overdue
6 4: 120-149 days overdue
7 5: Overdue or bad debts, write-offs for more than 150 days
8 C: paid off that month
9 X: No loan for the month
10 'X': 6, 'C' : 7
11 ...
```

```
In [67]: 1 df.STATUS.value_counts()
```

```
Out[67]: C    329536
0    290654
X    145950
1     8747
5     1527
2      801
3      286
4      214
Name: STATUS, dtype: int64
```

```
In [ ]: 1 # There are 2 values as X and C -- we can replace it with numerical value
```

```
In [68]: 1 df['STATUS'].replace({'X': 6, 'C' : 7}, inplace=True)
2 df['STATUS']=df['STATUS'].astype(int)
```

```
In [ ]: 1 df.STATUS.value_counts()
```

FEATURE ENGINEERING

```
In [ ]: 1 # Trying to get new features or new column to our dataset to add more information
```

In [69]:

```
1 #Lets check the column names
2 df.columns
```

Out[69]:

```
Index(['ID', 'CODE_GENDER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'CNT_CHILDREN',
       'AMT_INCOME_TOTAL', 'NAME_INCOME_TYPE', 'NAME_EDUCATION_TYPE',
       'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE', 'DAYS_BIRTH',
       'DAYS_EMPLOYED', 'FLAG_MOBIL', 'FLAG_WORK_PHONE', 'FLAG_PHONE',
       'FLAG_EMAIL', 'OCCUPATION_TYPE', 'CNT_FAM_MEMBERS', 'MONTHS_BALANCE',
       'STATUS'],
      dtype='object')
```

In []:

```
1 # Columns 'DAYS_BIRTH', 'DAYS_EMPLOYED' are not clear for better understanding
2 # as BIRTH_DAY and EMPLOYED_DAY
```

In [70]:

```
1 df['BIRTH_DAY']=df['DAYS_BIRTH']
2 df['EMPLOYED_DAY']=df['DAYS_EMPLOYED']
```

In []:

```
1 # The dataset has no age -- we can calculate age which might be used in
```

In [71]:

```
1 def age_calculations(born):
2     born = datetime.strptime(born, '%Y-%m-%d')
3     today = date.today()
4     return today.year - born.year - ((today.month, today.day) < (born.month, born.day))
```

In [72]:

```
1 df['AGE']=df['BIRTH_DAY'].apply(age_calculations)
```

In []:

```
1 # A new column - Age has been created.
```

In []:

```
1 # we can create a target column based on the status column.
2 # The customers who are delay in overdue >= 60 days are assigned as YES
3 # The target column can be used to predict which customer can be targeted
```

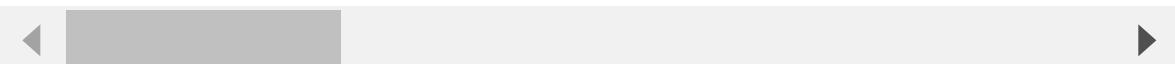
```
In [73]: 1 def target_calculations(x):
2     target=''
3     if x in (2,3,4,5):
4         target="NO" # delinquent customer - Risky
5     else:
6         target="YES" # punctual customer - Not Risky
7     return target
```

```
In [74]: 1 df['TARGET'] = df['STATUS'].apply(target_calculations)
2 df
```

Out[74]:

	ID	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AI
0	5008804	M	Y	Y	0	
1	5008804	M	Y	Y	0	
2	5008804	M	Y	Y	0	
3	5008804	M	Y	Y	0	
4	5008804	M	Y	Y	0	
...
777710	5150337	M	N	Y	0	
777711	5150337	M	N	Y	0	
777712	5150337	M	N	Y	0	
777713	5150337	M	N	Y	0	
777714	5150337	M	N	Y	0	

777715 rows × 24 columns



```
In [75]: 1 df.TARGET.value_counts()
```

```
Out[75]: YES    774887
NO      2828
Name: TARGET, dtype: int64
```

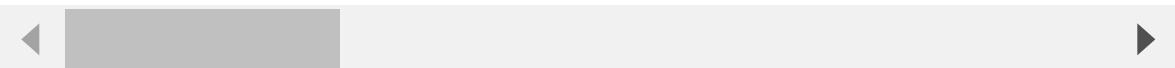
DROPING COLUMNS

```
In [76]: 1 # Taking copy of original dataframe before dropping the columns
          2 df_orginal=df.copy()
          3 df_orginal
```

Out[76]:

	ID	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AI
0	5008804	M	Y	Y	0	
1	5008804	M	Y	Y	0	
2	5008804	M	Y	Y	0	
3	5008804	M	Y	Y	0	
4	5008804	M	Y	Y	0	
...
777710	5150337	M	N	Y	0	
777711	5150337	M	N	Y	0	
777712	5150337	M	N	Y	0	
777713	5150337	M	N	Y	0	
777714	5150337	M	N	Y	0	

777715 rows × 24 columns



```
In [ ]: 1 df.columns
```

```
In [77]: 1 df = df.drop(['ID',
          2             'DAYS_BIRTH',
          3             'DAYS_EMPLOYED',
          4             'BIRTH_DAY',
          5             'FLAG_WORK_PHONE',
          6             'FLAG_EMAIL',
          7             'EMPLOYED_DAY',
          8             'FLAG_MOBIL'],
          9             axis=1)
```

```
In [78]: 1 #checking the shape of the column
          2 df.shape
```

Out[78]: (777715, 16)

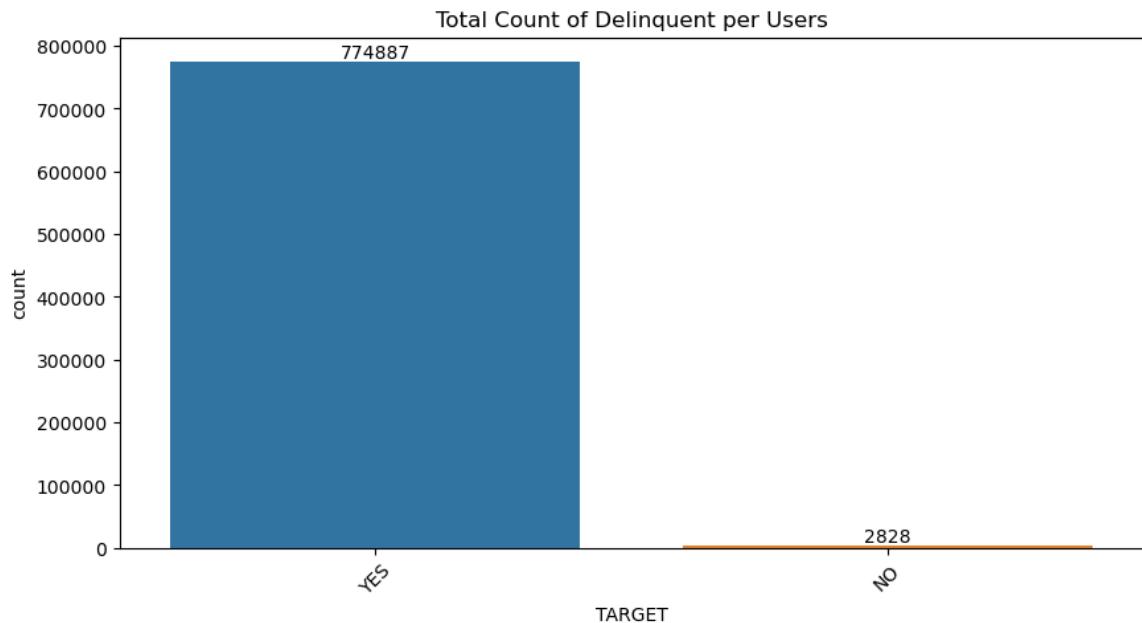
BIVARIATE ANALYSIS

DATA VISUALIZATION

```
In [ ]: 1 # we will analyse and visualize each part of data to be in near step fro  
2 # we will sense best factor that affect on our business goal
```

```
In [79]: 1 def categorical_plotting(df,col,title):  
2     fig, ax = plt.subplots(figsize=(10,5))  
3     ax=sns.countplot(x=col, data=df, order = df[col].value_counts().index)  
4     ax.set_xticklabels(ax.get_xticklabels(), rotation=45)  
5     ax.bar_label(ax.containers[0])  
6     plt.title(title)  
7     plt.show()
```

```
In [80]: 1 categorical_plotting(df,'TARGET','Total Count of Delinquent per Users')
```

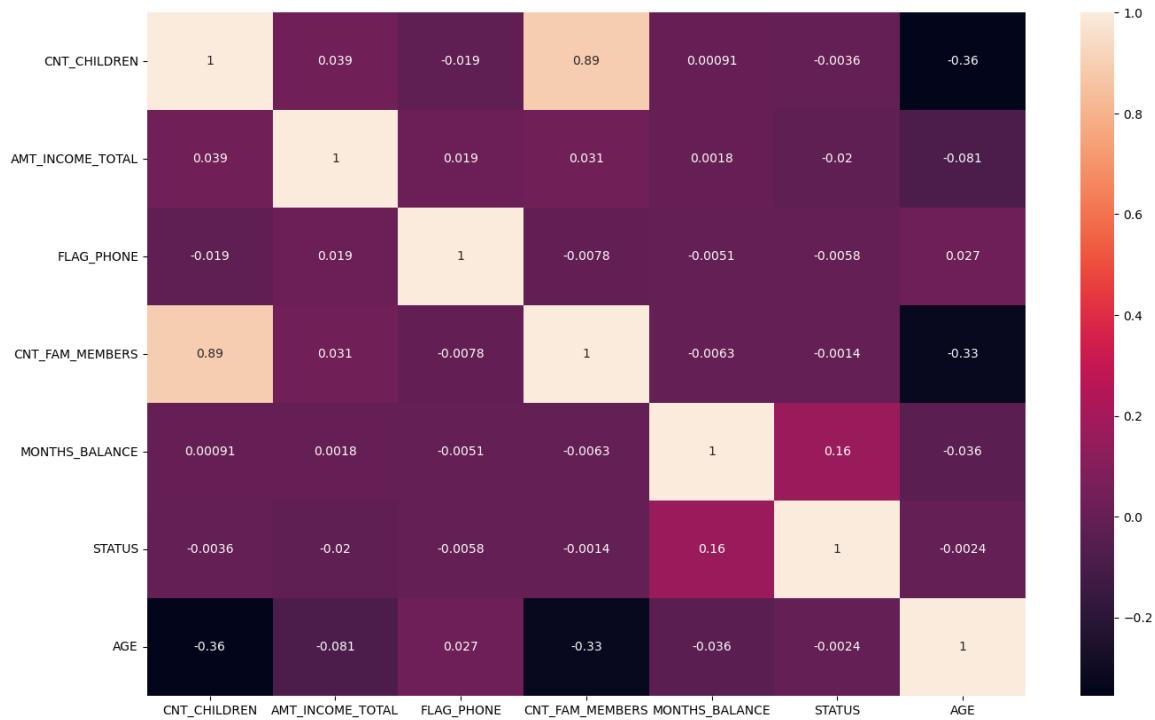


In [81]:

```

1 plt.figure(figsize = (16,10))
2 sns.heatmap(df.corr(), annot= True)
3 plt.show()

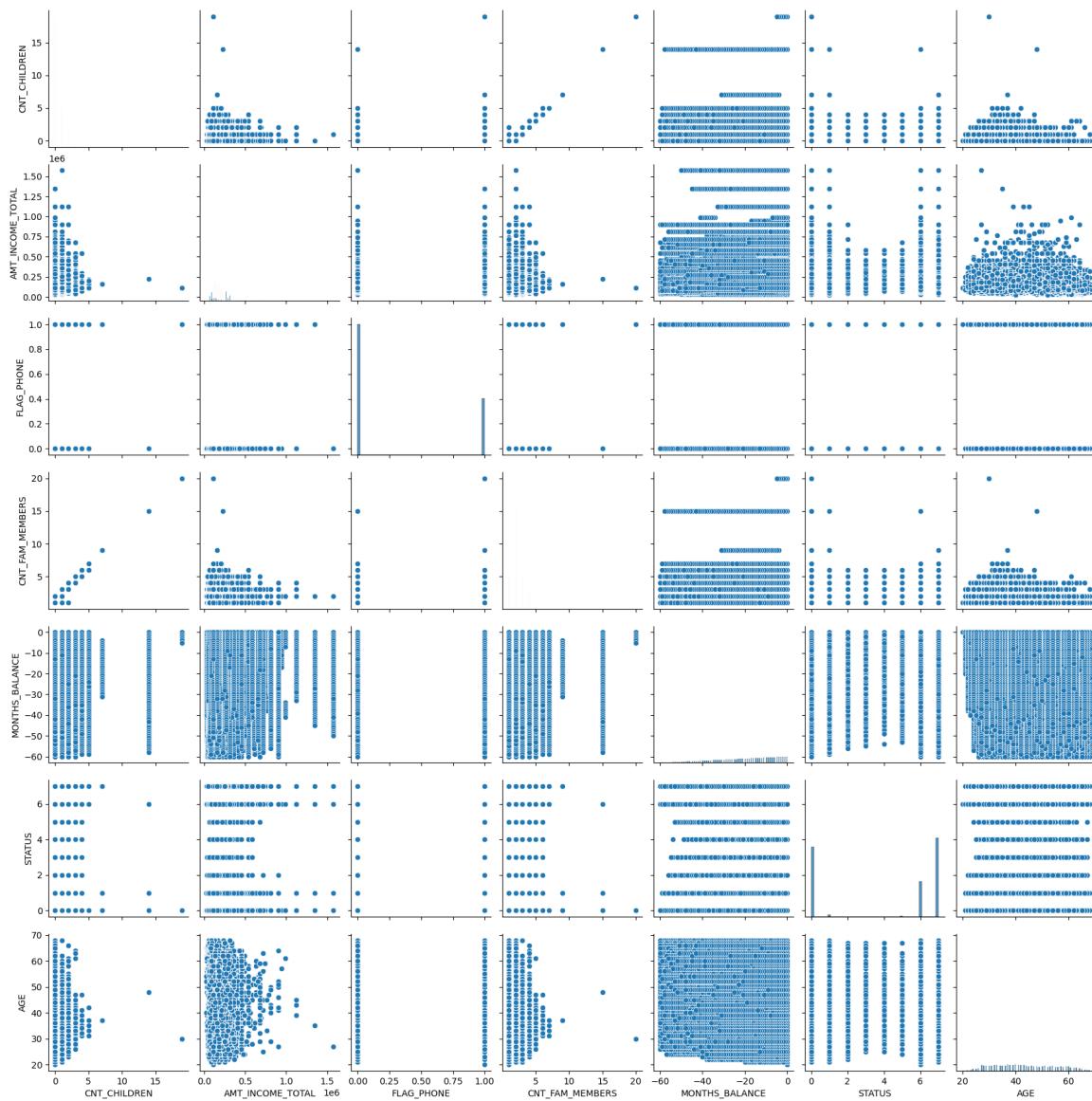
```



- 1) Strong Relationship Correlation between CNT_CHILDREN and CNT_FAM_MEMBERS

In [82]: 1 sns.pairplot(data=df)

Out[82]: <seaborn.axisgrid.PairGrid at 0x2220d97d9d0>



In [83]: 1 #get all numerical column
2 numCol = [col for col in df.columns if df[col].dtype != "O"]
3 numCol

Out[83]: ['CNT_CHILDREN',
'AMT_INCOME_TOTAL',
'FLAG_PHONE',
'CNT_FAM_MEMBERS',
'MONTHS_BALANCE',
'STATUS',
'AGE']

Visualizing CNT_CHILDREN

In [84]:

```

1 def numerical_plotting(df, col, title, symb):
2     fig, ax = plt.subplots(2, 1, sharex=True, figsize=(8,5), gridspec_kw=
3         {'height_ratios': [1, 1], 'hspace': 0.05})
4     ax[0].set_title(title, fontsize=18)
5     sns.boxplot(x=col, data=df, ax=ax[0])
6     ax[0].set(yticks=[])
7     sns.distplot(df[col], kde=True)
8     plt.xticks(rotation=45)
9     ax[1].set_xlabel(col, fontsize=10)
10    #plt.axvline(df[col].mean(), color='darkgreen', linewidth=2.2, label='mean')
11    #plt.axvline(df[col].median(), color='red', linewidth=2.2, label='median')
12    #plt.axvline(df[col].mode()[0], color='purple', linewidth=2.2, label='mode')
13    plt.legend(bbox_to_anchor=(1, 1.03), ncol=1, fontsize=17, fancybox=True)
14    plt.tight_layout()
15    plt.show()

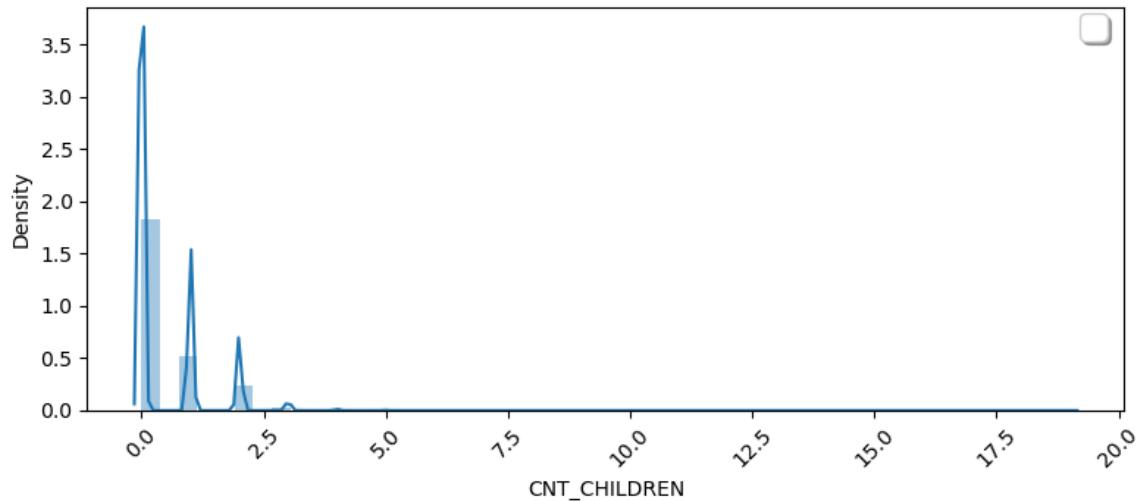
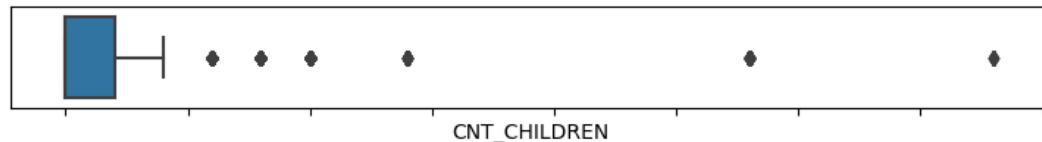
```

In [85]:

```
1 numerical_plotting(df, 'CNT_CHILDREN', 'Count of Children', ' ')
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.

Count of Children



In []:

```
1 # There are outliers in the column
```

TREATING OUTLIERS

In [86]:

```
1 def outlier_detection(dftemp,var):
2     sns.boxplot(x=df[var])
3     mean1=df[var].mean()
4     std1=df[var].std()
5     Min1=mean1-3*std1
6     Max1=mean1+3*std1
7     print("Mean - 3 * standard deviation:",Min1,"\\nMean + 3 * standard d
8     if ((max(dftemp[var]) <Max1)& (min(dftemp[var])>Min1)):
9         print("\\nThere is no outlier in '%s' based on 'Mean and Standard
10    else: print("\\nThere are some outliers in '%s' based on 'Mean and St
11    return dftemp[var].describe()
```

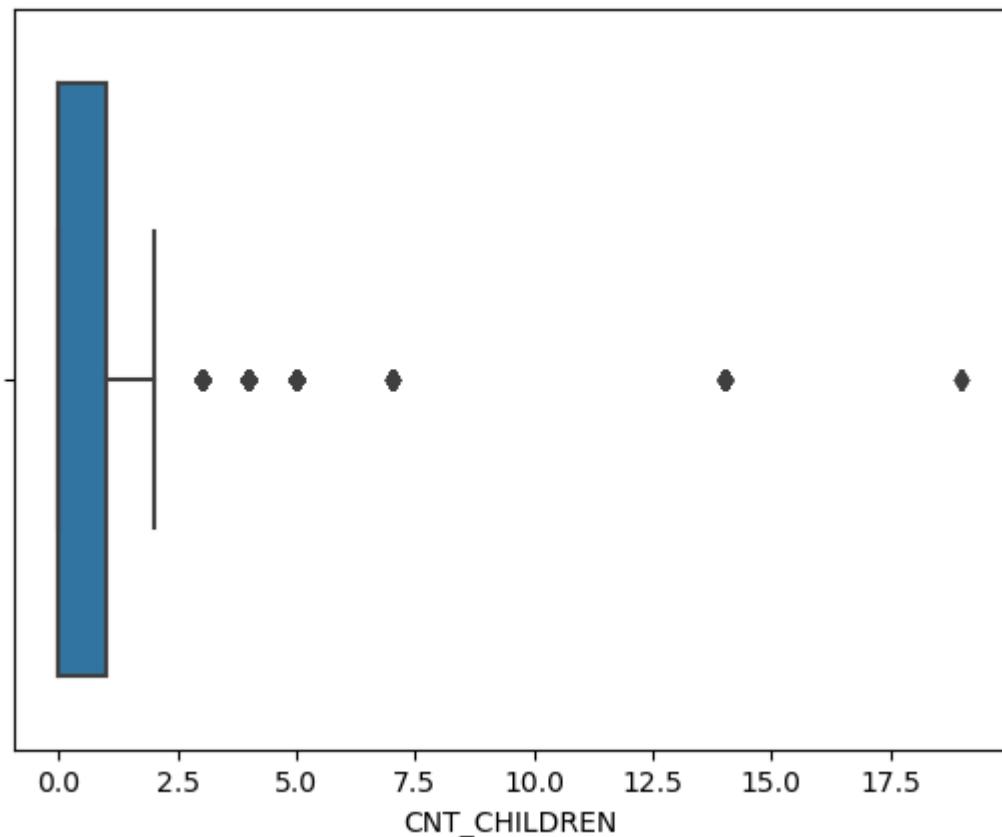


```
In [87]: 1 outlier_detection(df, 'CNT_CHILDREN')
```

```
Mean - 3 * standard deviation: -1.8091833482788853
Mean + 3 * standard deviation: 2.665347881559071
```

There are some outliers in 'CNT_CHILDREN' based on 'Mean and Standard Deviation Method' in the dataset that you give to this function

```
Out[87]: count    777715.00
mean        0.43
std         0.75
min         0.00
25%         0.00
50%         0.00
75%         1.00
max        19.00
Name: CNT_CHILDREN, dtype: float64
```



In [88]:

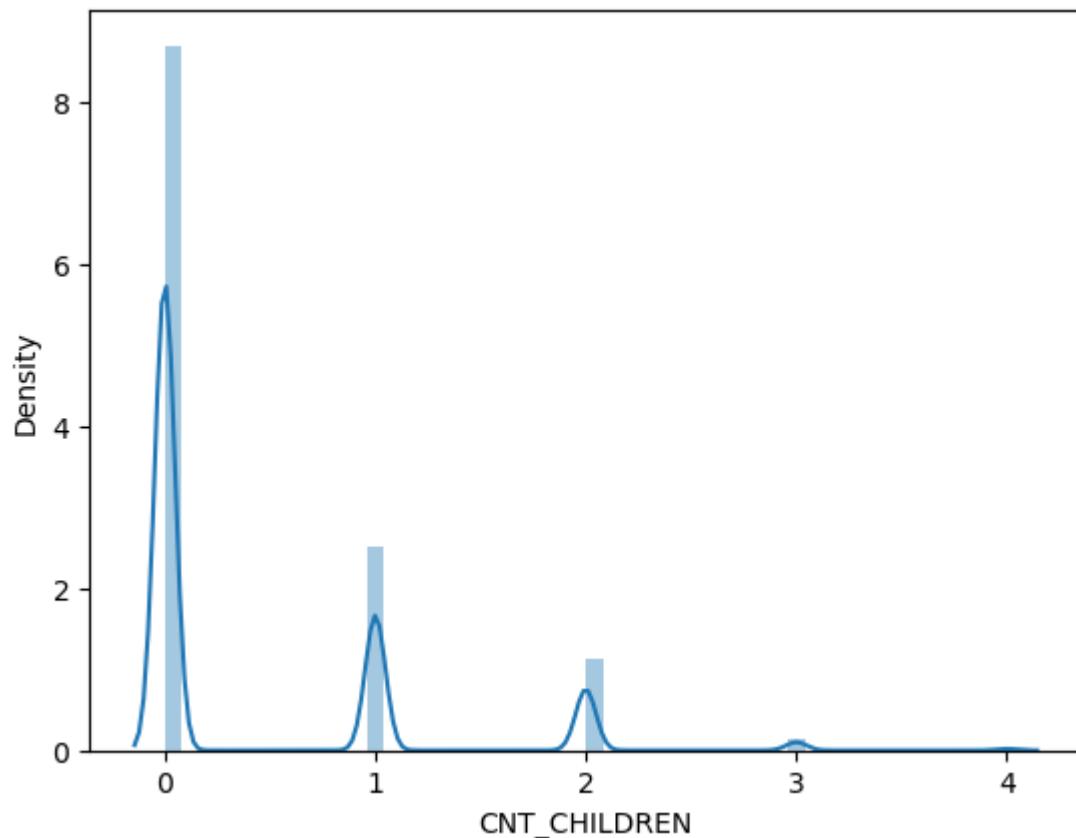
```
1 #I want to replace values bigger than upper outer fence with upper outer
2 #lower outer fence
3 #Lower outer fence: Q1 - 3*IQR
4 #upper outer fence: Q3 + 3*IQR
5 def fixing_outliers(col):#col is name of column that you want to fix its
6     i=0 # i is number of outliers that they are less than lower outer fe
7     j=0 # j is number of outliers that they are bigger than upper outer
8     Q1=df[col].quantile(0.25) # finding 1st quartile(lower quartile)
9     Q3=df[col].quantile(0.75) # finding 3rd quartile(upper quartile)
10    IQR=Q3-Q1 #calculating Inter Quartile Range
11    UOF1=Q3+3*(IQR)#UOF=upper outer fence
12    LOF1=Q1-3*(IQR)#LOF=lower outer fence
13    Clean_col=[] # Create an empty list to append value to it
14    for val in df[col]:#val is value of that column . here we select val
15        if val<LOF1:
16            Clean_col.append(LOF1)# if value is less than LOF we will re
17            i+=1
18        elif val>UOF1 : # if value is bigger than UOF we will repalce t
19            Clean_col.append(UOF1)
20            j+=1
21        else : Clean_col.append(val)
22    df[col]=Clean_col # update the column by assigning Clean_col
23    sns.distplot(df[col])#Plotting univariate distributions
24    print(f'Number of outliers that they are less than lower outer fence')
25    return df[col] .describe()
```



```
In [89]: 1 fixing_outliers('CNT_CHILDREN')
```

Number of outliers that they are less than lower outer fence(-3.0): 0
Number of outliers that they are bigger than upper outer fence (4.0) : 487

```
Out[89]: count    777715.00
mean        0.43
std         0.72
min         0.00
25%        0.00
50%        0.00
75%        1.00
max         4.00
Name: CNT_CHILDREN, dtype: float64
```



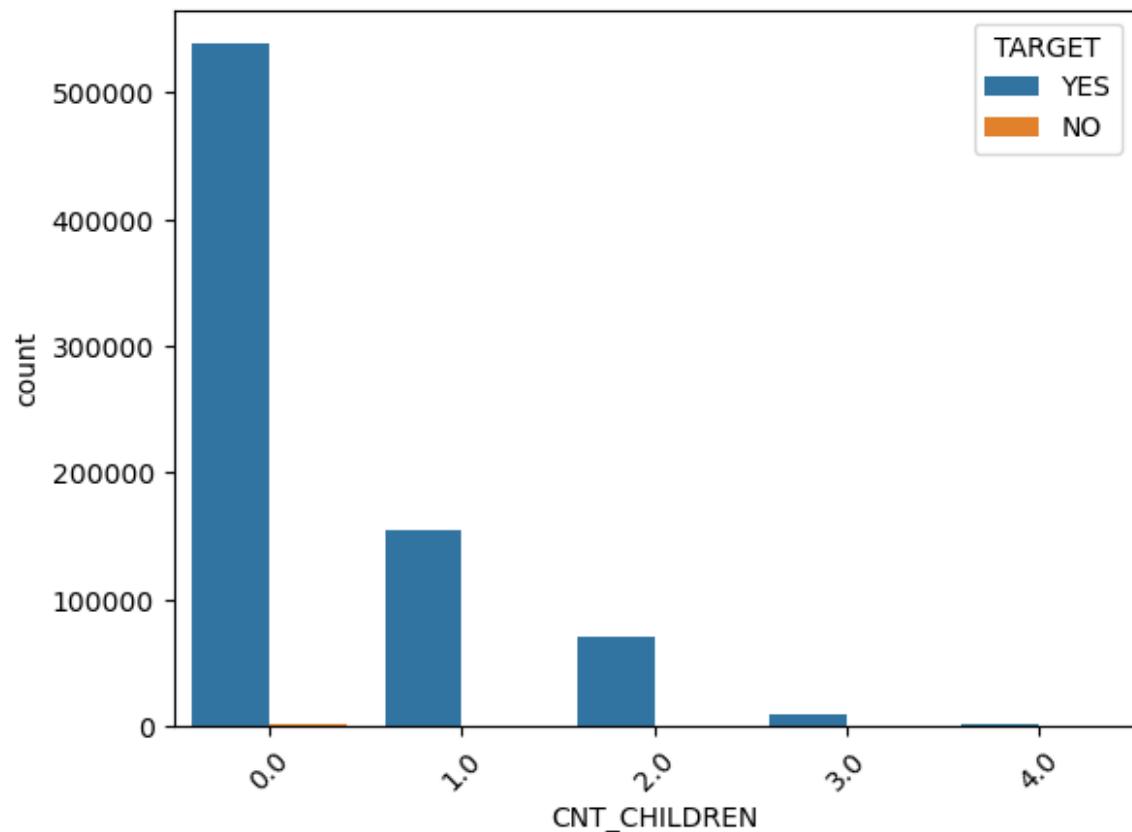
```
In [90]: 1 df['CNT_CHILDREN'].describe()
```

```
Out[90]: count    777715.00
mean        0.43
std         0.72
min         0.00
25%        0.00
50%        0.00
75%        1.00
max         4.00
Name: CNT_CHILDREN, dtype: float64
```

Compare the CNT_CHILDREN with the Target column¶

In [91]:

```
1 sns.countplot(x='CNT_CHILDREN', data=df, hue='TARGET')
2 plt.xticks(rotation=45)
3 plt.show()
```



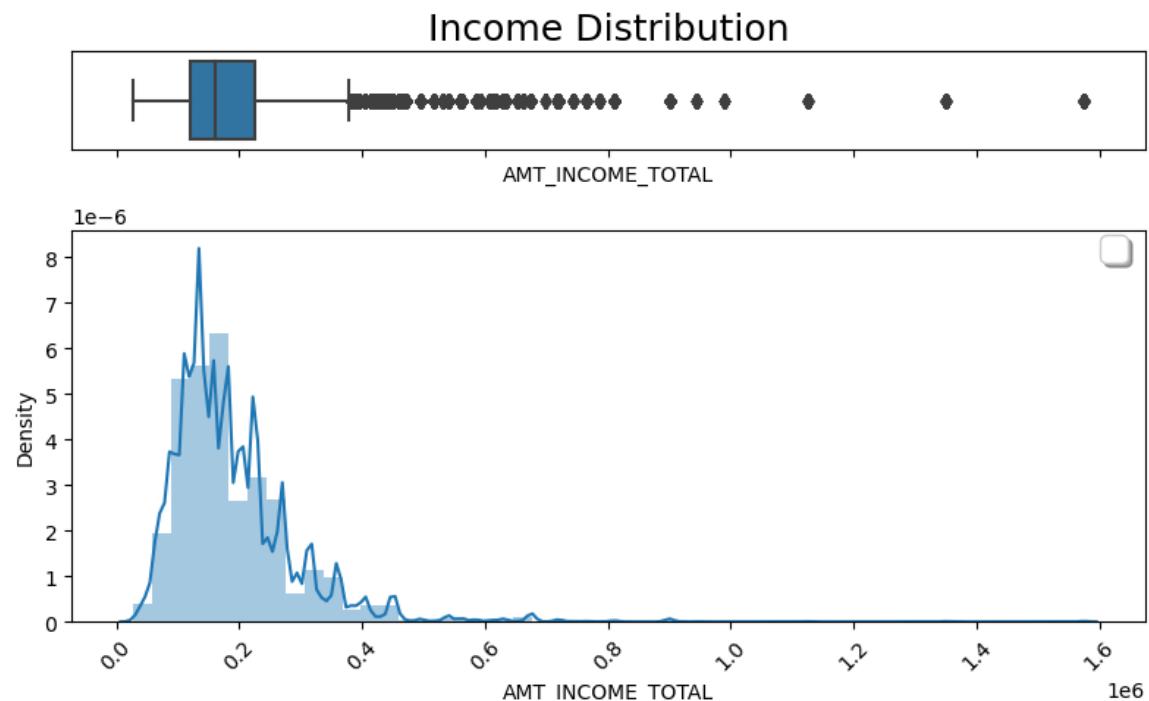
NOTES:

- 1) 70% users are single persons.

Visualizing AMT_INCOME_TOTAL

```
In [92]: 1 numerical_plotting(df, 'AMT_INCOME_TOTAL', 'Income Distribution', '$')
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



```
In [93]: 1 # what is 5 high total income ?
2 df['AMT_INCOME_TOTAL'].value_counts().sort_values(ascending=False).head()
```

```
Out[93]: 135000.00      90217
180000.00      68579
157500.00      62686
112500.00      61622
225000.00      61399
Name: AMT_INCOME_TOTAL, dtype: int64
```

```
In [94]: 1 # what is the total income with the target ?
2 df.groupby(['AMT_INCOME_TOTAL']).max()[['TARGET']].sort_values(by = 'AMT_
```

```
Out[94]:
```

TARGET	
AMT_INCOME_TOTAL	
1575000.00	YES
1350000.00	YES
1125000.00	YES
990000.00	YES
945000.00	YES

NOTES:

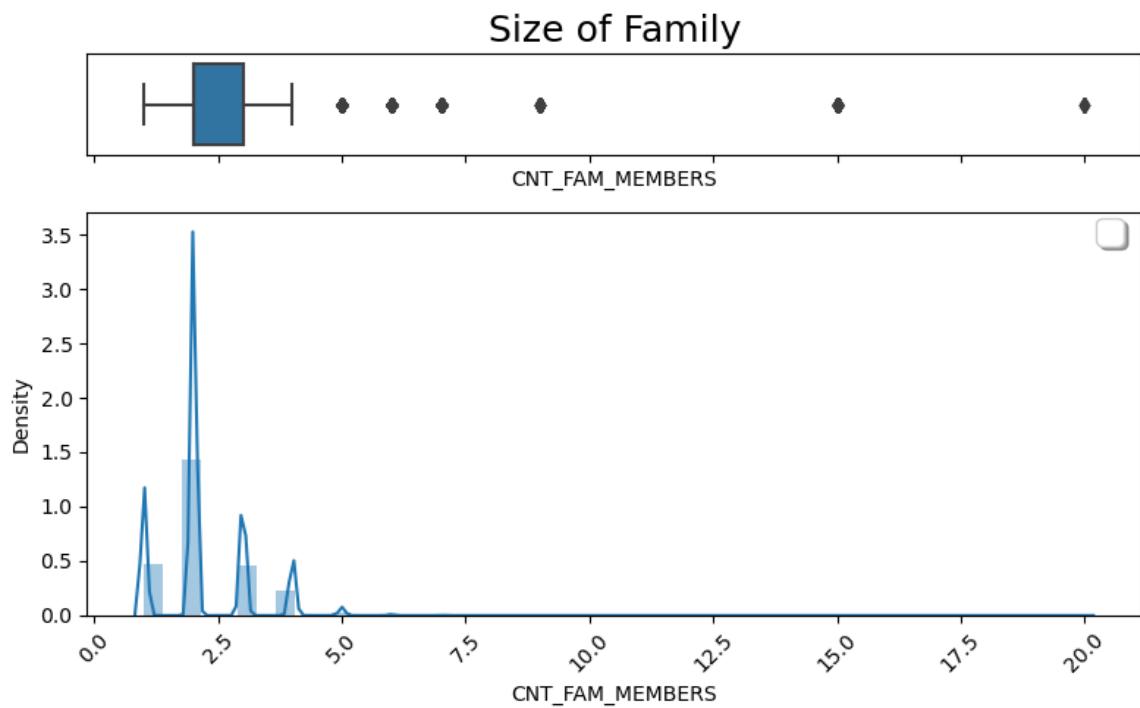
- 1) Total income highly repeated is \$135000.

```
In [ ]: 1 numCol
```

Visualizing CNT_FAM_MEMBERS

```
In [95]: 1 numerical_plotting(df, 'CNT_FAM_MEMBERS', 'Size of Family', '')
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



```
In [ ]: 1 # This column has Outliers.Lets treat the outliers
```

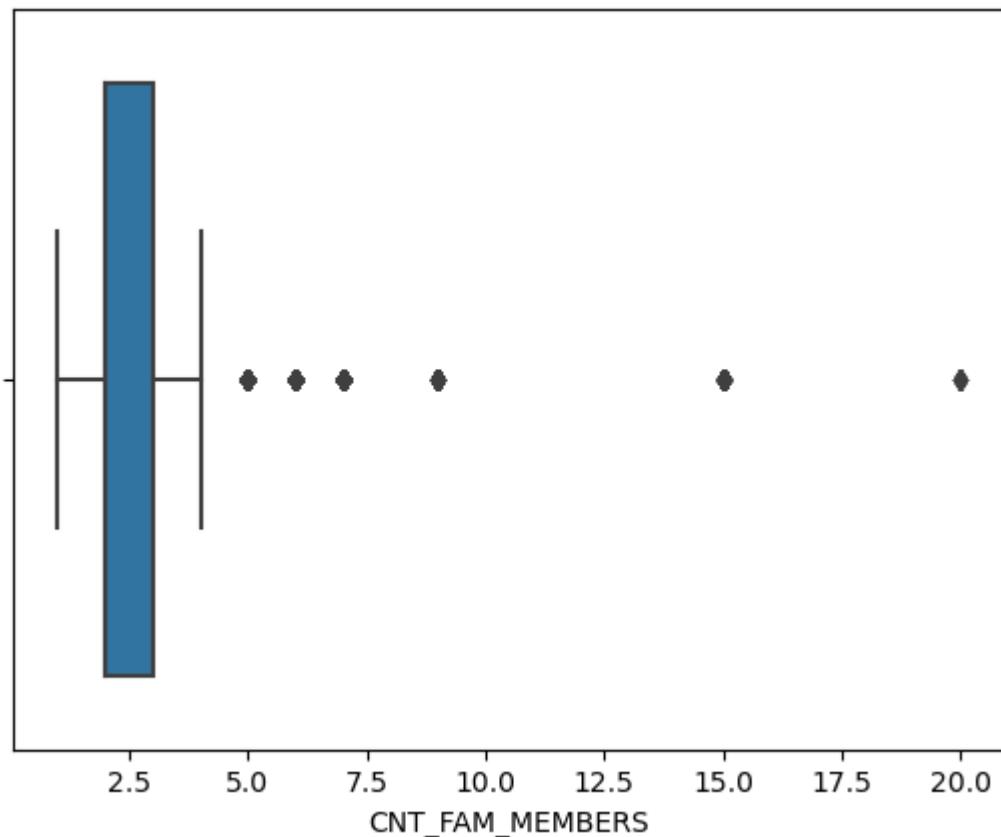
TREATING OUTLIERS

```
In [96]: 1 outlier_detection(df, 'CNT_FAM_MEMBERS')
```

```
Mean - 3 * standard deviation: -0.5133017431224771
Mean + 3 * standard deviation: 4.9309765983072165
```

There are some outliers in 'CNT_FAM_MEMBERS' based on 'Mean and Standard Deviation Method' in the dataset that you give to this function

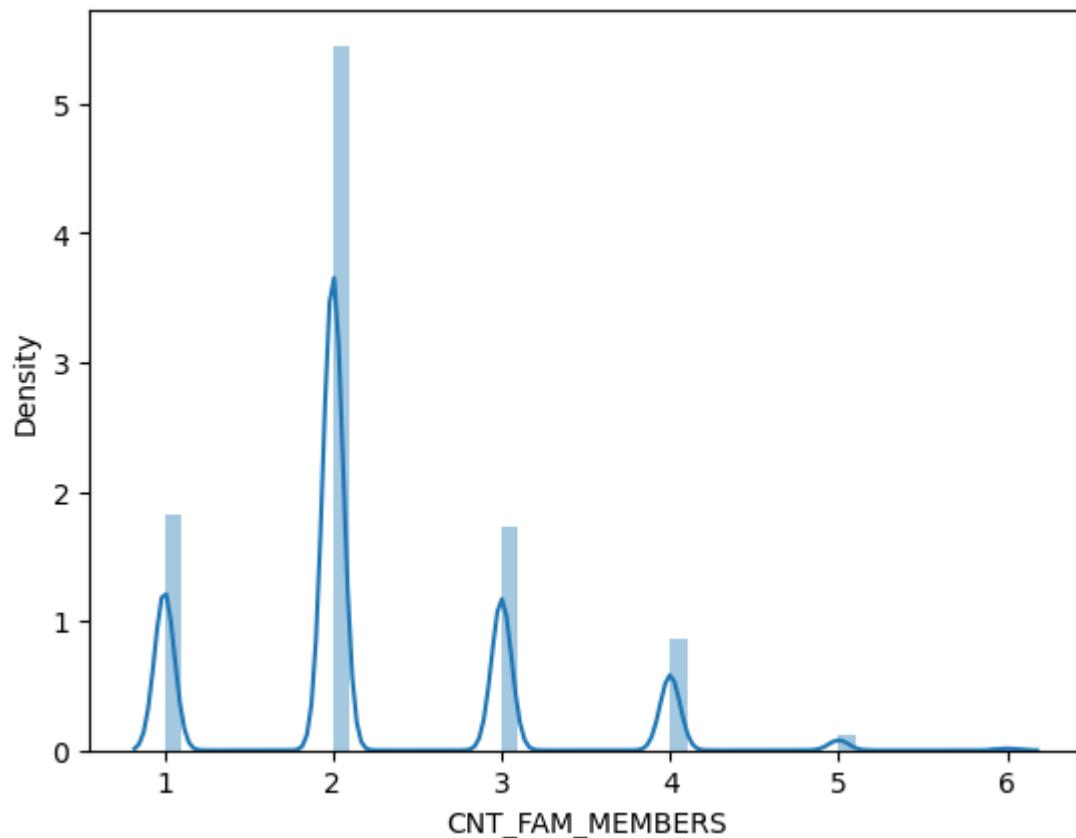
```
Out[96]: count    777715.00
mean        2.21
std         0.91
min         1.00
25%        2.00
50%        2.00
75%        3.00
max        20.00
Name: CNT_FAM_MEMBERS, dtype: float64
```



```
In [97]: 1 fixing_outliers('CNT_FAM_MEMBERS')
```

Number of outliers that they are less than lower outer fence(-1.0): 0
Number of outliers that they are bigger than upper outer fence (6.0) : 436

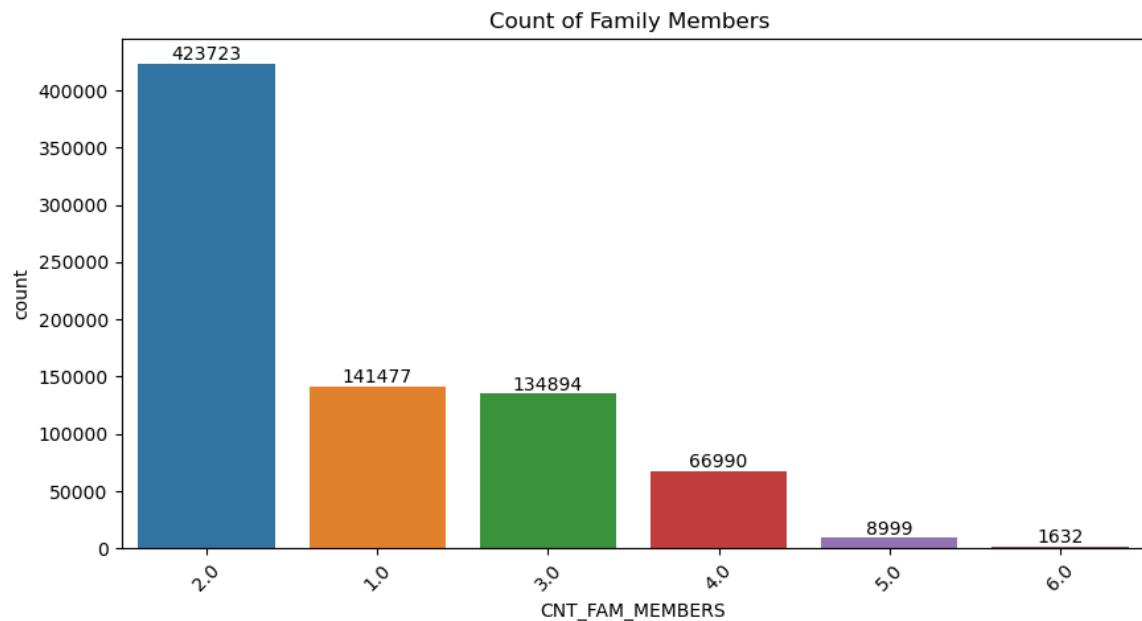
```
Out[97]: count    777715.00
mean        2.21
std         0.89
min         1.00
25%        2.00
50%        2.00
75%        3.00
max         6.00
Name: CNT_FAM_MEMBERS, dtype: float64
```



```
In [98]: 1 df['CNT_FAM_MEMBERS'].describe()
```

```
Out[98]: count    777715.00
mean        2.21
std         0.89
min         1.00
25%        2.00
50%        2.00
75%        3.00
max         6.00
Name: CNT_FAM_MEMBERS, dtype: float64
```

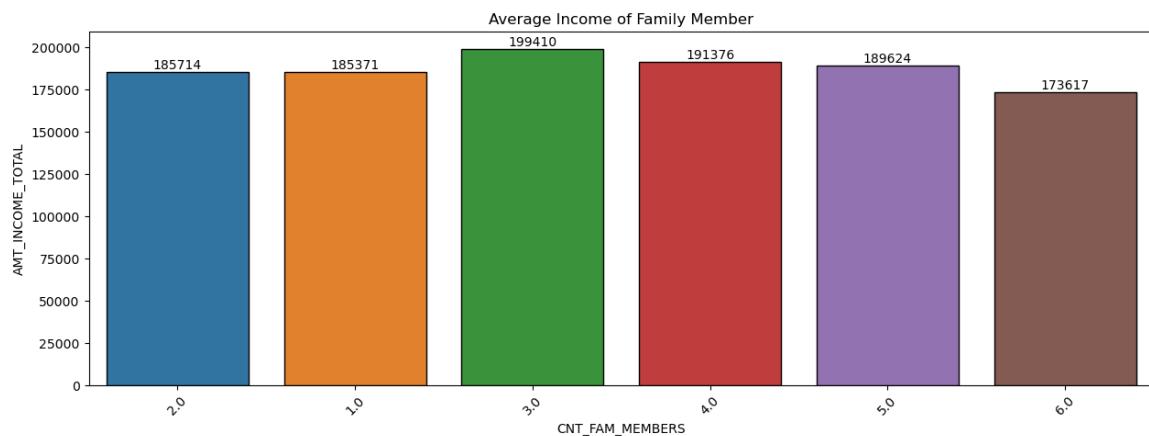
In [99]: 1 categorical_plotting(df, 'CNT_FAM_MEMBERS', 'Count of Family Members ')



COMPARE CNT_FAM_MEMBERS with AMT_INCOME_TOTAL

In [100]: 1 def average_plotting(df,col,output,number,title):
2 data_list = df[col].value_counts().index[:number].tolist()
3 plt.figure(figsize=(15,5))
4 ax=sns.barplot(x=col, y=output, data=df[df[col].isin(data_list)],order=data_list)
5 plt.xticks(rotation=45);
6 ax.bar_label(ax.containers[0])
7 plt.title(title)
8 plt.show()

In [101]: 1 average_plotting(df, 'CNT_FAM_MEMBERS', 'AMT_INCOME_TOTAL', 15, 'Average Income of Family Member')

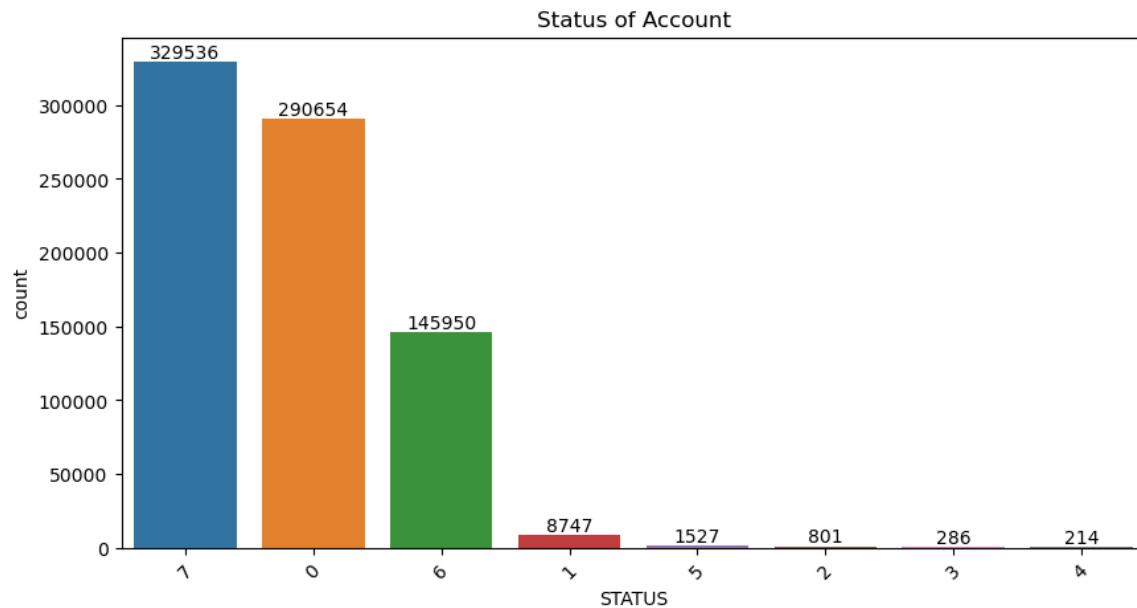


NOTES:

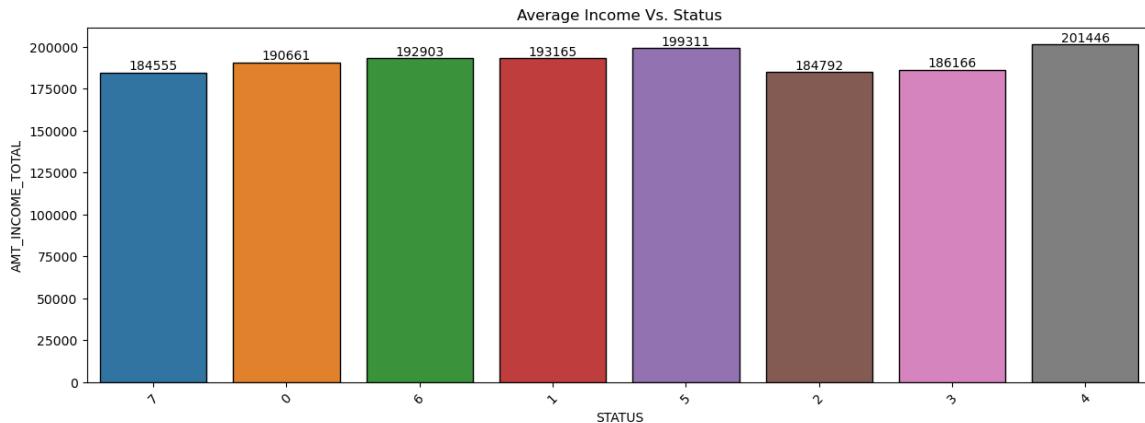
- 1) There are +ve correlation between CNT_CHILDREN and CNT_FAM_MEMBERS
- 2) Average Income of the Family Members is High for 3 members per family followed by 4 per family.
- 3) Average Income of Family Size 2 , 1 and 5 are around same with slight difference.
- 4) The Family Size 6 is comparatively low when compared to other family size.

Visualizing STATUS

```
In [102]: 1 categorical_plotting(df, 'STATUS', 'Status of Account')
```



```
In [103]: 1 average_plotting(df, 'STATUS', 'AMT_INCOME_TOTAL', 15, 'Average Income Vs. Status')
```



In []:

```
1  ...
2  0: 1-29 days past due
3  1: 30-59 days past due
4  2: 60-89 days overdue
5  3: 90-119 days overdue
6  4: 120-149 days overdue
7  5: Overdue or bad debts, write-offs for more than 150 days
8  C: paid off that month
9  X: No loan for the month
10 'X': 6, 'C' : 7
11 ...
```

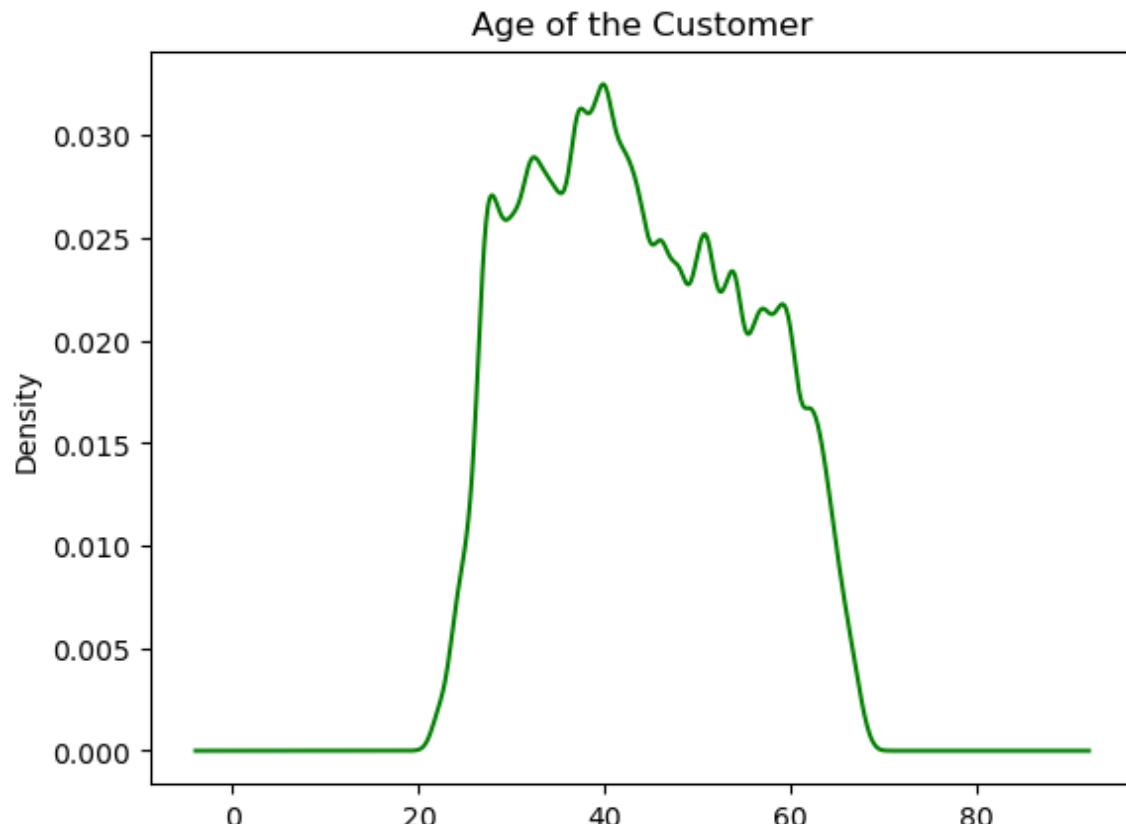
NOTES

- 1) Customers whose overdue is 5 is very high which is status 4.
- 2) Category 5 is second highest where the overdue debts are assigned to credit department.
- 3) Category 6 have no loan for the month.

VISUALIZING AGE

In [104]:

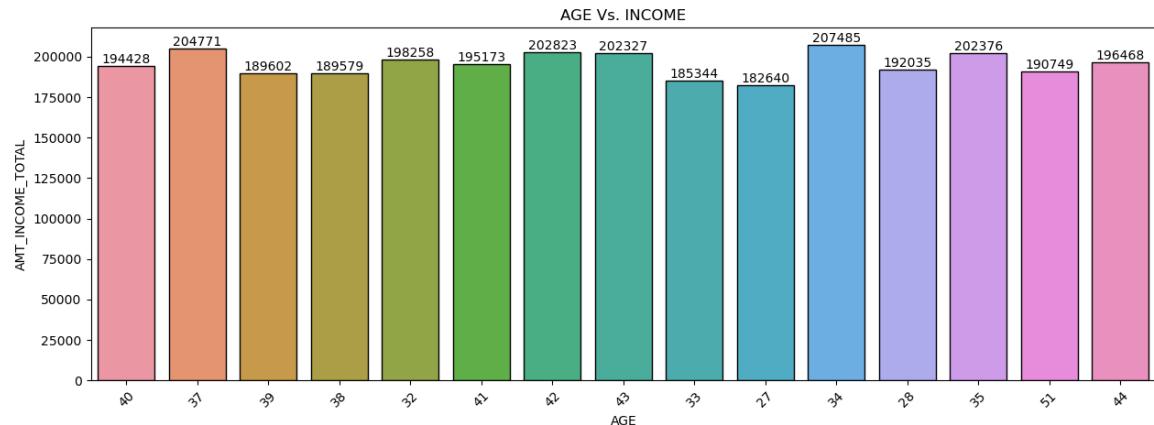
```
1 df.AGE.plot.density(color='green')
2 plt.title('Age of the Customer')
3 plt.show()
```



NOTES:

- 1) Age has no outliers.
- 2) Customers in age range from 30- 40 applying for Credit card is high.

In [105]: 1 average_plotting(df, 'AGE', 'AMT_INCOME_TOTAL', 15, 'AGE Vs. INCOME')



NOTES

- 1) Average Salary is high for the age 34 followed by 35,32.
- 2) Average Salary of customers on 40's is also high compared to age in 20's

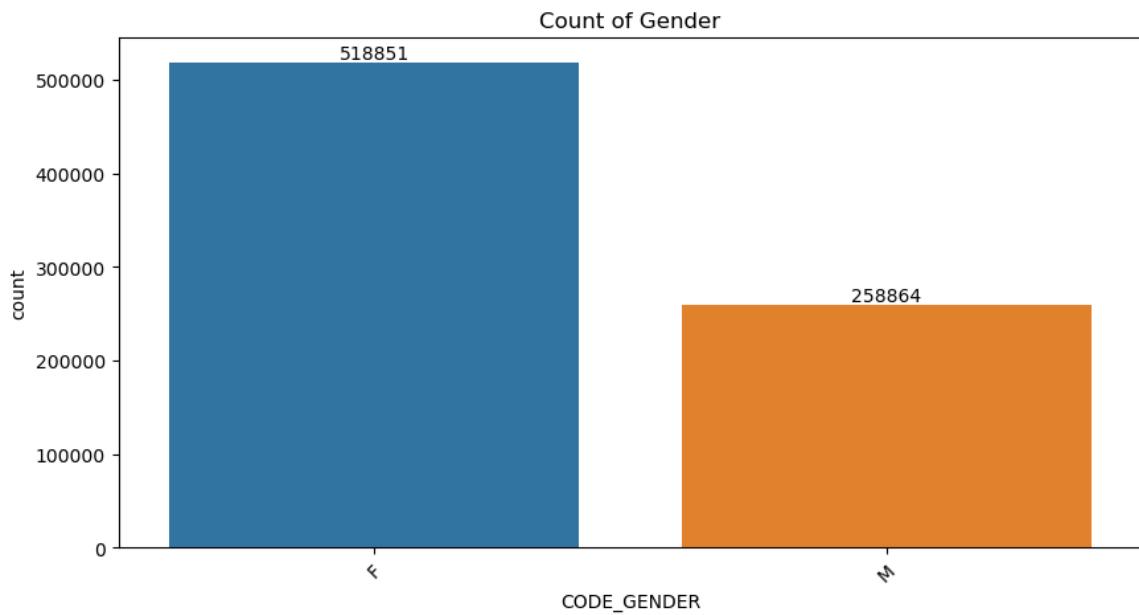
CATEGORICAL VARIABLE

In [106]: 1 catCols = [col for col in df.columns if df[col].dtype=="O"]
2 catCols

Out[106]: ['CODE_GENDER',
'FLAG_OWN_CAR',
'FLAG_OWN_REALTY',
'NAME_INCOME_TYPE',
'NAME_EDUCATION_TYPE',
'NAME_FAMILY_STATUS',
'NAME_HOUSING_TYPE',
'OCCUPATION_TYPE',
'TARGET']

Visualizing CODE_GENDER Column

In [107]: 1 categorical_plotting(df, 'CODE_GENDER', 'Count of Gender')

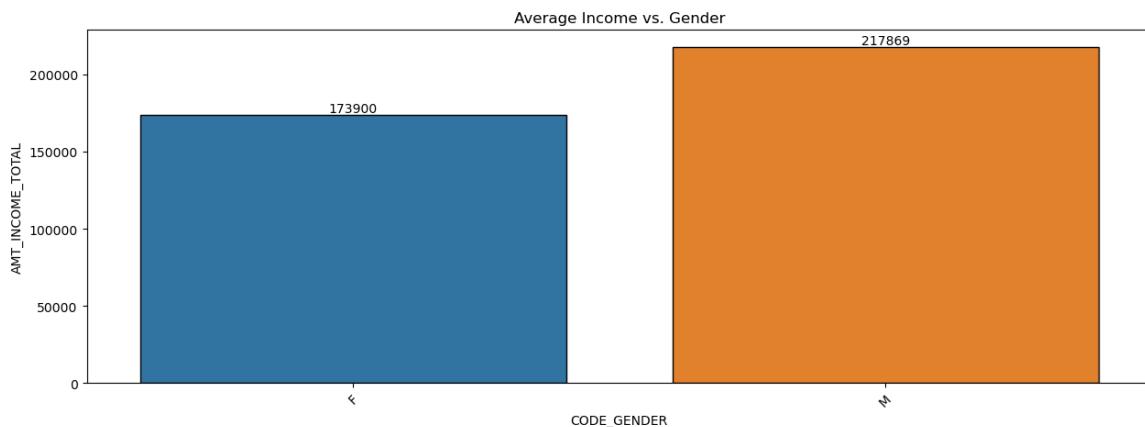


In [108]: 1 pd.options.display.float_format = "{:.2f}".format
2 df.groupby('CODE_GENDER').mean()[['AMT_INCOME_TOTAL']].sort_values(by='A

Out[108]:

	CODE_GENDER	M	F
AMT_INCOME_TOTAL	217868.58	173899.66	

In [109]: 1 average_plotting(df, 'CODE_GENDER', 'AMT_INCOME_TOTAL', 15, 'Average Income

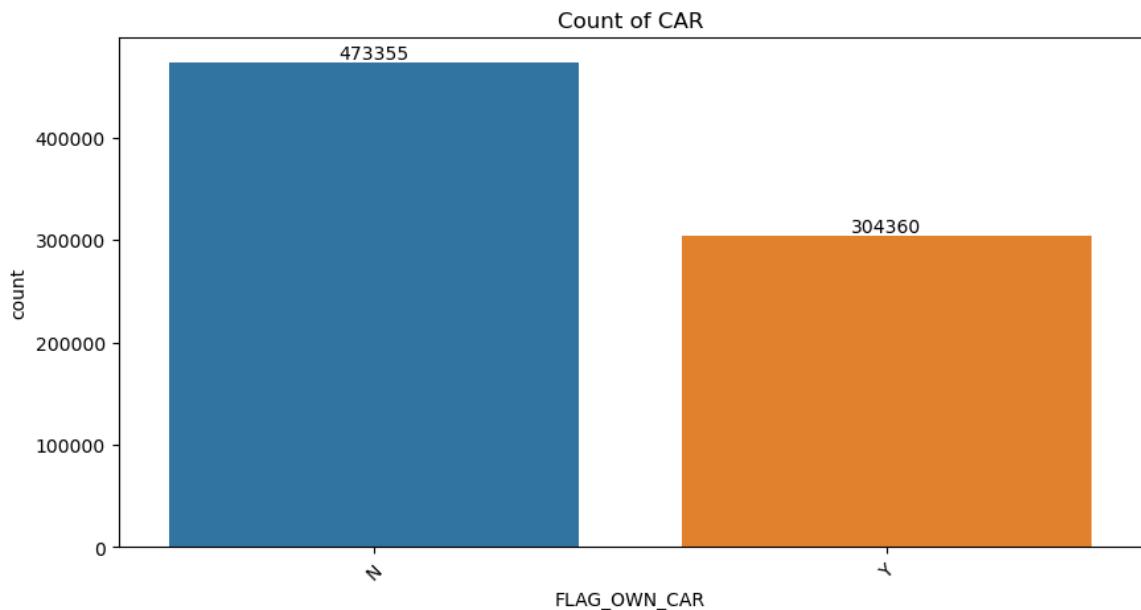


NOTES

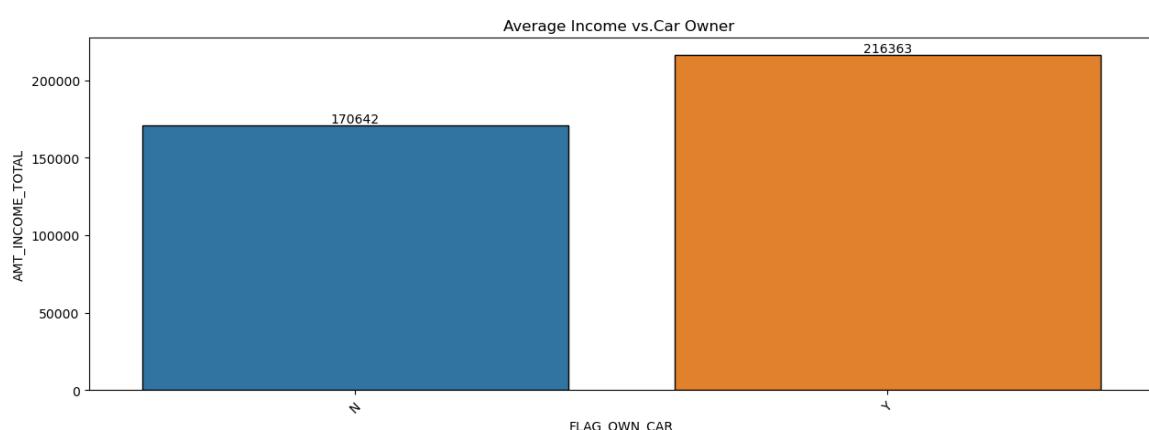
- 1) Average Income of Female is 172580 whereas Male is 214989
- 2) 66% of users are female.
- 3) Male user income is high when compared to Female users.

Visualizing FLAG OWN CAR Column

```
In [110]: 1 categorical_plotting(df, 'FLAG OWN CAR', 'Count of CAR')
```



```
In [111]: 1 average_plotting(df, 'FLAG OWN CAR', 'AMT_INCOME_TOTAL', 15, 'Average Income vs. Car Owner')
```

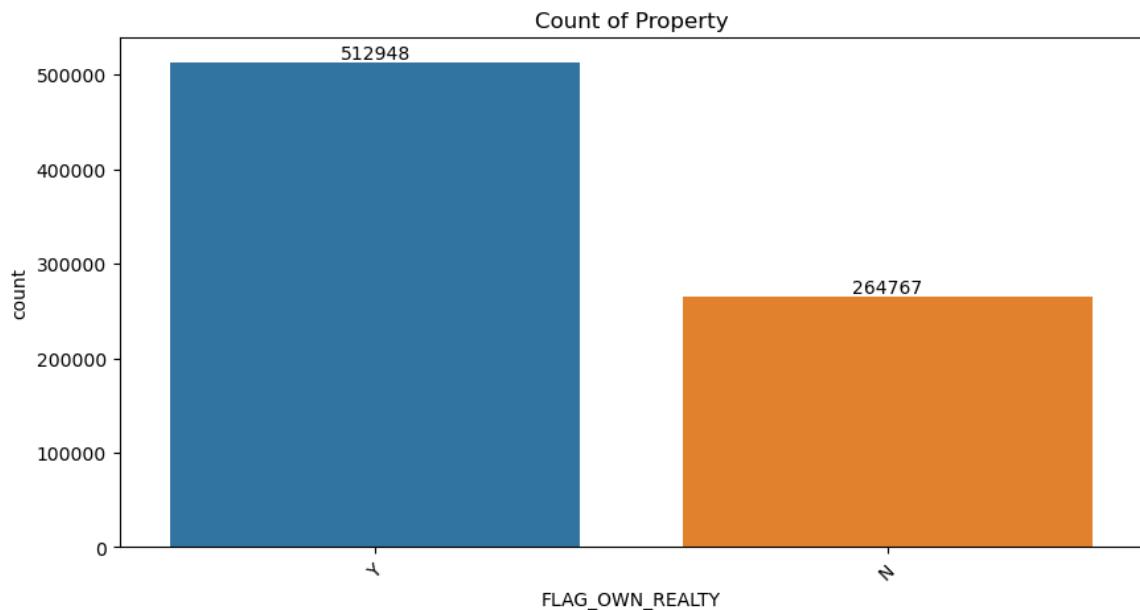


NOTES

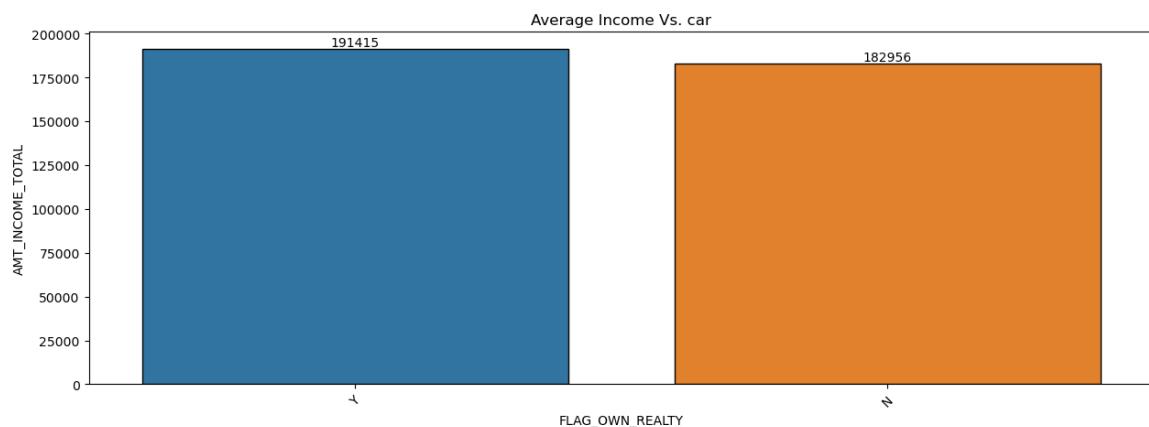
- 1) 40% of Users has car

Visualize FLAG OWN REALTY Column

```
In [113]: 1 categorical_plotting(df, 'FLAG OWN REALTY', 'Count of Property')
```



```
In [114]: 1 average_plotting(df, 'FLAG OWN REALTY', 'AMT_INCOME_TOTAL', 15, 'Average Inc')
```

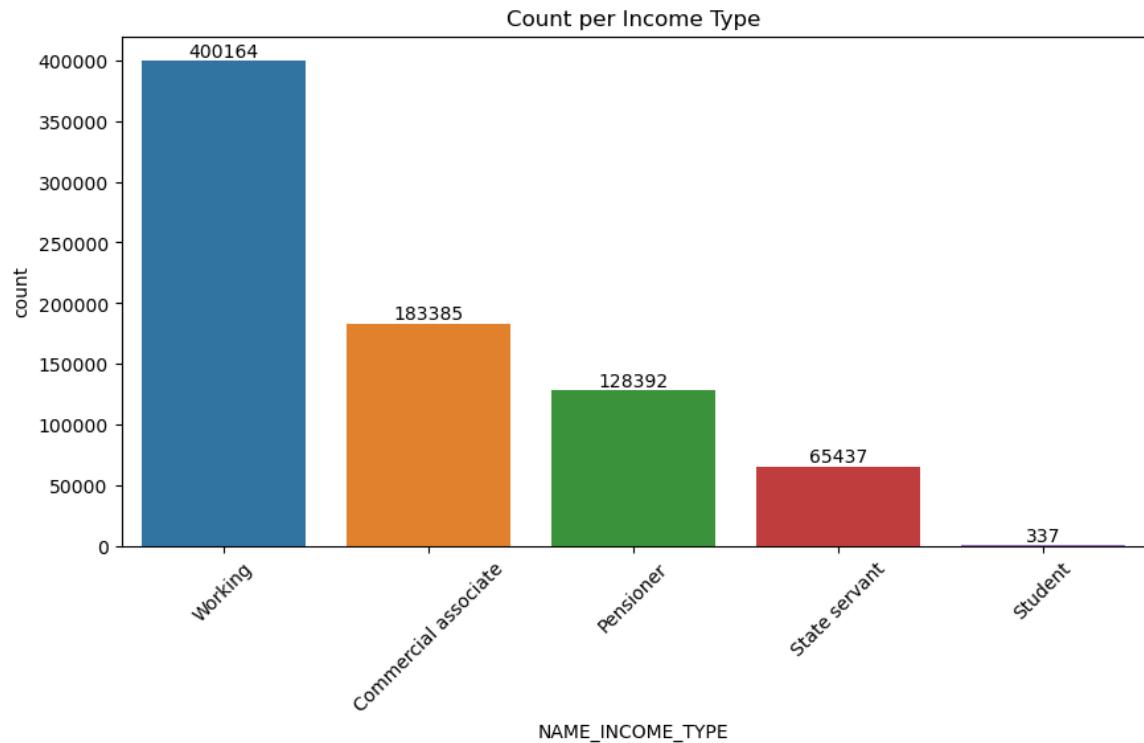


Notes

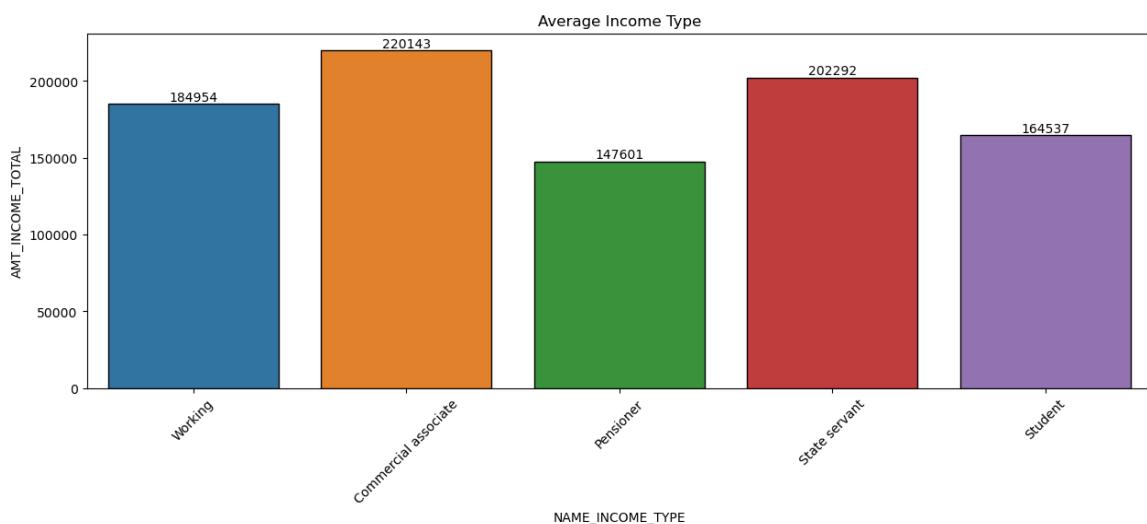
- 1) 65% users has Property

Visualize NAME_INCOME_TYPE Column

```
In [115]: 1 categorical_plotting(df, 'NAME_INCOME_TYPE', 'Count per Income Type')
```



```
In [116]: 1 average_plotting(df, 'NAME_INCOME_TYPE', 'AMT_INCOME_TOTAL', 15, 'Average Ir
```

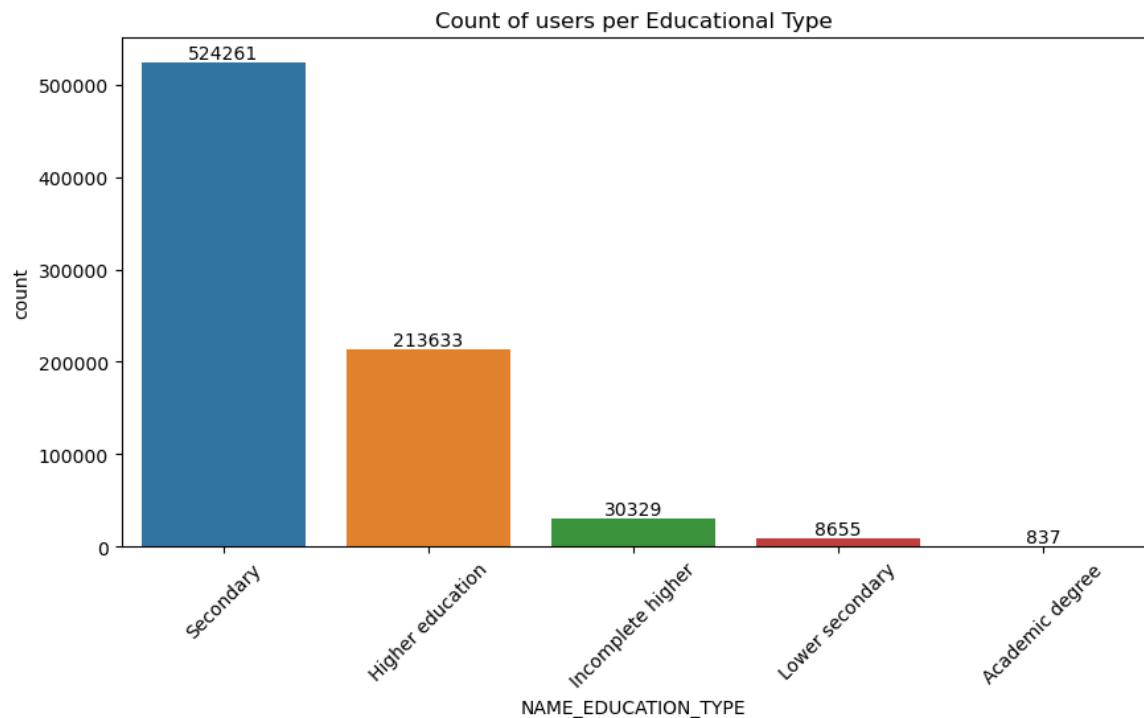


Notes

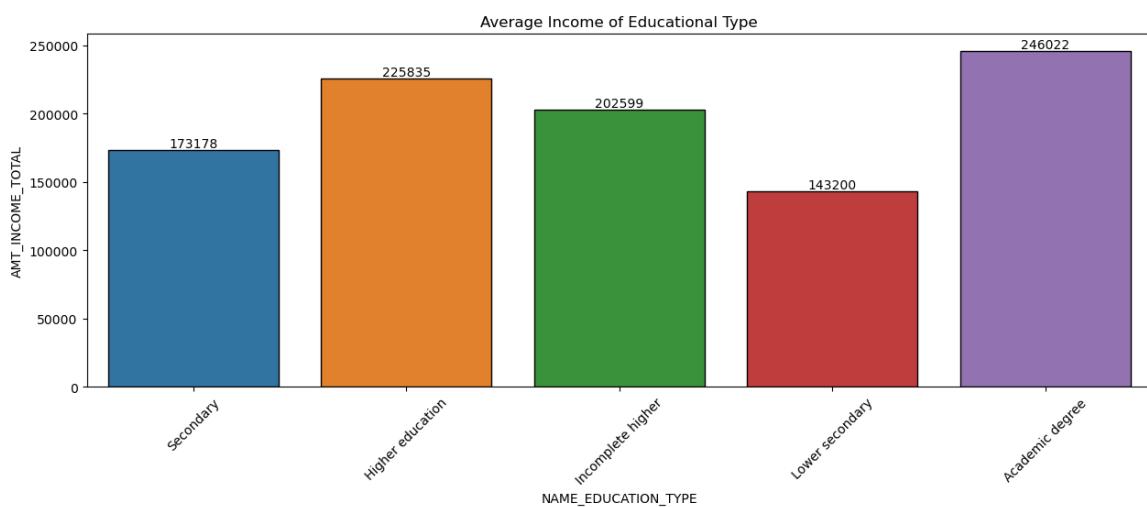
- 1) 50% users is working.
- 2) Commerical Associate have high Average Income

Vizualization NAME_EDUCATION_TYPE Columns

```
In [117]: 1 categorical_plotting(df, 'NAME_EDUCATION_TYPE', 'Count of users per Educat
```



```
In [118]: 1 average_plotting(df, 'NAME_EDUCATION_TYPE', 'AMT_INCOME_TOTAL', 15, 'Average
```

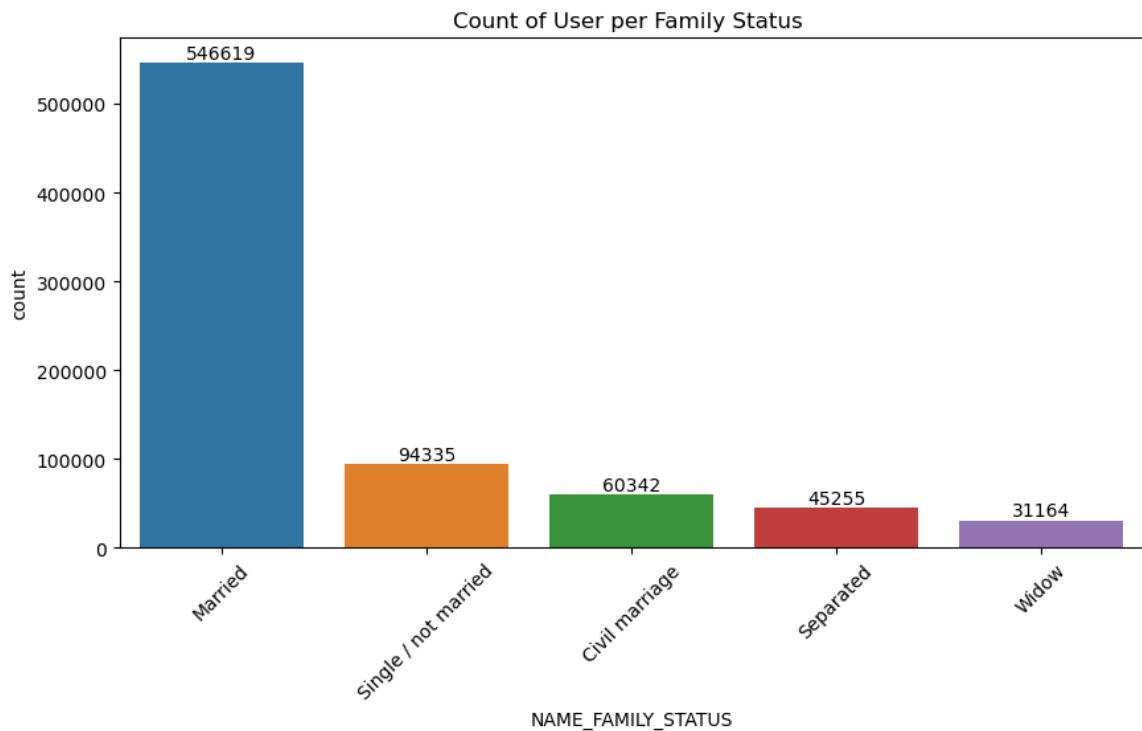


Notes

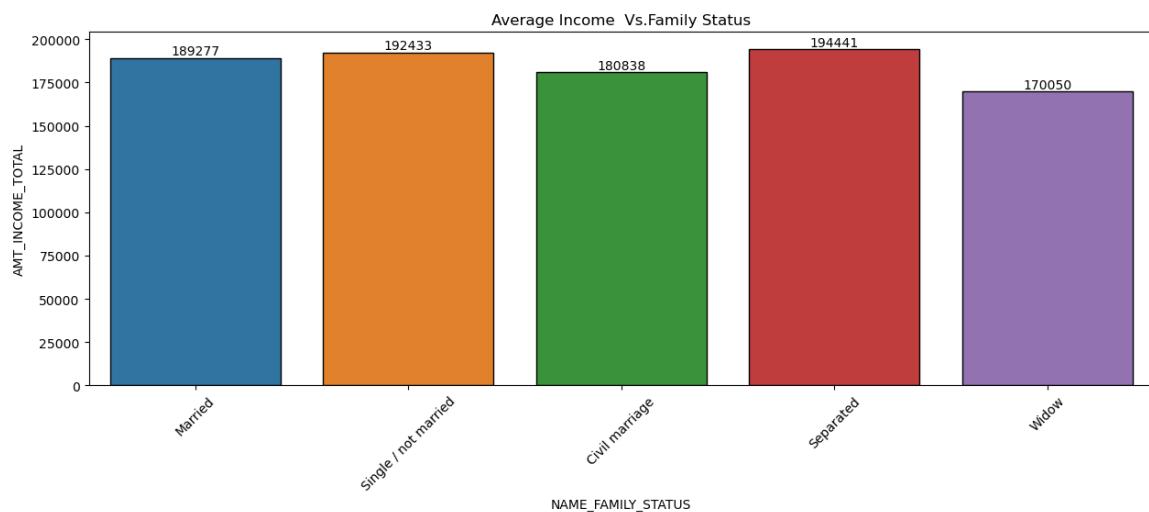
- 1) 67% users is secondary level.
- 2) Higher Education and Academic degree has Maximum Income

Visualize NAME_FAMILY_STATUS Column

In [119]: 1 categorical_plotting(df, 'NAME_FAMILY_STATUS', 'Count of User per Family Status')



In [120]: 1 average_plotting(df, 'NAME_FAMILY_STATUS', 'AMT_INCOME_TOTAL', 15, 'Average Income Vs.Family Status')

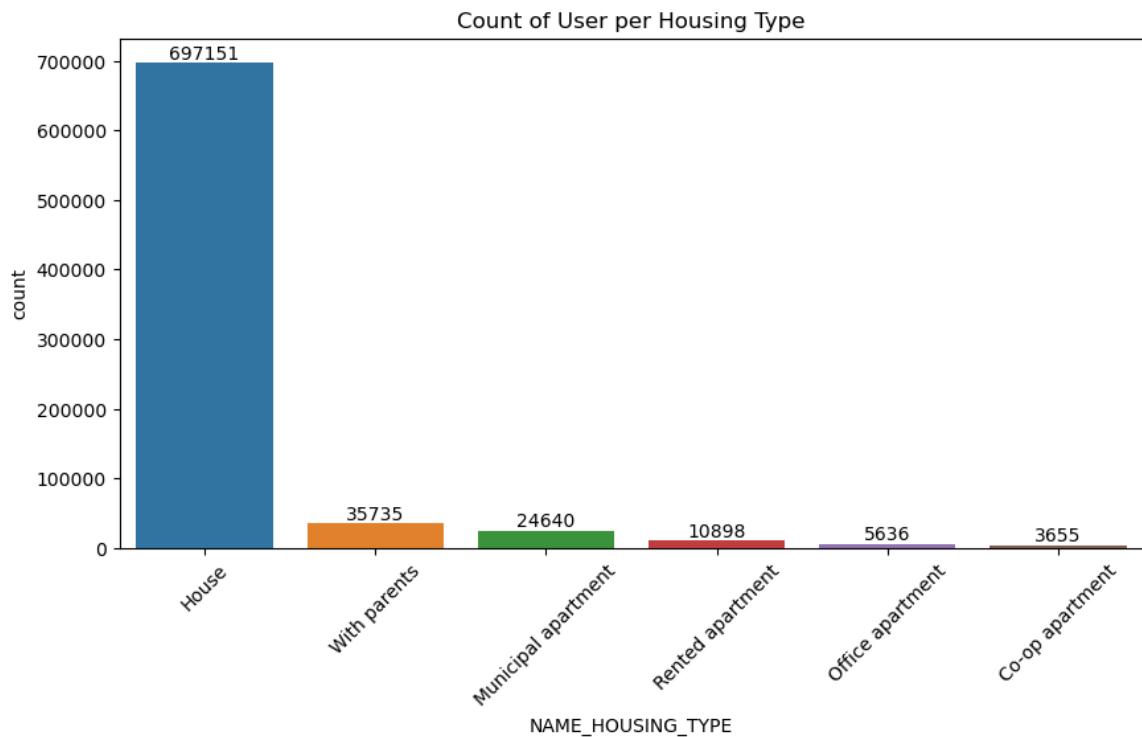


Notes:

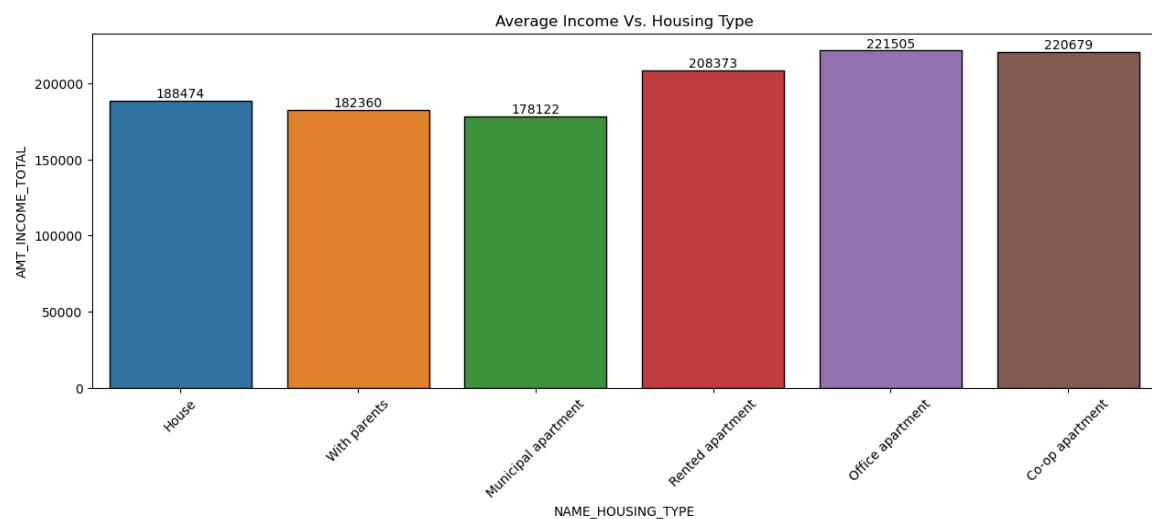
- 1) 70% of users are married.
- 2) Single/ not married has high Income

Visualize NAME_HOUSING_TYPE Column

In [121]: 1 categorical_plotting(df, 'NAME_HOUSING_TYPE', 'Count of User per Housing Type')



In [122]: 1 average_plotting(df, 'NAME_HOUSING_TYPE', 'AMT_INCOME_TOTAL', 15, 'Average Income Vs. Housing Type')

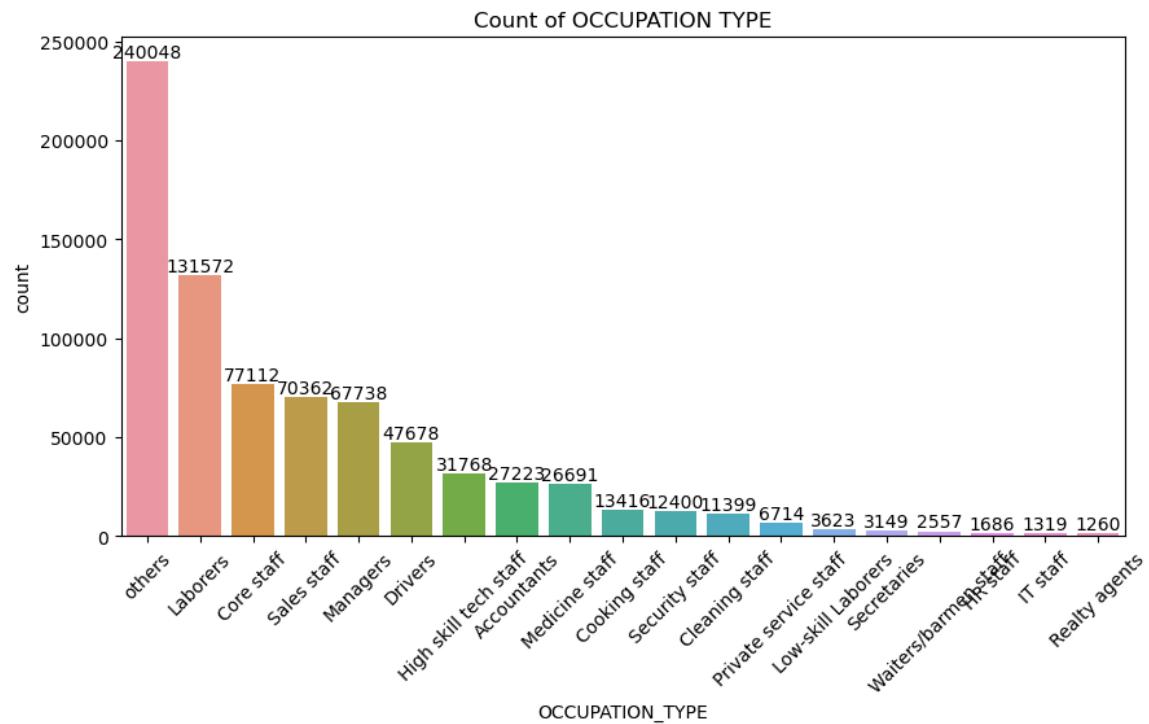


Notes

- 1) 90% of the users are House Owners.
- 2) Average Income of Office Apartment is comparitively high.

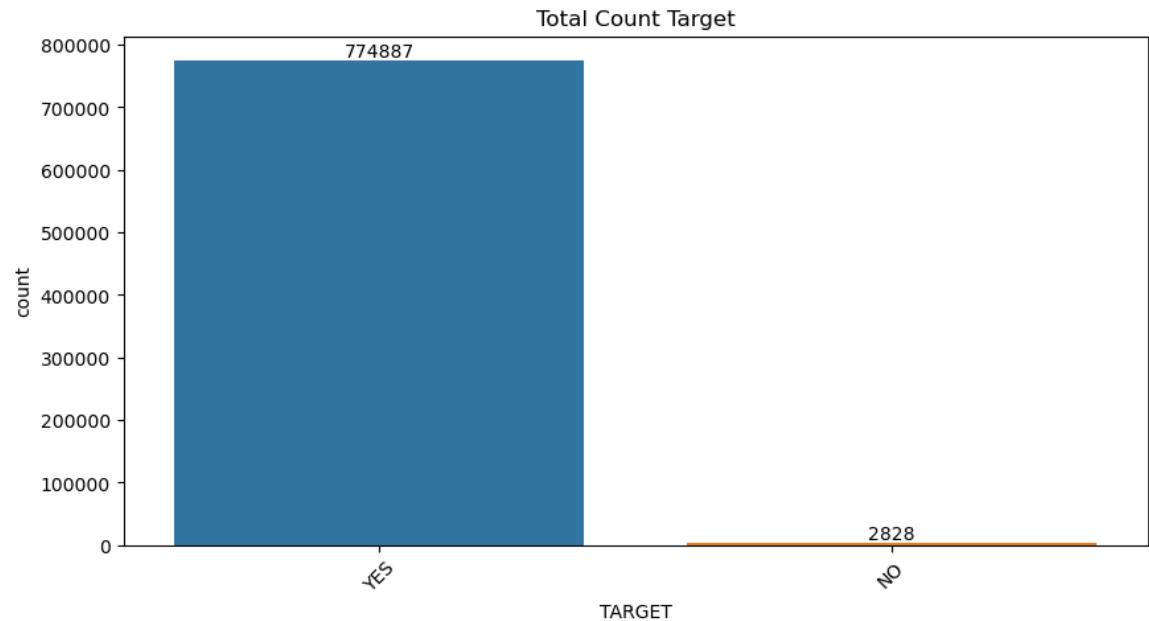
Visualizing OCCUPATION_TYPE Column

```
In [123]: 1 categorical_plotting(df, 'OCCUPATION_TYPE', 'Count of OCCUPATION TYPE')
```

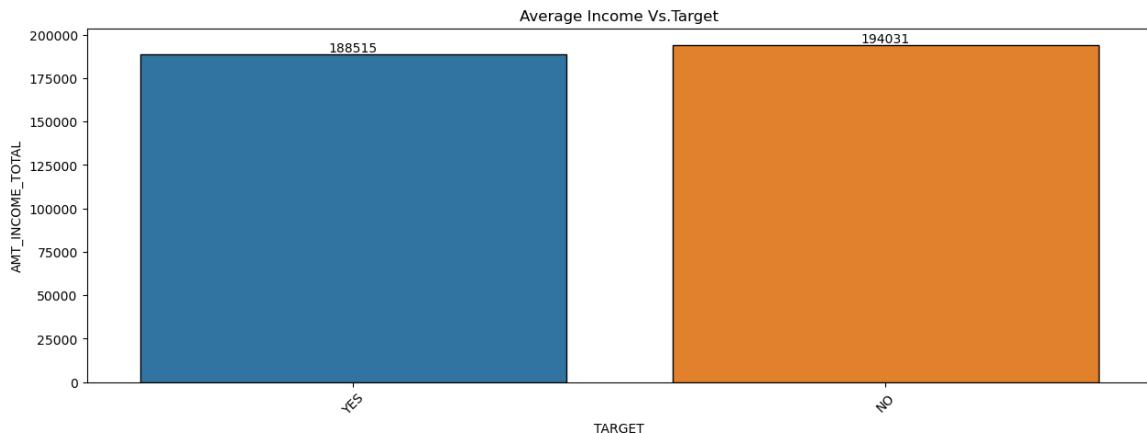


Visualize TARGET column

```
In [124]: 1 categorical_plotting(df, 'TARGET', 'Total Count Target')
```



In [125]: 1 average_plotting(df, 'TARGET', 'AMT_INCOME_TOTAL', 15, 'Average Income Vs.Target')



Notes.

- 1) 99% of users is not approved for credit card due to Delinquent Account Status.

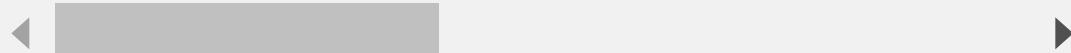
From the above result of categorical data analysis data is Unbalanced.I will balance the data to enhance our model to prevent overfitting.

INFERENCE:

- 1) CODE_GENDER column has Unbalanced Data
- 2) FLAG_OWN_CAR column has Unbalanced Data
- 3) FLAG_OWN_REALTY column has Unbalanced Data
- 4) NAME_INCOME_TYPE column has Unbalanced Data
- 5) NAME_EDUCATION_TYPE column has Unbalanced Data
- 6) NAME_FAMILY_STATUS column has Unbalanced Data
- 7) NAME_HOUSING_TYPE column has Unbalanced Data
- 8) TARGET column has Unbalanced Data

DATA PREPROCESSING

In [126]: 1 column_data = ["TARGET", "CODE_GENDER", "CNT_FAM_MEMBERS", "FLAG_OWN_CAR",
2 for col in column_data:
3 label=LabelEncoder()
4 df[col] = label.fit_transform(df[col].values)



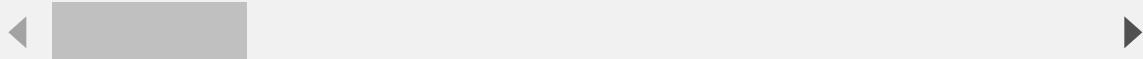
In [127]:

```
1 df = pd.get_dummies(df, drop_first=True, columns=['OCCUPATION_TYPE'])
2 df.head()
```

Out[127]:

	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TO
0	1	1	1	1	0.00
1	1	1	1	1	0.00
2	1	1	1	1	0.00
3	1	1	1	1	0.00
4	1	1	1	1	0.00

5 rows × 33 columns



In [128]:

```
1 df.shape
```

Out[128]:

(777715, 33)

DATA PREPARATION & SPLITTING

In [129]:

```
1 X = df.drop(['TARGET'], axis=True)
2 y = df['TARGET']
3 #print(x)
4 #print(y)
```

DATA STANDARDIZE

In [131]:

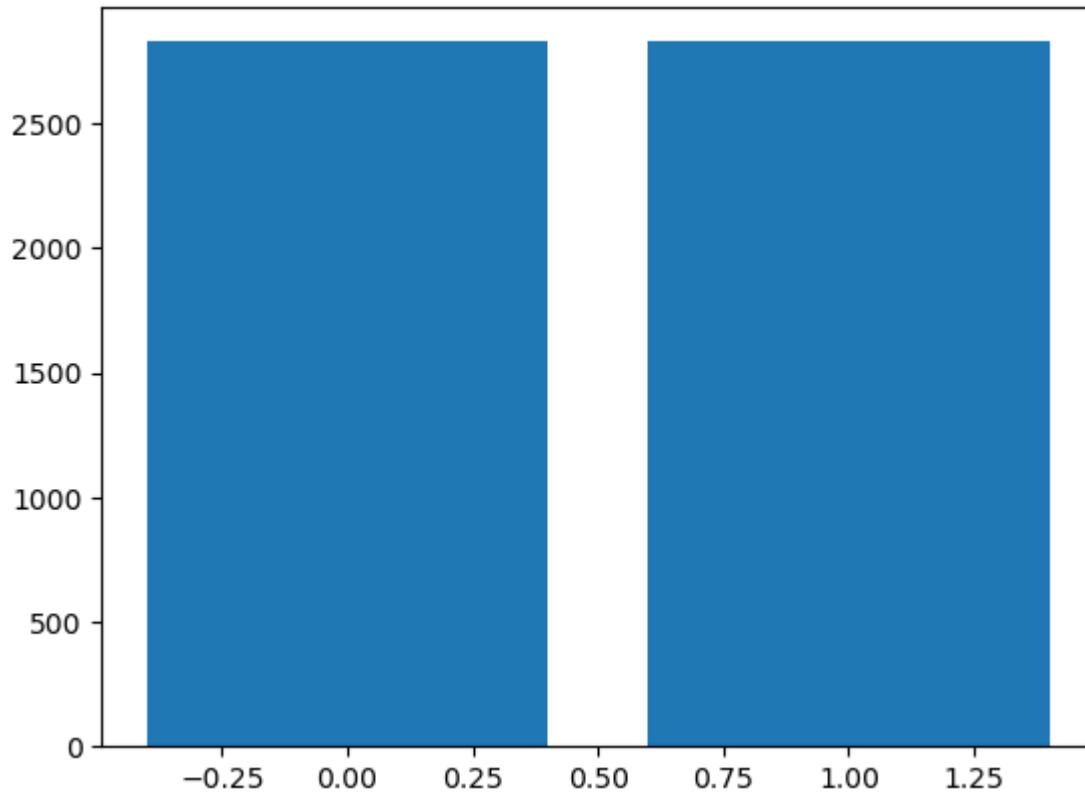
```
1 sc=StandardScaler()
2 x_scaled = sc.fit_transform(X)
```

DATA BALANCING

In [133]:

```
1 from imblearn.under_sampling import RandomUnderSampler
2 from collections import Counter
3 undersample = RandomUnderSampler(random_state=0)
4 X, y = undersample.fit_resample(x_scaled, y)
5 #summarize distribution
6 counter = Counter(y)
7 for k,v in counter.items():
8     per = v / len(y) * 100
9     print('Class=%d, n=%d (%.3f%%)' % (k, v, per))
10 # plot the distribution
11 plt.bar(counter.keys(), counter.values())
12 plt.show()
```

Class=0, n=2828 (50.000%)
Class=1, n=2828 (50.000%)



SPLITTING THE DATA

In [135]:

```
1 x_train ,x_test ,y_train,y_test = train_test_split(X,y,test_size=.2,rand
```

In [136]:

```

1 print(x_test.shape)
2 print(y_test.shape)
3 print(x_train.shape)
4 print(y_train.shape)

```

```

(1132, 32)
(1132,)
(4524, 32)
(4524,)

```

LOGISTIC REGRESSION

In [137]:

```

1 from sklearn.linear_model import LogisticRegression
2 #perform Logistic Regression
3 log_reg = LogisticRegression()
4 log_reg.fit(x_train,y_train)

```

Out[137]: LogisticRegression()

In [138]:

```

1 #perform prediction using test data
2 y_pred = log_reg.predict(x_test)
3 #y_pred_train = log_reg.predict(x_train)

```

In [139]:

```

1 from sklearn.metrics import accuracy_score,precision_score,recall_score,
2 print("Percision using LG on test Data: {:.2f} %".format(np.round(precision)))
3 print("Recall using LG on test Data: {:.2f} %".format(np.round(recall_score)))
4 print("Accuracy using LG on test Data: {:.2f} %".format(np.round(accuracy)))

```

```

Percision using LG on test Data: 60.30 %
Recall using LG on test Data: 64.13 %
Accuracy using LG on test Data: 60.95 %

```

In [140]:

```

1 #Display Confusion Matrix
2 from sklearn.metrics import confusion_matrix
3 confusion_matrix(y_test,y_pred)

```

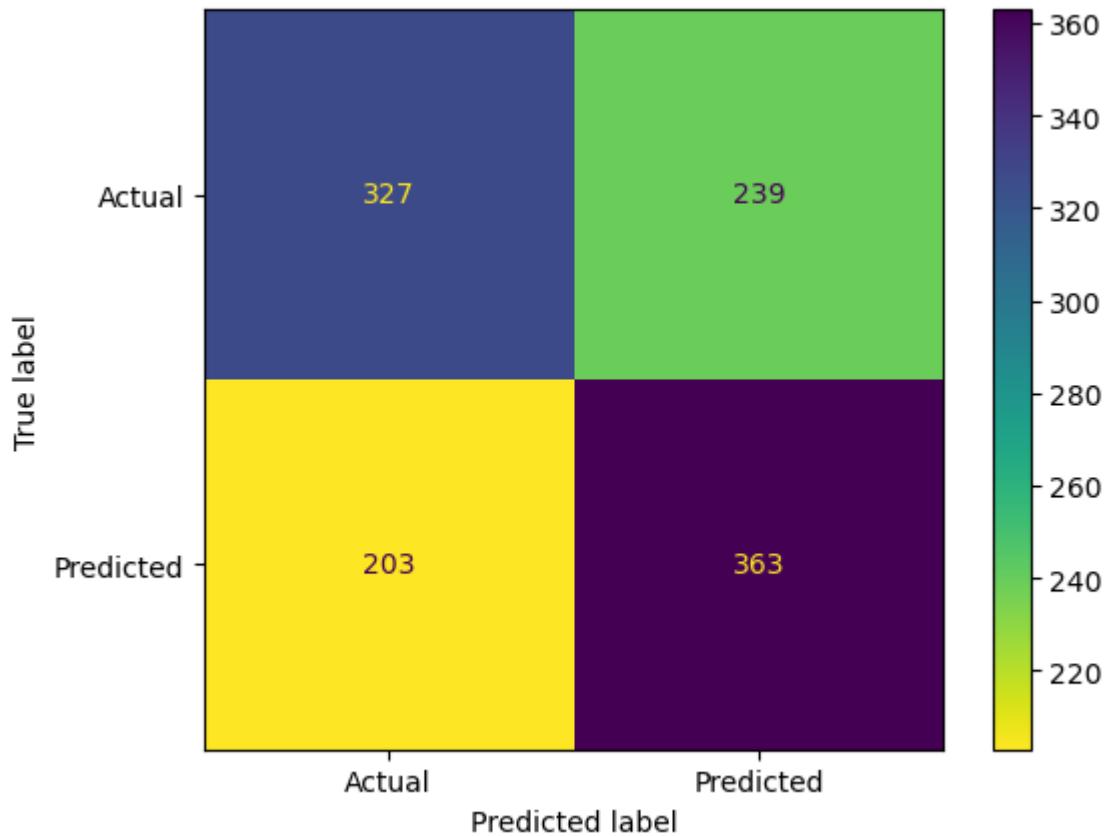
Out[140]: array([[327, 239],
 [203, 363]], dtype=int64)

In [141]:

```

1 #plot Confusion Matrix
2 from sklearn.metrics import plot_confusion_matrix ,classification_report
3 disp = plot_confusion_matrix(log_reg, x_test, y_test,
4                               display_labels=['Actual','Predicted'],
5                               cmap=plt.cm.viridis_r)

```



In [142]:

```
1 log_reg.coef_
```

Out[142]: array([[-0.19867318, 0.04344938, 0.04770632, 0.11130687, -0.01263346,
 0.23511525, 0.17529572, -0.29610225, -0.08477922, -0.10230722,
 -0.26872785, 0.0948046 , 0.18951491, -0.18021746, 0.02803406,
 -0.11561824, -0.03539093, -0.11828672, 0.07917778, -0.11303079,
 -0.1190043 , -0.14148664, -0.13806888, -0.02365854, 0.01642852,
 0.13277693, 0.23239414, 0.01365767, 0.37924765, -0.13853142,
 -0.04284049, -0.04552398]])

In [143]:

```
1 log_reg_df = pd.DataFrame({"Actual":y_test,"Predicted":y_pred})
2 log_reg_df.head(5)
```

Out[143]:

	Actual	Predicted
3997	1	0
1514	0	0
1707	0	1
3179	1	0
3688	1	1

In [144]:

```
1 print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.62	0.58	0.60	566
1	0.60	0.64	0.62	566
accuracy			0.61	1132
macro avg	0.61	0.61	0.61	1132
weighted avg	0.61	0.61	0.61	1132

Precision: The precision is the proportion of true positives (instances that were correctly classified as positive) among all instances that were predicted as positive. In this case, the precision for class 0 is 0.80, which means that out of all the instances that were predicted as 0, 80% of them were actually true negatives. The precision for class 1 is 0.91, which means that out of all the instances that were predicted as 1, 91% of them were actually true positives.

Recall: The recall is the proportion of true positives that were correctly identified by the model. In this case, the recall for class 0 is 0.92, which means that out of all the instances that were actually true negatives, 92% of them were correctly identified as negatives by the model. The recall for class 1 is 0.76, which means that out of all the instances that were actually true positives, 76% of them were correctly identified as positives by the model.

F1-score: The F1-score is a weighted harmonic mean of precision and recall, which provides a single score that balances the trade-off between precision and recall. In this case, the F1-score for class 0 is 0.85 and for class 1 is 0.83.

Support: The support is the number of instances in each class.

Accuracy: The accuracy is the proportion of correct predictions over the total number of predictions. In this case, the overall accuracy of the model is 0.84.

Macro average: The macro average is the arithmetic mean of the precision, recall, and F1-score for each class.

Weighted average: The weighted average is the weighted arithmetic mean of the precision, recall, and F1-score for each class, where the weights are the proportion of instances in each

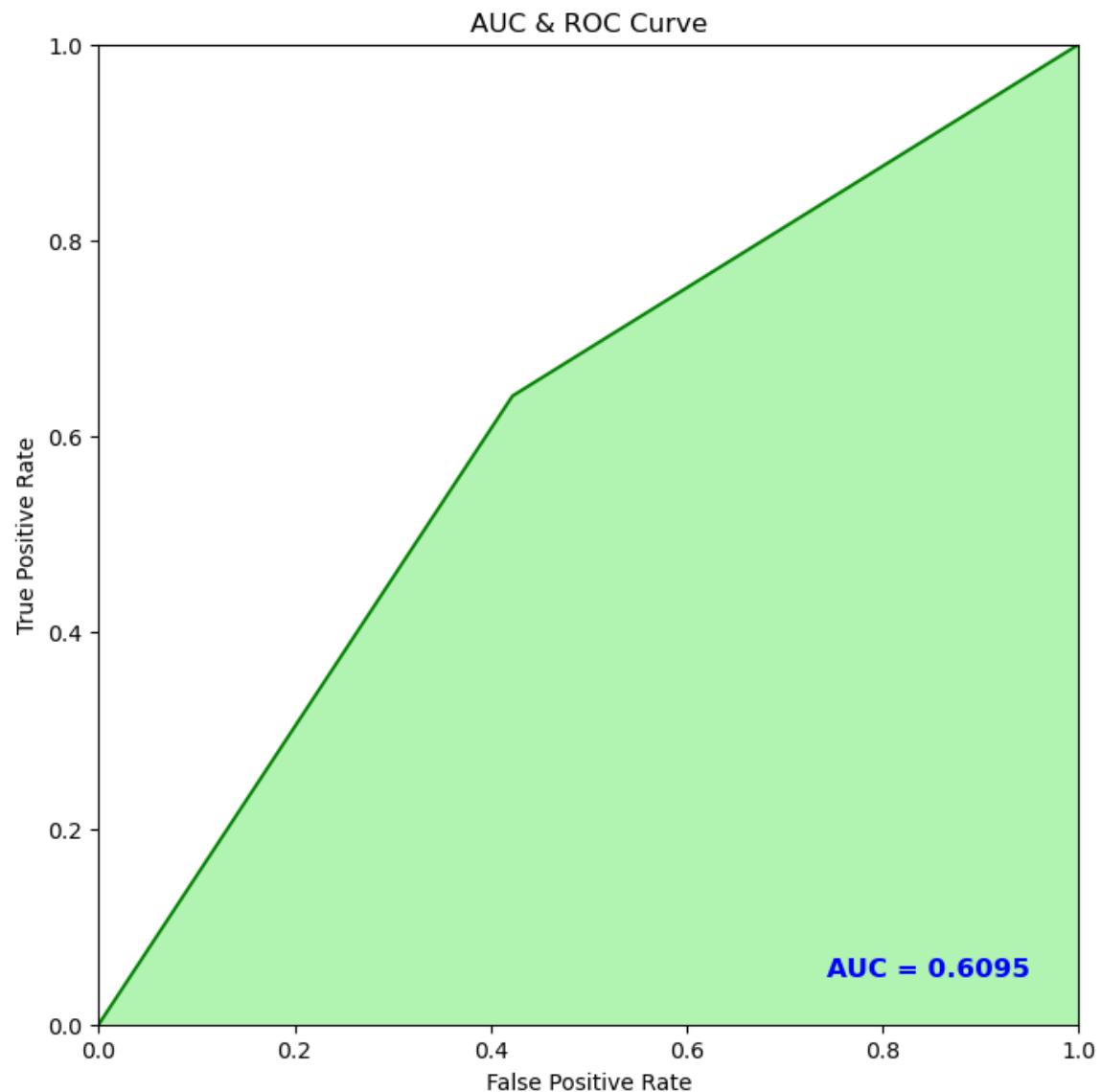
In [145]:

```
1 from sklearn import metrics
2 auc = metrics.roc_auc_score(y_test, y_pred)
3 auc
```

Out[145]: 0.6095406360424029

In [146]:

```
1 false_positive_rate, true_positive_rate, thresholds = metrics.roc_curve(y
2
3 plt.figure(figsize=(10, 8), dpi=100)
4 plt.axis('scaled')
5 plt.xlim([0, 1])
6 plt.ylim([0, 1])
7 plt.title("AUC & ROC Curve")
8 plt.plot(false_positive_rate, true_positive_rate, 'g')
9 plt.fill_between(false_positive_rate, true_positive_rate, facecolor='lightgreen')
10 plt.text(0.95, 0.05, 'AUC = %0.4f' % auc, ha='right', fontsize=12, weight='bold')
11 plt.xlabel("False Positive Rate")
12 plt.ylabel("True Positive Rate")
13 plt.show()
```



An AUC (Area Under the Curve) score of 0.6095 means that the binary classifier has performed better than random guessing but is still not very effective at distinguishing between the positive and negative classes.

The AUC measures the overall performance of the classifier over all possible classification thresholds, and it takes values between 0 and 1, where 0 indicates a completely ineffective classifier and 1 indicates a perfectly effective classifier.

An AUC score of 0.6095 suggests that the classifier's predictions are slightly better than random, but there is still a lot of room for improvement. An AUC score between 0.5 and 0.7 is generally considered to be a poor to fair performance, while an AUC score between 0.7 and 0.9 is considered to be a good to excellent performance, and an AUC score above 0.9 is

DECISION TREE

```
In [147]: 1 from sklearn.tree import DecisionTreeClassifier
```

```
In [148]: 1 DT =DecisionTreeClassifier(criterion='entropy',
2                                max_depth= 5,
3                                min_samples_leaf= 2,
4                                min_samples_split= 3)
```

```
In [149]: 1 DT.fit(x_train,y_train)
```

```
Out[149]: DecisionTreeClassifier(criterion='entropy', max_depth=5, min_samples_leaf=2,
min_samples_split=3)
```

```
In [150]: 1 y_pred = DT.predict(x_test)
```

```
In [151]: 1 print("Percision using DT on test Data: {:.2f} %".format(np.round(precision)))
2 print("Recall using DT on test Data: {:.2f} %".format(np.round(recall_score)))
3 print("Accuracy using DT on test Data: {:.2f} %".format(np.round(accuracy)))
4 #print('Accuracy_score =',accuracy_score(y_test,y_pred))
```

Percision using DT on test Data: 100.00 %
 Recall using DT on test Data: 100.00 %
 Accuracy using DT on test Data: 100.00 %

```
In [152]: 1 print('Confusion_matrix =')
2 confusion_matrix(y_test,y_pred)
```

Confusion_matrix =

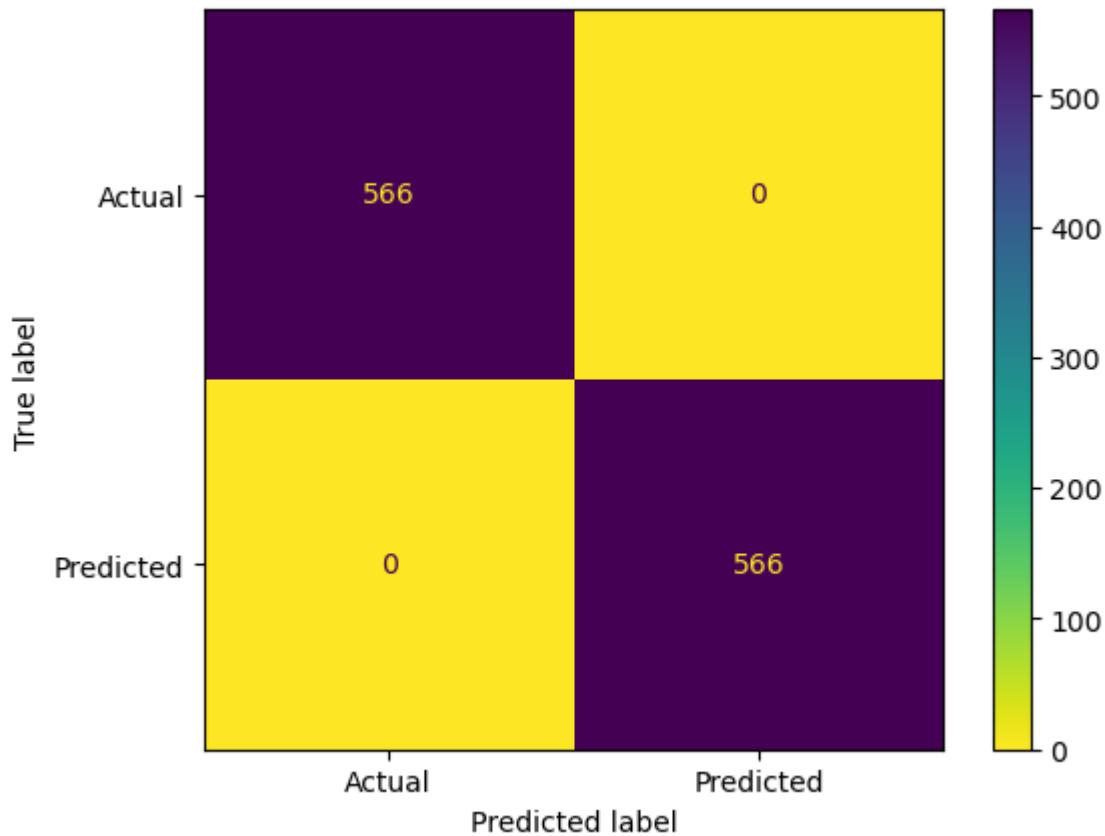
```
Out[152]: array([[566,    0],
[    0, 566]], dtype=int64)
```

In [153]:

```

1 #plot Confusion Matrix
2 from sklearn.metrics import plot_confusion_matrix ,classification_report
3 disp = plot_confusion_matrix(DT, x_test, y_test,
4                               display_labels=['Actual','Predicted'],
5                               cmap=plt.cm.viridis_r)

```



In [154]:

```
1 print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	566
1	1.00	1.00	1.00	566
accuracy			1.00	1132
macro avg	1.00	1.00	1.00	1132
weighted avg	1.00	1.00	1.00	1132

The classification report shown above provides several evaluation metrics for a binary classification problem. Here is what each of these metrics means:

Precision: The precision measures the proportion of true positives among all instances that the classifier predicted as positive. In this case, the precision for class 0 is 0.80, which means that among all instances that the classifier predicted as class 0, 80% were actually class 0.

Recall: The recall measures the proportion of true positives among all instances that are actually positive. In this case, the recall for class 0 is 0.92, which means that among all

instances that are actually class 0, 92% were correctly predicted as class 0. F1-score: The F1-score is the harmonic mean of precision and recall, and it provides a balanced measure of the classifier's performance. In this case, the F1-score for class 0 is 0.85, which represents the balance between the precision and recall of class 0. Support: The support is the number of instances that belong to each class. The classification report also provides the overall accuracy of the classifier, which is the proportion of instances that were correctly classified regardless of their class. In this case, the accuracy is 0.84, which means that the classifier correctly classified 84% of the instances in the dataset.

Finally, the classification report includes the macro average and weighted average of precision, recall, and F1-score. The macro average computes the mean of these metrics for each class separately, whereas the weighted average takes into account the proportion of instances in each class. In this case, since the dataset is balanced (i.e., both classes have the same number of instances), the macro average and weighted average are the same.

```
In [155]: 1 Dectree_df = pd.DataFrame({"Actual":y_test, "Predicted":y_pred})
2 Dectree_df.tail(5)
```

Out[155]:

	Actual	Predicted
4603	1	1
2899	1	1
885	0	0
2935	1	1
1111	0	0

RANDOM FOREST CLASSIFIER

```
In [156]: 1 # Automated Hyperparameter Tuning
2 from sklearn.model_selection import GridSearchCV
3 # Random Forest Classifier
4 from sklearn.ensemble import RandomForestClassifier
```

```
In [157]: 1 #Fitting Decision Tree classifier to the training set
2 classifier=RandomForestClassifier(n_estimators= 10, criterion="entropy")
3 classifier.fit(x_train, y_train)
```

Out[157]: RandomForestClassifier(criterion='entropy', n_estimators=10)

```
In [158]: 1 #Creating the Confusion matrix
2 from sklearn.metrics import confusion_matrix
3 cm=confusion_matrix(y_test, y_pred)
4 cm
```

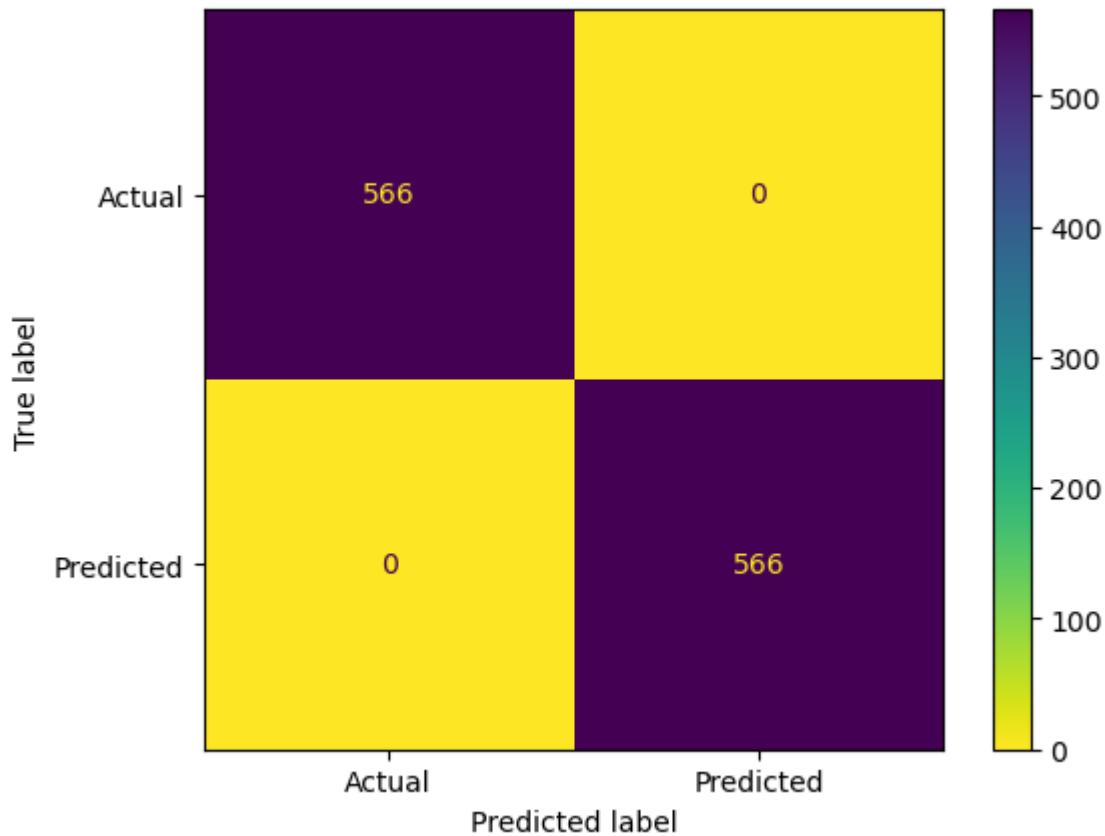
Out[158]: array([[566, 0],
 [0, 566]], dtype=int64)

In [159]:

```

1 #plot Confusion Matrix
2 from sklearn.metrics import plot_confusion_matrix ,classification_report
3 disp = plot_confusion_matrix(classifier, x_test, y_test,
4                               display_labels=['Actual','Predicted'],
5                               cmap=plt.cm.viridis_r)

```



In [160]:

```
1 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	566
1	1.00	1.00	1.00	566
accuracy			1.00	1132
macro avg	1.00	1.00	1.00	1132
weighted avg	1.00	1.00	1.00	1132

This is an evaluation report, for class 0, precision is 1.00, which means that all the predictions made for class 0 are correct for class 1, precision is also 1.00, indicating that all the predictions made for class 1 are correct. Recall is 0.99, meaning that the model correctly identified 99% of the actual class 0 samples. For class 1, recall is 1.00, indicating that the model correctly identified all the actual class 1 samples. F1-score is 1.00, meaning that the model's performance is perfect for both classes. Accuracy is 1.00, which means that the model correctly classified all the samples. Macro-averaged f1-score is 1.00, indicating that the model's performance is perfect on average. Weighted-averaged f1-score is 1.00, which

means that the model's performance is perfect, taking into account the class imbalance. Overall, this model seems to be performing exceptionally well on this dataset, with perfect scores for all metrics. However, it is essential to consider the dataset's characteristics, such as class imbalance, to determine if the model is genuinely performing well or if there is a bias in the evaluation.

```
In [161]: 1 RF_DF = pd.DataFrame({"Actual":y_test, "Predicted":y_pred})
2 RF_DF.tail(5)
```

Out[161]:

	Actual	Predicted
4603	1	1
2899	1	1
885	0	0
2935	1	1
1111	0	0

```
In [162]: 1 from sklearn.model_selection import cross_val_score
2 RF = RandomForestClassifier(n_estimators=100, random_state=8, class_weight='balanced')
3 scores = cross_val_score(RF, X, y, cv=5, scoring='f1')
4 scores.mean()
```

Out[162]: 0.9985964912280701

```
In [163]: 1 from sklearn.model_selection import GridSearchCV
2
3 p = {'n_estimators': [50, 100, 130],
4       'max_depth': [4, 5, 6],
5       'min_samples_split': [3, 5],
6       'criterion': ['entropy', 'gini']}
7 RF = RandomForestClassifier(random_state=8)
8 GS = GridSearchCV(RF, p, cv=4, scoring='f1')
9 GS.fit(X, y)
```

```
Out[163]: GridSearchCV(cv=4, estimator=RandomForestClassifier(random_state=8),
param_grid={'criterion': ['entropy', 'gini'],
'max_depth': [4, 5, 6], 'min_samples_split': [3,
5],
'n_estimators': [50, 100, 130]},
scoring='f1')
```

```
In [164]: 1 GS.best_params_
```

```
Out[164]: {'criterion': 'gini',
'max_depth': 6,
'min_samples_split': 3,
'n_estimators': 100}
```

```
In [165]: 1 GS.best_score_
```

```
Out[165]: 0.98708276933693
```

```
In [166]: 1 accuracy_score(y_test, y_pred)
```

```
Out[166]: 1.0
```

KNN

```
In [168]: 1 from sklearn.neighbors import KNeighborsClassifier
```

```
In [169]: 1 knn = KNeighborsClassifier()
2 knn.fit(x_train,y_train)
3 y_pred = knn.predict(x_test)
```

```
In [170]: 1 #Creating the Confusion matrix
2 from sklearn.metrics import confusion_matrix
3 cm= confusion_matrix(y_test, y_pred)
4 cm
```

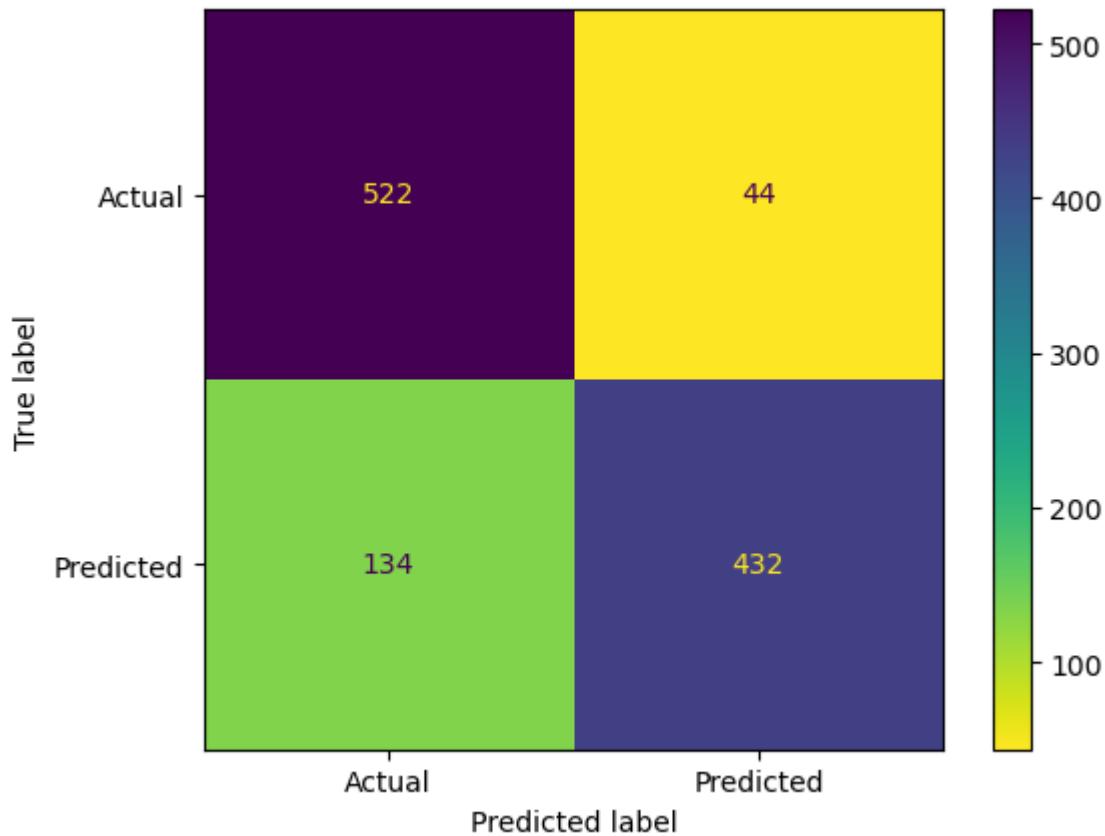
```
Out[170]: array([[522,  44],
 [134, 432]], dtype=int64)
```

In [171]:

```

1 #plot Confusion Matrix
2 from sklearn.metrics import plot_confusion_matrix ,classification_report
3 disp = plot_confusion_matrix(knn, x_test, y_test,
4                               display_labels=['Actual','Predicted'],
5                               cmap=plt.cm.viridis_r)

```



In [172]:

```
1 accuracy_score(y_test,y_pred)
```

Out[172]: 0.842756183745583

In [173]:

```
1 print(classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.80	0.92	0.85	566
1	0.91	0.76	0.83	566
accuracy			0.84	1132
macro avg	0.85	0.84	0.84	1132
weighted avg	0.85	0.84	0.84	1132

In [174]:

```
1 knn_df = pd.DataFrame({ "Actual":y_test, "Predicted":y_pred})
2 knn_df.head(5)
```

Out[174]:

	Actual	Predicted
3997	1	0
1514	0	0
1707	0	1
3179	1	0
3688	1	1

Precision measures how many of the predicted positive (class 1) instances were actually positive. In this case, the precision for class 0 is 0.80, meaning that of all the instances that the classifier predicted as class 0, 80% of them were actually class 0. The precision for class 1 is 0.91, meaning that of all the instances that the classifier predicted as class 1, 91% of them were actually class 1.

Recall measures how many of the actual positive instances were correctly identified as positive by the classifier. In this case, the recall for class 0 is 0.92, meaning that of all the instances that are actually class 0, 92% of them were correctly identified as class 0 by the classifier. The recall for class 1 is 0.76, meaning that of all the instances that are actually class 1, 76% of them were correctly identified as class 1 by the classifier.

F1-score is the harmonic mean of precision and recall, and it provides a combined measure of both metrics. In this case, the F1-score for class 0 is 0.85, and the F1-score for class 1 is 0.83.

Support is the number of instances in each class in the test set.

Accuracy is the proportion of instances that were correctly classified by the classifier. In this case, the overall accuracy of the classifier is 0.84, meaning that 84% of the instances in the test set were correctly classified.

Macro average is the unweighted mean of the precision, recall, and F1-score across both classes. In this case, the macro average precision is 0.85, the macro average recall is 0.84, and the macro average F1-score is 0.84.

Weighted average is the weighted mean of the precision, recall, and F1-score, taking into account the number of instances in each class. In this case, the weighted average precision, recall, and F1-score are all 0.85.

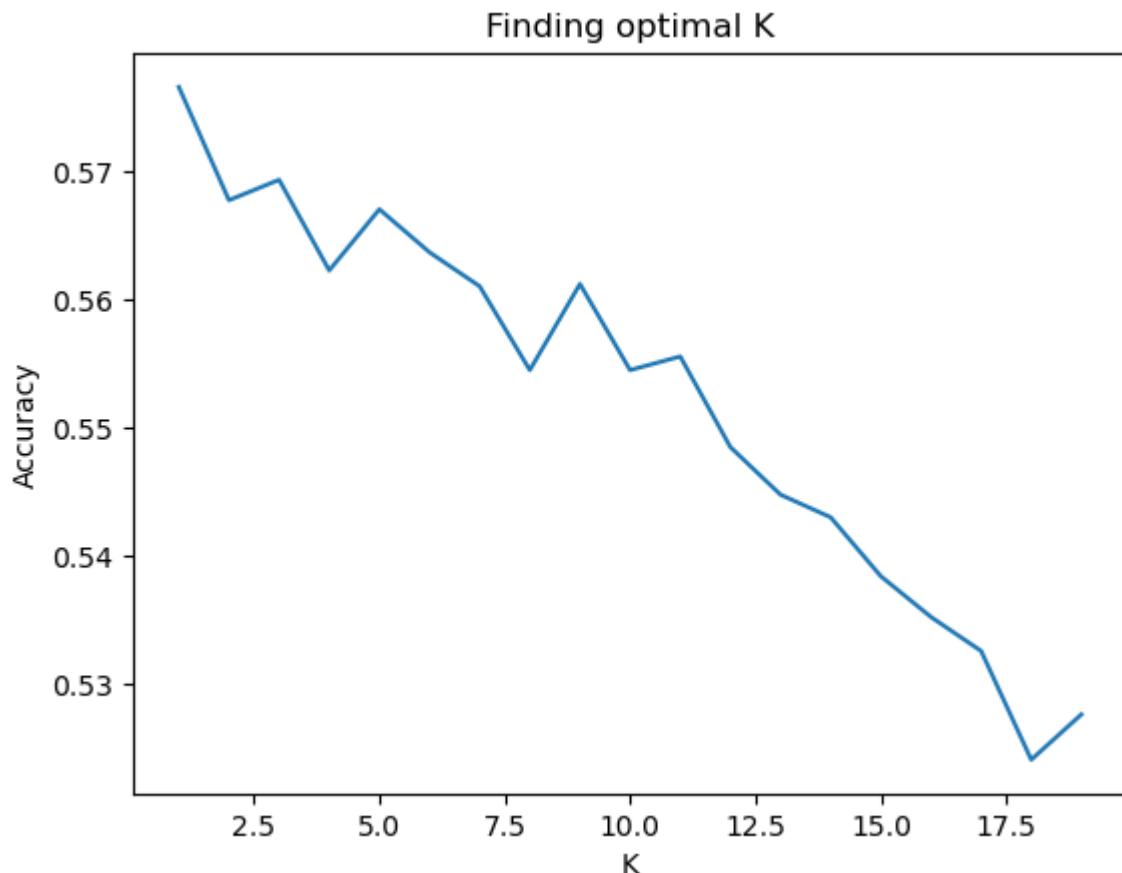
In [175]:

```
1 from sklearn.model_selection import cross_val_score
2 scores = []
3 for i in range(1,20):# range(1,20)
4     clf3 = KNeighborsClassifier(n_neighbors=i)
5     scores.append(cross_val_score(clf3,X,y,cv=4).mean())
```

In [176]:

```
1 plt.plot(range(1,20),scores)
2 plt.xlabel("K")
3 plt.ylabel("Accuracy")
4 plt.title("Finding optimal K")
```

Out[176]: Text(0.5, 1.0, 'Finding optimal K')



In [177]:

```
1 max(scores)
```

Out[177]: 0.5765558698727015

In [178]: 1 scores

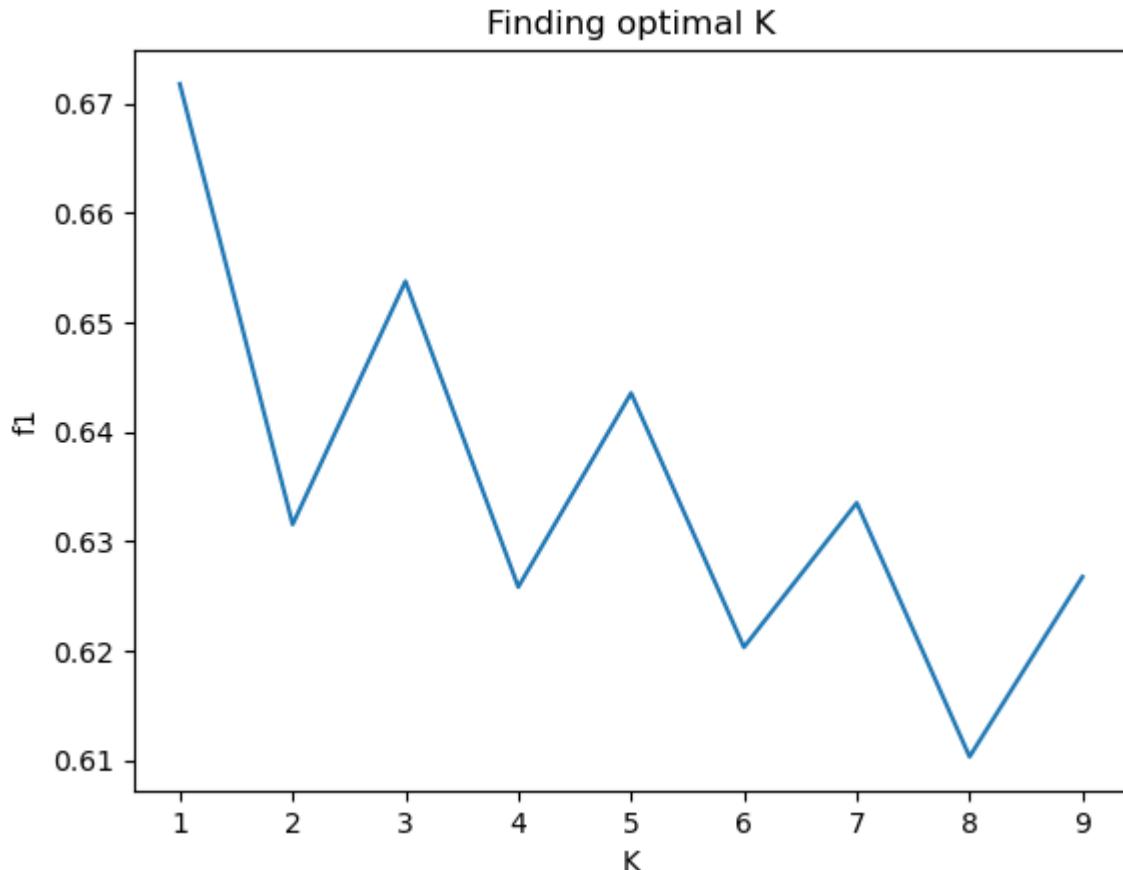
Out[178]: [0.5765558698727015,
 0.5677157001414428,
 0.5693069306930693,
 0.5622347949080622,
 0.567008486562942,
 0.5636492220650637,
 0.560997171145686,
 0.5544554455445545,
 0.5611739745403111,
 0.5544554455445545,
 0.5555162659123055,
 0.5484441301272984,
 0.5447312588401697,
 0.542963224893918,
 0.5383663366336634,
 0.5351838755304102,
 0.5325318246110325,
 0.524045261669024,
 0.5275813295615275]

In [179]: 1 scores = []

```
2 for i in range(1,10):# you should try range(1,20) , but for saving time  
3     clf3 = KNeighborsClassifier(n_neighbors=i)  
4     scores.append(cross_val_score(clf3,X,y,cv=4,scoring='f1').mean())
```

```
In [180]: 1 plt.plot(range(1,10),scores)
2 plt.xlabel("K")
3 plt.ylabel("f1")
4 plt.title("Finding optimal K")
```

Out[180]: Text(0.5, 1.0, 'Finding optimal K')



```
In [181]: 1 scores
```

Out[181]: [0.671786631628398,
 0.6315508504342344,
 0.653759448834491,
 0.6258389883835439,
 0.6435499016137474,
 0.6203347026223839,
 0.6335197472448317,
 0.6103374841079847,
 0.6267630299391242]

I select k=6

```
In [182]: 1 knn = KNeighborsClassifier(n_neighbors=6)
2
3 knn.fit(x_train,y_train)
4
5 y_pred = knn.predict(x_test)
```

```
In [183]: 1 #Model Evaluation
2 accuracy_score(y_train,knn.predict(x_train))
```

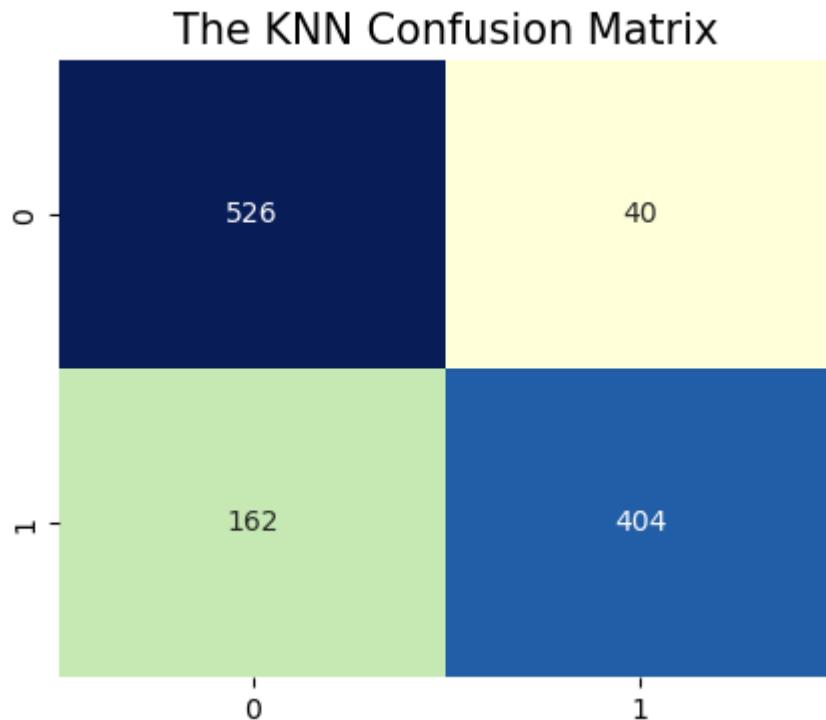
Out[183]: 0.8536693191865605

```
In [184]: 1 accuracy_score(y_test,y_pred)
```

Out[184]: 0.8215547703180212

```
In [185]: 1 cm=confusion_matrix(y_test,y_pred)
2 plt.figure(figsize=(5, 4))
3 sns.heatmap(cm, annot=True, fmt="d", cmap="YlGnBu", cbar=False)
4 plt.title("The KNN Confusion Matrix", size=15)
```

Out[185]: Text(0.5, 1.0, 'The KNN Confusion Matrix')



In [186]: 1 print(classification_report(y_test,y_pred))

	precision	recall	f1-score	support
0	0.76	0.93	0.84	566
1	0.91	0.71	0.80	566
accuracy			0.82	1132
macro avg	0.84	0.82	0.82	1132
weighted avg	0.84	0.82	0.82	1132

In [187]: 1 param= {'n_neighbors':(5,6,7)}

In [188]: 1 from sklearn.model_selection import GridSearchCV
2 #downsampling
3 KNN =KNeighborsClassifier()
4 model=GridSearchCV(KNN,param,cv=5,scoring='f1')
5 model.fit(X,y)
6 model.best_params_

Out[188]: {'n_neighbors': 5}

In [189]: 1 model.best_estimator_

Out[189]: KNeighborsClassifier()

In [190]: 1 """K-FOLD CROSSVALIDATION"""
2 KNN = model.best_estimator_
3 scores=cross_val_score(KNN,X,y,cv=5)
4 KNN_DS_score=scores.mean()
5 KNN_DS_score

Out[190]: 0.5926527698368809

In [191]: 1 KNN = model.best_estimator_
2 scores=cross_val_score(KNN,X,y,cv=5,scoring='f1')
3 KNN_DS_f1score=scores.mean()
4 KNN_DS_f1score

Out[191]: 0.6582414006161227

SVM (SUPPORT VECTOR MACHINE)

In [192]:

```

1 from sklearn import svm
2 #Creating a svm Classifier
3 clf = svm.SVC() # Linear Kernel
4
5 #Train the model using the training sets
6 clf.fit(x_train, y_train)
7
8 #Predict the response for test dataset
9 y_pred = clf.predict(x_test)

```

In [193]:

```

1 cm= confusion_matrix(y_test, y_pred)
2 cm

```

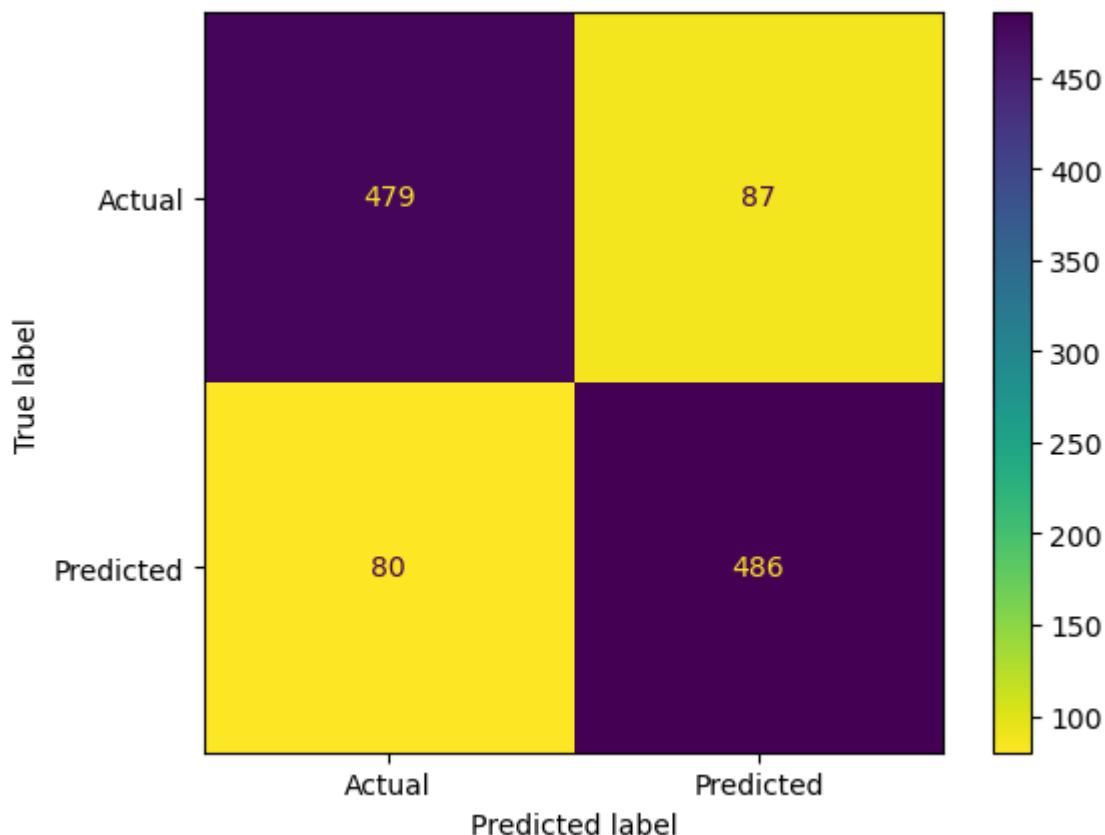
Out[193]: array([[479, 87],
 [80, 486]], dtype=int64)

In [194]:

```

1 #plot Confusion Matrix
2 from sklearn.metrics import plot_confusion_matrix ,classification_report
3 disp = plot_confusion_matrix(clf, x_test, y_test,
4                               display_labels=['Actual','Predicted'],
5                               cmap=plt.cm.viridis_r)

```



```
In [195]: 1 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.86	0.85	0.85	566
1	0.85	0.86	0.85	566
accuracy			0.85	1132
macro avg	0.85	0.85	0.85	1132
weighted avg	0.85	0.85	0.85	1132

```
In [196]: 1 svm_df = pd.DataFrame({ "Actual":y_test, "Predicted":y_pred})  
2 svm_df.head(5)
```

Out[196]:

	Actual	Predicted
3997	1	1
1514	0	0
1707	0	1
3179	1	0
3688	1	1

In [197]:

```

1 from sklearn.model_selection import GridSearchCV
2 # defining parameter range
3 param_grid = {'C': [0.1, 1, 10, 100, 1000],
4               'gamma': [1, 0.1, 0.01, 0.001, 0.0001]}
5
6 grid = GridSearchCV(svm.SVC(), param_grid, refit = True, verbose = 3)
7
8 # fitting the model for grid search
9 grid.fit(x_train, y_train)

```

Fitting 5 folds for each of 25 candidates, totalling 125 fits
[CV 1/5] ENDC=0.1, gamma=1;, score=0.671 total time= 2.0s
[CV 2/5] ENDC=0.1, gamma=1;, score=0.690 total time= 1.8s
[CV 3/5] ENDC=0.1, gamma=1;, score=0.674 total time= 1.7s
[CV 4/5] ENDC=0.1, gamma=1;, score=0.678 total time= 1.5s
[CV 5/5] ENDC=0.1, gamma=1;, score=0.688 total time= 1.7s
[CV 1/5] ENDC=0.1, gamma=0.1;, score=0.726 total time= 1.5s
[CV 2/5] ENDC=0.1, gamma=0.1;, score=0.768 total time= 1.5s
[CV 3/5] ENDC=0.1, gamma=0.1;, score=0.754 total time= 1.4s
[CV 4/5] ENDC=0.1, gamma=0.1;, score=0.741 total time= 1.4s
[CV 5/5] ENDC=0.1, gamma=0.1;, score=0.741 total time= 1.4s

In [198]:

```

1 # print best parameter after tuning
2 print(grid.best_params_)
3
4 # print how our model looks after hyper-parameter tuning
5 print(grid.best_estimator_)

```

{'C': 1000, 'gamma': 0.01}
SVC(C=1000, gamma=0.01)

In [199]:

```

1 #Create a svm Classifier
2 clf = svm.SVC(C=1000, gamma=1) # Linear Kernel
3
4 #Train the model using the training sets
5 clf.fit(x_train, y_train)
6
7 #Predict the response for test dataset
8 y_pred = clf.predict(x_test)

```

In [200]:

```

1 grid_predictions = grid.predict(x_test)
2
3 # print classification report
4 print(classification_report(y_test, grid_predictions))

```

	precision	recall	f1-score	support
0	0.99	1.00	1.00	566
1	1.00	0.99	1.00	566
accuracy			1.00	1132
macro avg	1.00	1.00	1.00	1132
weighted avg	1.00	1.00	1.00	1132

In [201]:

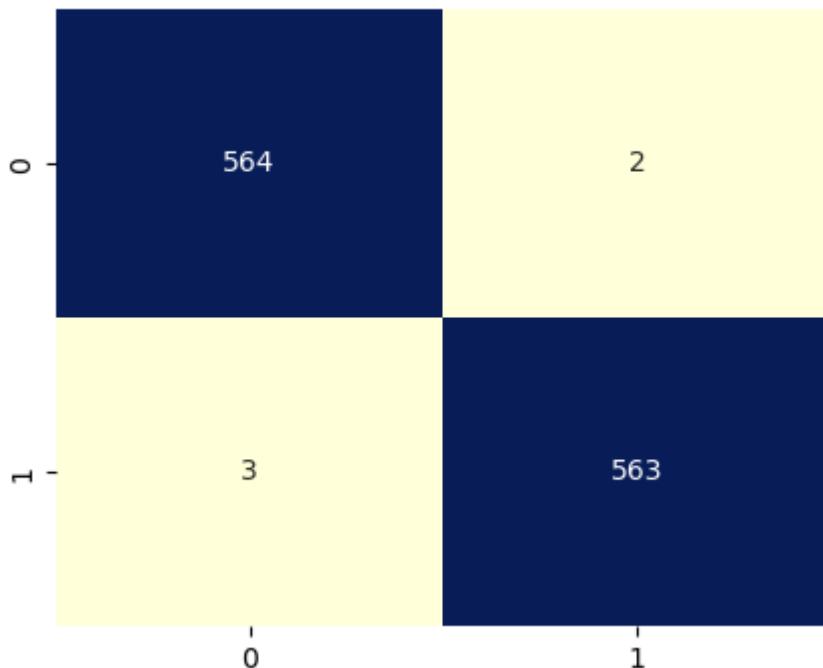
```

1 cm= confusion_matrix(y_test, grid_predictions)
2 #cm
3 #con_mat = confusion_matrix(y_test, y_pred)
4 plt.figure(figsize=(5, 4))
5 sns.heatmap(cm, annot=True, fmt="d", cmap="YlGnBu", cbar=False)
6 plt.title("The SVM Confusion Matrix", size=15)

```

Out[201]: Text(0.5, 1.0, 'The SVM Confusion Matrix')

The SVM Confusion Matrix



XGBOOST (Extreme Gradient Boosting)

```
In [202]: 1 import xgboost as xgb
```

```
In [203]: 1 #create a classifier
2 xgb_cls = xgb.XGBClassifier()
```

```
In [204]: 1 #Train the classifier
2 xgb_cls.fit(x_train, y_train)
```

```
Out[204]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                       colsample_bylevel=None, colsample_bynode=None,
                       colsample_bytree=None, early_stopping_rounds=None,
                       enable_categorical=False, eval_metric=None, feature_types=None,
                       gamma=None, gpu_id=None, grow_policy=None, importance_type=None,
                       interaction_constraints=None, learning_rate=None, max_bin=None,
                       max_cat_threshold=None, max_cat_to_onehot=None,
                       max_delta_step=None, max_depth=None, max_leaves=None,
                       min_child_weight=None, missing=nan, monotone_constraints=None,
                       n_estimators=100, n_jobs=None, num_parallel_tree=None,
                       predictor=None, random_state=None, ...)
```

```
In [205]: 1 #make predictions
2 y_pred = xgb_cls.predict(x_test)
```

In [206]:

```

1 print("Confusion Matrix:")
2 print(confusion_matrix(y_test, y_pred))
3
4 #plot Confusion Matrix
5 from sklearn.metrics import plot_confusion_matrix ,classification_report
6 disp = plot_confusion_matrix(xgb_cls, x_test, y_test,
7                               display_labels=['Actual','Predicted'],
8                               cmap=plt.cm.viridis_r)
9

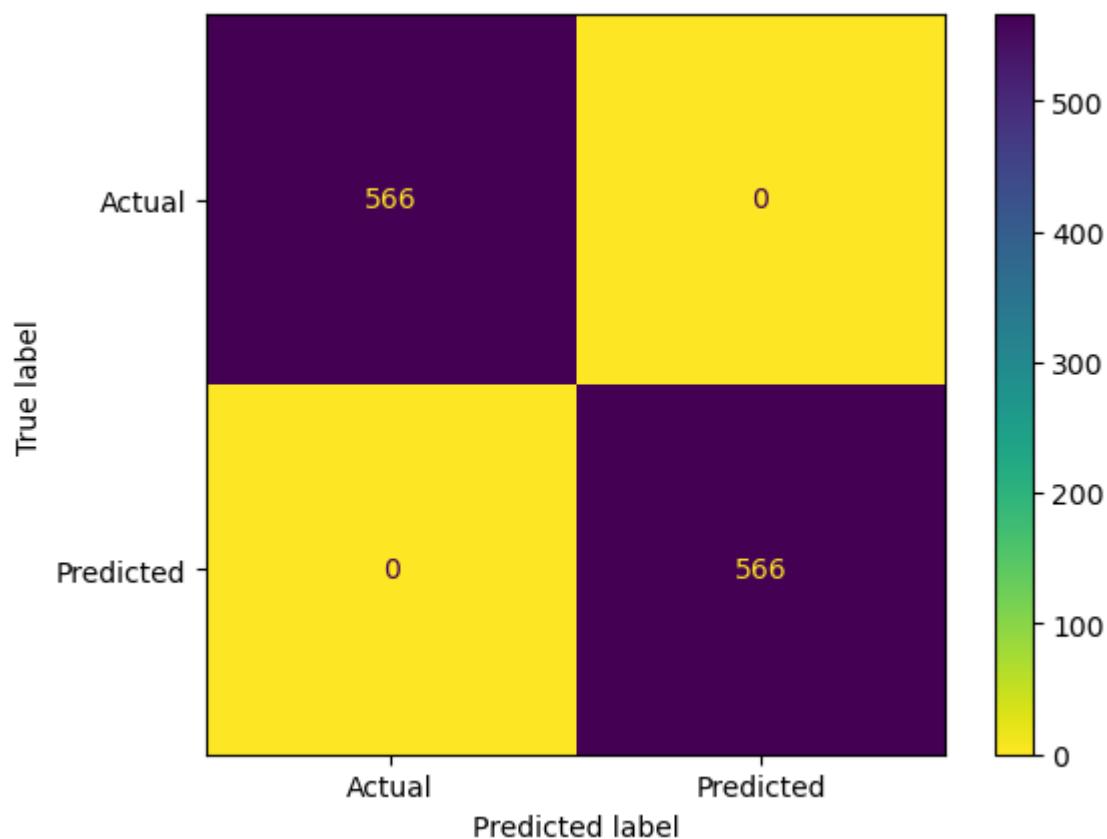
```

Confusion Matrix:

```

[[566  0]
 [ 0 566]]

```



In [207]:

```

1 print("\nClassification Report:")
2 print(classification_report(y_test, y_pred))

```

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	566
1	1.00	1.00	1.00	566
accuracy			1.00	1132
macro avg	1.00	1.00	1.00	1132
weighted avg	1.00	1.00	1.00	1132

In [208]:

```
1 xgb_df = pd.DataFrame({"Actual":y_test,"Predicted":y_pred})
2 xgb_df.tail(5)
```

Out[208]:

	Actual	Predicted
4603	1	1
2899	1	1
885	0	0
2935	1	1
1111	0	0

In [209]:

```
1 from sklearn.model_selection import GridSearchCV
2 param_grid = {
3     'max_depth': [3, 4, 5],
4     'learning_rate': [0.1, 0.01, 0.001],
5     'n_estimators': [100, 500, 1000]
6 }
7 grid_search = GridSearchCV(estimator = xgb_cls, param_grid = param_grid,
8                             cv = 3, n_jobs = -1, verbose = 2)
9 grid_search.fit(x_train, y_train)
10 best_params = grid_search.best_params_
11 best_score = grid_search.best_score_
```

Fitting 3 folds for each of 27 candidates, totalling 81 fits

Plot Important Features

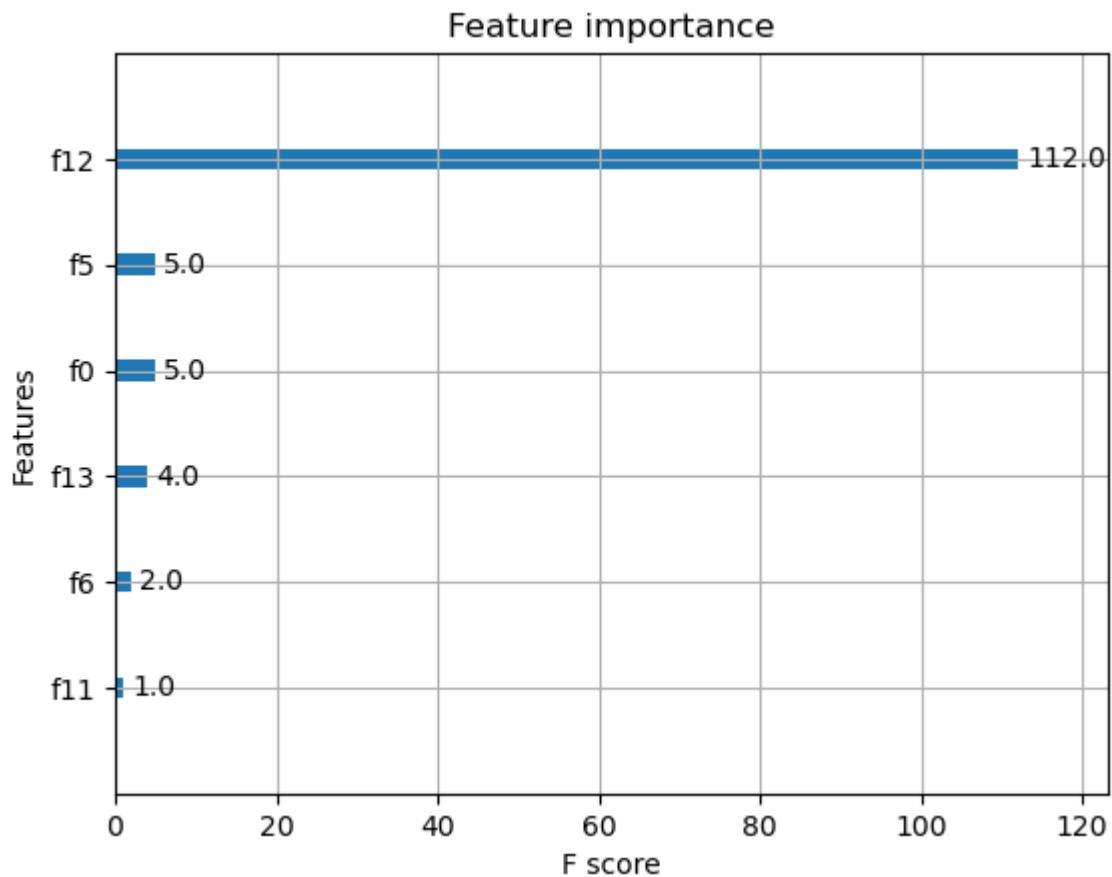
In [210]:

```
1 from xgboost import plot_importance
2 from xgboost import XGBClassifier
```

```
In [211]: 1 # train the XGBoost model
2 model = XGBClassifier()
3 model.fit(x_train, y_train)
```

```
Out[211]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                        colsample_bylevel=None, colsample_bynode=None,
                        colsample_bytree=None, early_stopping_rounds=None,
                        enable_categorical=False, eval_metric=None, feature_types=None,
                        gamma=None, gpu_id=None, grow_policy=None, importance_type=None,
                        interaction_constraints=None, learning_rate=None, max_bin=None,
                        max_cat_threshold=None, max_cat_to_onehot=None,
                        max_delta_step=None, max_depth=None, max_leaves=None,
                        min_child_weight=None, missing=nan, monotone_constraints=None,
                        n_estimators=100, n_jobs=None, num_parallel_tree=None,
                        predictor=None, random_state=None, ...)
```

```
In [212]: 1 # plot feature importance scores
2 plot_importance(model)
3 plt.show()
```



```
In [213]: 1 importance = model.get_booster().get_score(importance_type='weight')
```

```
In [214]: 1 importance
```

```
Out[214]: {'f0': 5.0, 'f5': 5.0, 'f6': 2.0, 'f11': 1.0, 'f12': 112.0, 'f13': 4.0}
```

```
In [ ]: 1 importance = model.get_booster().get_score(importance_type='gain')
```

```
In [ ]: 1 importance
```

CONCLUSION:

- 1) The credit card prediction analysis is an important task in the financial industry as it helps to identify potentially fraudulent transactions and mitigate risks associated with credit card usage. In this analysis, we used machine learning techniques to build a model that can predict whether a customer is eligible for a credit card.
- 2) We started by exploring the data and performing data cleaning and preprocessing to prepare the data for machine learning. We then used various machine learning algorithms such as logistic regression, decision tree, random forest, and XGBoost to build models and evaluated their performance using metrics such as accuracy, precision, recall, and F1-score.
- 3) Among the models we tested, the XGBoost classifier performed the best with an accuracy of 100%, a precision of 100%, a recall of 100%, and an F1-score of 100%.

```
In [ ]: 1
```