*Name: Nikhil Zodape*

*Email Id.: Nikszodape@gmail.com*

*Project Name: Bank Loan Case Study*

*Using: Python*

## Problem Statement:

As a data analyst at a finance company specializing in urban lending, the primary challenge is addressing default rates among customers with insufficient credit history. Some applications exploit this gaps, resulting in loan defaults. The objective is to conduct Exploratory Data Analysis (EDA) to uncover patterns in the data, ensuring that deserving application are not unfairly rejected.

When customers apply for loans, four potential outcomes exist: Approval, Cancellation during the approval process, Rejection and Approval with the loan remaining, unused. The overarching goal is to identify patterns that signal a customer's likelihood of struggling with instalment payments. This insight can inform decision such as loan denial, reducing loan amount, or offering loans with higher interest rates to higher-risk application. Ultimately, the company aims to discern the key factors influencing loan default, enabling more information decisions in the loan approval process.

## Task:

1. Identify the missing data in the dataset and decide on an appropriate method to deal with it.
2. Detect and identify outliers in the dataset.
3. Determine if there is data imbalance in the loan application dataset and calculate the ration of data imbalance.
4. Perform univariate analysis to understand the distribution of individual variables, segmented univariate analysis to compare variable distributions for difference scenarios and bivariate analysis to explore relationships between variables and the target variable.
5. Identify the top correlations for each segments.

## Data Understanding:

The dataset contain 3 files:
1. Application_data.csv: It contains all the information of the client at the time of application. The data is about the client is having any payment difficulty.
2. Previous_application.csv: It contain information about clients previous loan data. It contains the data whether the previous application had been Approved, Cancelled, Refused or Unused offer.
3. Columns_description.csv: It is the data dictionary which describe the meaning of the variables.
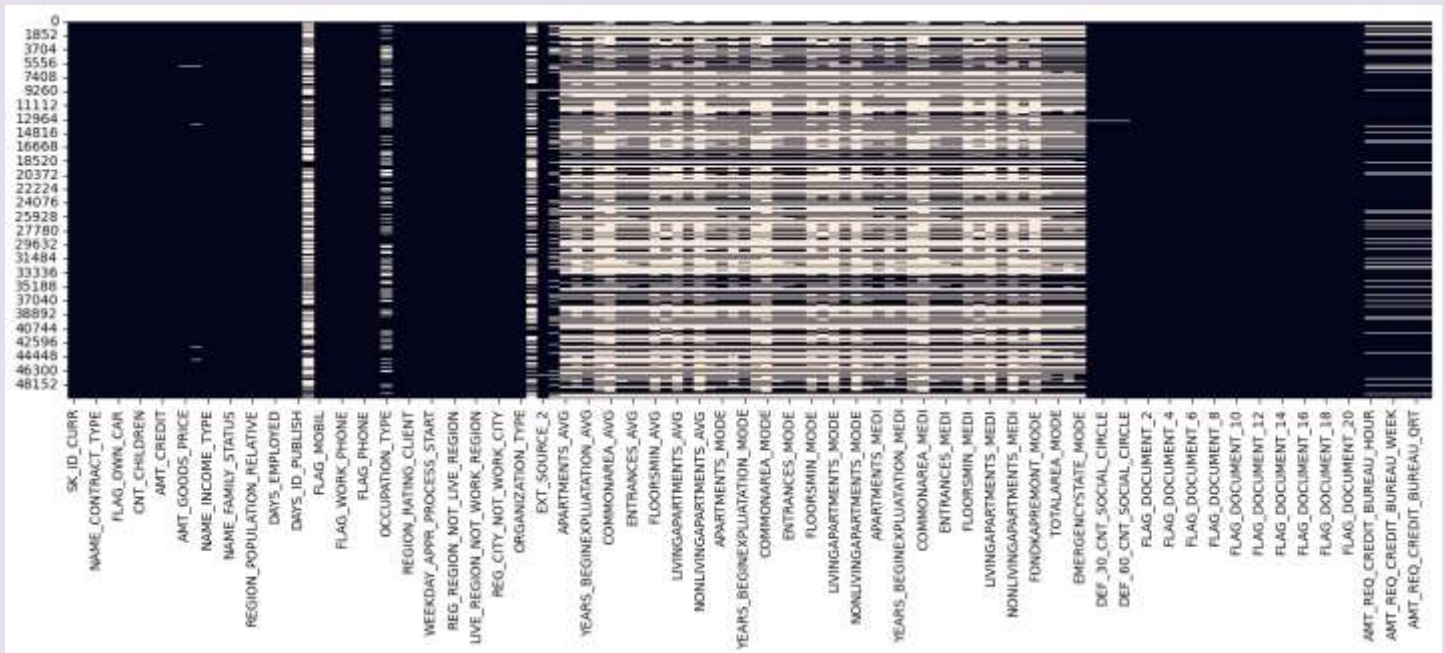
-→ Application data:

The application dataset contains 49999 rows and 122 columns.

Checking the Null values in application dataset.

| | column_name | null_count | null_percentage |
|---|---|---|---|
| 76 | COMMONAREA_MEDI | 34960 | 69.92 |
| 48 | COMMONAREA_AVG | 34960 | 69.92 |
| 62 | COMMONAREA_MODE | 34960 | 69.92 |
| 70 | NONLIVINGAPARTMENTS_MODE | 34714 | 69.43 |
| 56 | NONLIVINGAPARTMENTS_AVG | 34714 | 69.43 |
| 84 | NONLIVINGAPARTMENTS_MEDI | 34714 | 69.43 |
| 68 | LIVINGAPARTMENTS_MODE | 34226 | 68.45 |
| 54 | LIVINGAPARTMENTS_AVG | 34226 | 68.45 |
| 82 | LIVINGAPARTMENTS_MEDI | 34226 | 68.45 |
| 86 | FONDKAPREMONT_MODE | 34191 | 68.38 |
| 52 | FLOORSMIN_AVG | 33894 | 67.79 |
| 66 | FLOORSMIN_MODE | 33894 | 67.79 |
| 80 | FLOORSMIN_MEDI | 33894 | 67.79 |
| 75 | YEARS_BUILD_MEDI | 33239 | 66.48 |
| 61 | YEARS_BUILD_MODE | 33239 | 66.48 |
| 47 | YEARS_BUILD_AVG | 33239 | 66.48 |
| 21 | OWN_CAR_AGE | 32950 | 65.90 |
| 81 | LANDAREA_MEDI | 29721 | 59.44 |
| 67 | LANDAREA_MODE | 29721 | 59.44 |
| 53 | LANDAREA_AVG | 29721 | 59.44 |
| 73 | BASEMENTAREA_MEDI | 29199 | 58.40 |
| 45 | BASEMENTAREA_AVG | 29199 | 58.40 |
| 59 | BASEMENTAREA_MODE | 29199 | 58.40 |
| 41 | EXT_SOURCE_1 | 28172 | 56.35 |
| 71 | NONLIVINGAREA_MODE | 27572 | 55.15 |
| 57 | NONLIVINGAREA_AVG | 27572 | 55.15 |
| 85 | NONLIVINGAREA_MEDI | 27572 | 55.15 |
| 77 | ELEVATORS_MEDI | 26651 | 53.30 |
| 49 | ELEVATORS_AVG | 26651 | 53.30 |
| 63 | ELEVATORS_MODE | 26651 | 53.30 |
| 89 | WALLSMATERIAL_MODE | 25459 | 50.92 |
| 72 | APARTMENTS_MEDI | 25385 | 50.77 |
| 44 | APARTMENTS_AVG | 25385 | 50.77 |
| 58 | APARTMENTS_MODE | 25385 | 50.77 |
| 78 | ENTRANCES_MEDI | 25195 | 50.39 |
| 50 | ENTRANCES_AVG | 25195 | 50.39 |
| 64 | ENTRANCES_MODE | 25195 | 50.39 |

| | | | |
|---|---|---|---|
| 55 | LIVINGAREA_AVG | 25137 | 50.28 |
| 69 | LIVINGAREA_MODE | 25137 | 50.28 |
| 83 | LIVINGAREA_MEDI | 25137 | 50.28 |
| 87 | HOUSETYPE_MODE | 25075 | 50.15 |
| 65 | FLOORSMAX_MODE | 24875 | 49.75 |
| 79 | FLOORSMAX_MEDI | 24875 | 49.75 |
| 51 | FLOORSMAX_AVG | 24875 | 49.75 |
| 60 | YEARS_BEGINEXPLUATATION_MODE | 24394 | 48.79 |
| 74 | YEARS_BEGINEXPLUATATION_MEDI | 24394 | 48.79 |
| 46 | YEARS_BEGINEXPLUATATION_AVG | 24394 | 48.79 |
| 88 | TOTALAREA_MODE | 24148 | 48.30 |
| 90 | EMERGENCYSTATE_MODE | 23698 | 47.40 |
| 28 | OCCUPATION_TYPE | 15654 | 31.31 |
| 43 | EXT_SOURCE_3 | 9944 | 19.89 |
| 116 | AMT_REQ_CREDIT_BUREAU_HOUR | 6734 | 13.47 |
| 117 | AMT_REQ_CREDIT_BUREAU_DAY | 6734 | 13.47 |
| 118 | AMT_REQ_CREDIT_BUREAU_WEEK | 6734 | 13.47 |
| 119 | AMT_REQ_CREDIT_BUREAU_MON | 6734 | 13.47 |
| 120 | AMT_REQ_CREDIT_BUREAU_QRT | 6734 | 13.47 |
| 121 | AMT_REQ_CREDIT_BUREAU_YEAR | 6734 | 13.47 |
| 11 | NAME_TYPE_SUITE | 192 | 0.38 |
| 92 | DEF_30_CNT_SOCIAL_CIRCLE | 168 | 0.34 |
| 91 | OBS_30_CNT_SOCIAL_CIRCLE | 168 | 0.34 |
| 93 | OBS_60_CNT_SOCIAL_CIRCLE | 168 | 0.34 |
| 94 | DEF_60_CNT_SOCIAL_CIRCLE | 168 | 0.34 |
| 42 | EXT_SOURCE_2 | 126 | 0.25 |
| 10 | AMT_GOODS_PRICE | 38 | 0.08 |
| 6 | CNT_CHILDREN | 0 | 0.00 |
| 102 | FLAG_DOCUMENT_8 | 0 | 0.00 |
| 2 | NAME_CONTRACT_TYPE | 0 | 0.00 |
| 3 | CODE_GENDER | 0 | 0.00 |
| 4 | FLAG_OWN_CAR | 0 | 0.00 |
| 95 | DAYS_LAST_PHONE_CHANGE | 1 | 0.00 |
| 96 | FLAG_DOCUMENT_2 | 0 | 0.00 |
| 97 | FLAG_DOCUMENT_3 | 0 | 0.00 |
| 98 | FLAG_DOCUMENT_4 | 0 | 0.00 |
| 99 | FLAG_DOCUMENT_5 | 0 | 0.00 |
| 100 | FLAG_DOCUMENT_6 | 0 | 0.00 |

-→ Application data:

Plotting the heatmap to check the visualization of NaN values in columns.



Removing the NaN values columns whose columns are having NaN values Greater than 40% in it.

```
    Click here to ask Blackbox to help you code faster
1   # displaying all the collumn who are having more than 40% NaN values in it.
2   null_value_gre40 = null_percentage[null_percentage > 40].index
3   null_value_gre40
✓ 0.0s
```

```
Index(['OWN_CAR_AGE', 'EXT_SOURCE_1', 'APARTMENTS_AVG', 'BASEMENTAREA_AVG',
       'YEARS_BEGINEXPLUATATION_AVG', 'YEARS_BUILD_AVG', 'COMMONAREA_AVG',
       'ELEVATORS_AVG', 'ENTRANCES_AVG', 'FLOORSMAX_AVG', 'FLOORSMIN_AVG',
       'LANDAREA_AVG', 'LIVINGAPARTMENTS_AVG', 'LIVINGAREA_AVG',
       'NONLIVINGAPARTMENTS_AVG', 'NONLIVINGAREA_AVG', 'APARTMENTS_MODE',
       'BASEMENTAREA_MODE', 'YEARS_BEGINEXPLUATATION_MODE', 'YEARS_BUILD_MODE',
       'COMMONAREA_MODE', 'ELEVATORS_MODE', 'ENTRANCES_MODE', 'FLOORSMAX_MODE',
       'FLOORSMIN_MODE', 'LANDAREA_MODE', 'LIVINGAPARTMENTS_MODE',
       'LIVINGAREA_MODE', 'NONLIVINGAPARTMENTS_MODE', 'NONLIVINGAREA_MODE',
       'APARTMENTS_MEDI', 'BASEMENTAREA_MEDI', 'YEARS_BEGINEXPLUATATION_MEDI',
       'YEARS_BUILD_MEDI', 'COMMONAREA_MEDI', 'ELEVATORS_MEDI',
       'ENTRANCES_MEDI', 'FLOORSMAX_MEDI', 'FLOORSMIN_MEDI', 'LANDAREA_MEDI',
       'LIVINGAPARTMENTS_MEDI', 'LIVINGAREA_MEDI', 'NONLIVINGAPARTMENTS_MEDI',
       'NONLIVINGAREA_MEDI', 'FONDKAPREMONT_MODE', 'HOUSETYPE_MODE',
       'TOTALAREA_MODE', 'WALLSMATERIAL_MODE', 'EMERGENCYSTATE_MODE'],
      dtype='object')
```

Also dropping all the Flag Document Columns from the dataset.

```
    Click here to ask Blackbox to help you code faster
1   # Function to see all the documents columns in one list and to remove them.
2
3   def col_name(data):
4       columns_to_drop = []
5       for column in data.columns:
6           if "DOCUMENT" in column:
7               columns_to_drop.append(column)
8
9       return columns_to_drop
10
11
12  # Applying the function to the dataset
13  document_col = col_name(appli_data)
  ✓ 0.0s
```

```
['FLAG_DOCUMENT_2',
 'FLAG_DOCUMENT_3',
 'FLAG_DOCUMENT_4',
 'FLAG_DOCUMENT_5',
 'FLAG_DOCUMENT_6',
 'FLAG_DOCUMENT_7',
 'FLAG_DOCUMENT_8',
 'FLAG_DOCUMENT_9',
 'FLAG_DOCUMENT_10',
 'FLAG_DOCUMENT_11',
 'FLAG_DOCUMENT_12',
 'FLAG_DOCUMENT_13',
 'FLAG_DOCUMENT_14',
 'FLAG_DOCUMENT_15',
 'FLAG_DOCUMENT_16',
 'FLAG_DOCUMENT_17',
 'FLAG_DOCUMENT_18',
 'FLAG_DOCUMENT_19',
 'FLAG_DOCUMENT_20',
 'FLAG_DOCUMENT_21']
```

Again Checking all the FlAG Column from the dataset:

```python
# Click here to ask Blackbox to help you code faster
def flag_name(data):
    # This function takes a DataFrame 'data' as input and identifies columns containing the substring "FLAG".
    # It then returns a list of such columns.
    columns_to_drop = []
    for column in data.columns:
        if "FLAG" in column:
            columns_to_drop.append(column)

    return columns_to_drop

# Call the function with a DataFrame named 'appli_data' to get a list of columns containing 'FLAG'
fla = flag_name(appli_data)
print(fla)
print("\n")

# Create a dictionary to store unique values for each identified 'FLAG' column
unique_value = {}
for column in fla:
    # For each 'FLAG' column, store its unique values in the dictionary
    unique_value[column] = appli_data[column].unique()

# Print the dictionary containing unique values for 'FLAG' columns
unique_value
```
✓ 0.0s

```
['FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'FLAG_MOBIL', 'FLAG_EMP_PHONE', 'FLAG_WORK_PHONE', 'FLAG_CONT_MOBILE', 'FLAG_PHONE', 'FLAG_EMAIL']

{'FLAG_OWN_CAR': array(['N', 'Y'], dtype=object),
 'FLAG_OWN_REALTY': array(['Y', 'N'], dtype=object),
 'FLAG_MOBIL': array([1, 0], dtype=int64),
 'FLAG_EMP_PHONE': array([1, 0], dtype=int64),
 'FLAG_WORK_PHONE': array([0, 1], dtype=int64),
 'FLAG_CONT_MOBILE': array([1, 0], dtype=int64),
 'FLAG_PHONE': array([1, 0], dtype=int64),
 'FLAG_EMAIL': array([0, 1], dtype=int64)}
```

Now we dropped all the FLAG Document and Remaining FLAG Columns from the dataset Except FLAG_OWN_CAR and FLAG_OWN_REALTY.

Now we will check all the Unique values of all the object columns to count the frequency of each columns and fill the Nan values as well as XNA and XAP values.

```python
# Click here to ask Blackbox to help you code faster
# Get the columns with data type 'object'
object_column = appli_data.select_dtypes(include = ['object']).columns

# Create a dictionary to store unique values for each object column
unique_values_dict = {}

# Iterate through each object column and store its unique values.
for column in object_column:
    unique_values_dict[column] = appli_data[column].unique()

for column, values in unique_values_dict.items():
    print(f"Column: {column}")
    print(f"unique_value: {values}")
    print()
```
✓ 0.0s

```
Column: NAME_CONTRACT_TYPE
unique_value: ['Cash loans' 'Revolving loans']

Column: CODE_GENDER
unique_value: ['M' 'F' 'XNA']

Column: FLAG_OWN_CAR
unique_value: ['N' 'Y']

Column: FLAG_OWN_REALTY
unique_value: ['Y' 'N']

Column: NAME_TYPE_SUITE
unique_value: ['Unaccompanied' 'Family' 'Spouse, partner' 'Children' 'Other_A' nan
 'Other_B' 'Group of people']

Column: NAME_INCOME_TYPE
unique_value: ['Working' 'State servant' 'Commercial associate' 'Pensioner' 'Unemployed'
 'Student' 'Businessman' 'Maternity leave']

Column: NAME_EDUCATION_TYPE
unique_value: ['Secondary / secondary special' 'Higher education' 'Incomplete higher'
 'Lower secondary' 'Academic degree']

Column: NAME_FAMILY_STATUS
unique_value: ['Single / not married' 'Married' 'Civil marriage' 'Widow' 'Separated'
 'Unknown']

Column: NAME_HOUSING_TYPE
unique_value: ['House / apartment' 'Rented apartment' 'With parents'
 'Municipal apartment' 'Office apartment' 'Co-op apartment']
```

```
Column: OCCUPATION_TYPE
unique_value: ['Laborers' 'Core staff' 'Accountants' 'Managers' nan 'Drivers'
 'Sales staff' 'Cleaning staff' 'Cooking staff' 'Private service staff'
 'Medicine staff' 'Security staff' 'High skill tech staff'
 'Waiters/barmen staff' 'Low-skill Laborers' 'Realty agents' 'Secretaries'
 'IT staff' 'HR staff']

Column: WEEKDAY_APPR_PROCESS_START
unique_value: ['WEDNESDAY' 'MONDAY' 'THURSDAY' 'SUNDAY' 'SATURDAY' 'FRIDAY' 'TUESDAY']

Column: ORGANIZATION_TYPE
unique_value: ['Business Entity Type 3' 'School' 'Government' 'Religion' 'Other' 'XNA'
 'Electricity' 'Medicine' 'Business Entity Type 2' 'Self-employed'
 'Transport: type 2' 'Construction' 'Housing' 'Kindergarten'
 'Trade: type 7' 'Industry: type 11' 'Military' 'Services'
 'Security Ministries' 'Transport: type 4' 'Industry: type 1' 'Emergency'
 'Security' 'Trade: type 2' 'University' 'Transport: type 3' 'Police'
 'Business Entity Type 1' 'Postal' 'Industry: type 4' 'Agriculture'
 'Restaurant' 'Culture' 'Hotel' 'Industry: type 7' 'Trade: type 3'
 'Industry: type 3' 'Bank' 'Industry: type 9' 'Insurance' 'Trade: type 6'
 'Industry: type 2' 'Transport: type 1' 'Industry: type 12' 'Mobile'
 'Trade: type 1' 'Industry: type 5' 'Industry: type 10' 'Legal Services'
 'Advertising' 'Trade: type 5' 'Cleaning' 'Industry: type 13'
 'Trade: type 4' 'Telecom' 'Industry: type 8' 'Realtor' 'Industry: type 6']
```

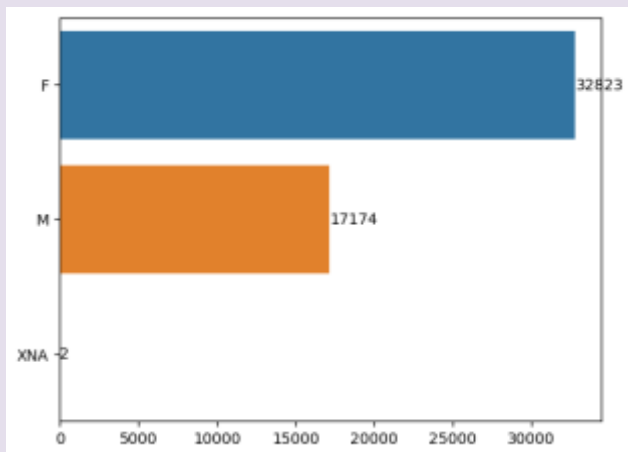Replacing all the XNA and XAP values of all the columns.'
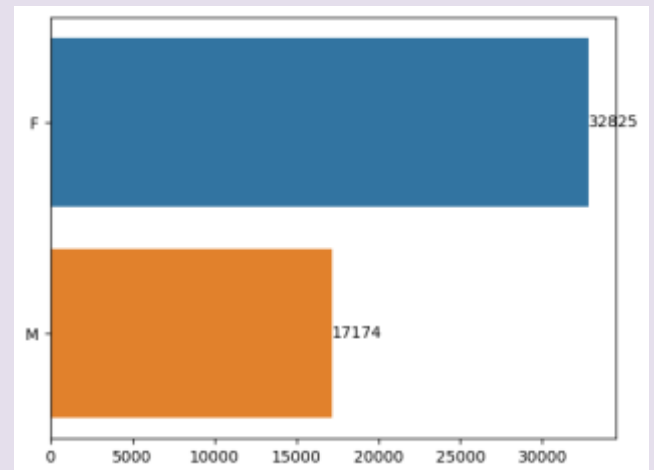
1. CODE_GENDER

```python
# Calculate the count of each unique value in the 'CODE_GENDER' column and store it in 'gender_count'
gender_count = appli_data.CODE_GENDER.value_counts()
print("Before Replacing XNA values: ")
fig = sns.barplot(y = gender_count.index, x = gender_count.values, estimator = 'sum', errorbar = None)
fig.bar_label(fig.containers[0], fontsize=10)
plt.show()


# Replacing the 'XNA' values with the most frequent gender (top most count) in the 'CODE_GENDER' column
# The assumption here is that the most frequent gender is a reasonable replacement for 'XNA'
appli_data['CODE_GENDER'] = appli_data['CODE_GENDER'].replace("XNA", "F")
gen = appli_data.CODE_GENDER.value_counts()

print('\nAfter Replacing the values: ')
fig = sns.barplot(y = gen.index, x = gen.values, estimator = 'sum', errorbar = None)
fig.bar_label(fig.containers[0], fontsize=10)
plt.show()
```

Before Replacing XNA values:                    After Replacing the values:



2. ORGANIZATION_TYPE

```python
# Checking the Orginization type counts
org_count = appli_data.ORGANIZATION_TYPE.value_counts()
print(org_count)

# Count of 'XNA' values
xna_count = org_count['XNA']

# Values to fill evenly for 'XNA'
fill_values = [
    'Business Entity Type 3',
    'Self-employed',
    'Other',
    'Medicine',
    'Government',
    'Business Entity Type 2',
    'School',
    'Trade: type 7',
    'Kindergarten'
]

# Create a list of values to replace 'XNA'
replace_values = np.tile(fill_values, int(np.ceil(xna_count / len(fill_values))))

# Trim the list to match the number of 'XNA' values
replace_values = replace_values[:xna_count]

# Replace 'XNA' values with the calculated list of values
appli_data.loc[appli_data['ORGANIZATION_TYPE'] == 'XNA', 'ORGANIZATION_TYPE'] = replace_values

# Checking the updated count
org_count_update = appli_data['ORGANIZATION_TYPE'].value_counts()
print("\n#####################################\n")
print(org_count_update)
```

- → Handling the NaN values:

```python
1  nan_col = {}
2  for column in object_column:
3      nan_col[column] = appli_data[column].isnull().sum()
4
5  nan_col
```
✓ 0.0s

```
{'NAME_CONTRACT_TYPE': 0,
 'CODE_GENDER': 0,
 'FLAG_OWN_CAR': 0,
 'FLAG_OWN_REALTY': 0,
 'NAME_TYPE_SUITE': 192,
 'NAME_INCOME_TYPE': 0,
 'NAME_EDUCATION_TYPE': 0,
 'NAME_FAMILY_STATUS': 0,
 'NAME_HOUSING_TYPE': 0,
 'OCCUPATION_TYPE': 15654,
 'WEEKDAY_APPR_PROCESS_START': 0,
 'ORGANIZATION_TYPE': 0}
```

```python
1  # Calculate and print the count of each unique value in the 'NAME_TYPE_SUITE' column
2  print(appli_data['NAME_TYPE_SUITE'].value_counts(),"\n")
3
4  # Fill missing (NaN) values in the 'NAME_TYPE_SUITE' column with the most frequent value (mode)
5  # The mode is accessed using .mode()[0]
6  appli_data['NAME_TYPE_SUITE'].fillna(appli_data['NAME_TYPE_SUITE'].mode()[0], inplace=True)
7
8  # Check and print the updated count of missing (NaN) values in the 'NAME_TYPE_SUITE' column
9  nan_count_after_fill = appli_data['NAME_TYPE_SUITE'].isna().sum()
10 print(nan_count_after_fill)
```
✓ 0.0s

```
Unaccompanied     40435
Family             6549
Spouse, partner    1849
Children            542
Other_B             259
Other_A             137
Group of people      36
Name: NAME_TYPE_SUITE, dtype: int64

0
```

```python
1  # Print the count of NaN values in the 'OCCUPATION_TYPE' column
2  print(f'There are {appli_data.OCCUPATION_TYPE.isnull().sum()} NaN values present in the `Occupation_type`.')
3
4  # Print the number of unique non-NaN values in the 'OCCUPATION_TYPE' column
5  print(f'The unique values present in the `Occupational_type` is: \n{appli_data.OCCUPATION_TYPE.value_counts()}.')
6
7  # Get the top 12 most frequent non-NaN values in the 'OCCUPATION_TYPE' column
8  top_12_frequency = appli_data.OCCUPATION_TYPE.value_counts(ascending=False).nlargest(12).index
9
10 # Create a boolean index for NaN values in the 'OCCUPATION_TYPE' column
11 nan_index = appli_data['OCCUPATION_TYPE'].isna()
12
13 # Replace NaN values with a random choice from the top 12 most frequent non-NaN values
14 appli_data.loc[nan_index, 'OCCUPATION_TYPE'] = np.random.choice(top_12_frequency, nan_index.sum())
15
16 # Print the updated count of NaN values in the 'OCCUPATION_TYPE' column
17 print(f'There are {appli_data['OCCUPATION_TYPE'].isna().sum()} NaN values present in the `Occuption_type`.')
```
✓ 0.0s

```
There are 15654 NaN values present in the `Occupation_type`.
The unique values present in the `Occupational_type` is:
Laborers              8952
Sales staff           5160
Core staff            4434
Managers              3489
Drivers               3044
High skill tech staff 1852
Accountants           1621
Medicine staff        1403
Security staff        1140
Cooking staff          963
Cleaning staff         739
Private service staff  447
Low-skill Laborers     357
Waiters/barmen staff   228
Secretaries            212
Realty agents          123
HR staff               101
IT staff                80
Name: OCCUPATION_TYPE, dtype: int64.
There are 0 NaN values present in the `Occuption_type`.
```

Checking all the remaining NaN values columns with the dtypes and percentage.

```
Top 16 columns with highest missing values and their data types:

                             Missing Values   Data Types   Nan_values_percentage
            EXT_SOURCE_3              9944       float64              19.888398
 AMT_REQ_CREDIT_BUREAU_YEAR          6734       float64              13.468269
  AMT_REQ_CREDIT_BUREAU_QRT          6734       float64              13.468269
  AMT_REQ_CREDIT_BUREAU_MON          6734       float64              13.468269
 AMT_REQ_CREDIT_BUREAU_WEEK          6734       float64              13.468269
  AMT_REQ_CREDIT_BUREAU_DAY          6734       float64              13.468269
 AMT_REQ_CREDIT_BUREAU_HOUR          6734       float64              13.468269
    OBS_60_CNT_SOCIAL_CIRCLE          168       float64               0.336007
    OBS_30_CNT_SOCIAL_CIRCLE          168       float64               0.336007
    DEF_30_CNT_SOCIAL_CIRCLE          168       float64               0.336007
    DEF_60_CNT_SOCIAL_CIRCLE          168       float64               0.336007
              EXT_SOURCE_2            126       float64               0.252005
           AMT_GOODS_PRICE            38       float64               0.076002
              AMT_ANNUITY             1       float64               0.002000
     DAYS_LAST_PHONE_CHANGE           1       float64               0.002000
           CNT_FAM_MEMBERS           1       float64               0.002000
```

Replacing the NaN values.



```
The unique values present in the 'CNT_FAM_MEMBERS' column are [ 1.  2.  3.  4.  5.  6.  9.  7.  8. 10. 13. nan]
_____
2.0      25807
1.0      10873
3.0       8635
4.0       4000
5.0        592
6.0         68
7.0         12
8.0          6
9.0          2
10.0         2
13.0         1
Name: CNT_FAM_MEMBERS, dtype: int64
_____
There are 0 values present in the 'CNT_FAM_MEMBERS'.
_____
[ 1.  2.  3.  4.  5.  6.  9.  7.  8. 10. 13.]
```

```
1  round(appli_data[['AMT_ANNUITY','DAYS_LAST_PHONE_CHANGE','AMT_GOODS_PRICE']].describe(),2)
```
✓ 0.0s

|       | AMT_ANNUITY | DAYS_LAST_PHONE_CHANGE | AMT_GOODS_PRICE |
|-------|-------------|------------------------|------------------|
| count | 49998.00    | 49998.00               | 49961.00         |
| mean  | 27107.38    | -964.30                | 539060.04        |
| std   | 14562.94    | 829.49                 | 369853.25        |
| min   | 2052.00     | -4002.00               | 45000.00         |
| 25%   | 16456.50    | -1573.00               | 238500.00        |
| 50%   | 24939.00    | -755.00                | 450000.00        |
| 75%   | 34596.00    | -270.00                | 679500.00        |
| max   | 258025.50   | 0.00                   | 4050000.00       |

```
1  # Fill missing values in the 'AMT_ANNUITY' column with the median value of the column
2  appli_data['AMT_ANNUITY'].fillna(appli_data['AMT_ANNUITY'].median(), inplace = True)
3
4  # Fill missing values in the 'DAYS_LAST_PHONE_CHANGE' column with the mode (most frequent value) of the column
5  appli_data['DAYS_LAST_PHONE_CHANGE'].fillna(appli_data['DAYS_LAST_PHONE_CHANGE'].mode()[0], inplace = True)
6
7  # Fill missing values in the 'AMT_GOODS_PRICE' column with the median value of the column
8  appli_data['AMT_GOODS_PRICE'].fillna(appli_data['AMT_GOODS_PRICE'].median(), inplace = True)
```
✓ 0.0s

```
1  # List of columns to iterate through
2  columns_to_check = ['DEF_60_CNT_SOCIAL_CIRCLE', 'OBS_30_CNT_SOCIAL_CIRCLE', 'OBS_60_CNT_SOCIAL_CIRCLE', 'DEF_30_CNT_SOCIAL_CIRCLE']
3
4  # Set up the subplots
5  fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(18, 8))
6
7  # Flatten the axes to make it easier to index
8  axes = axes.flatten()
9
10 # Iterate through each specified column in the list 'columns_to_check'
11 for i, column in enumerate(columns_to_check):
12     # Fill missing (NaN) values in the current column with its mode (most frequent value)
13     appli_data[column].fillna(appli_data[column].mode()[0], inplace=True)
14
15     # Plot the distribution using Seaborn's distplot
16     sns.kdeplot(appli_data[column], ax=axes[i], shade = True)
17     axes[i].set_title(f'Distribution of {column}')
18
19 # Adjust layout for better visualization
20 plt.tight_layout()
21 plt.show()
```
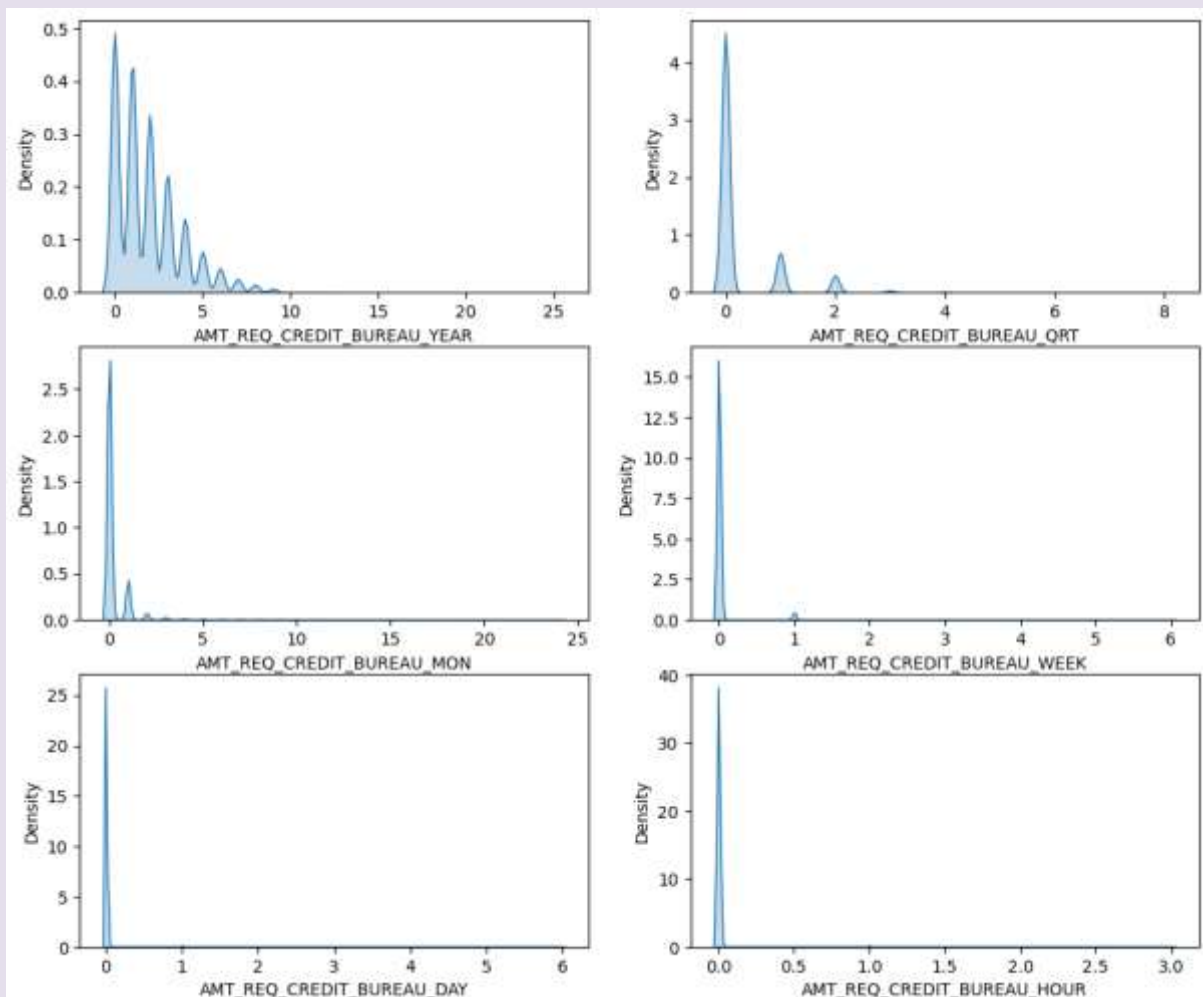✓ 2.0s

```python
# Iterate through each column and print its value counts
for column in columns_to_check:
    print(appli_data[column].value_counts())
    print("_____")

# Iterate through each specified column in the list 'columns_to_check'
for column in columns_to_check:

    # Fill missing (NaN) values in the current column with its mode (most frequent value)
    appli_data[column].fillna(appli_data[column].mode()[0], inplace=True)

    # Print the count of remaining NaN values after filling for the current column
    print(f'There are {appli_data[column].isna().sum()} NaN values present in {column}.')
    print("_____")
```

```python
plt.figure(figsize=(12, 10))

plt.subplot(3, 2, 1)
sns.kdeplot(appli_data['AMT_REQ_CREDIT_BUREAU_YEAR'], shade = True)

plt.subplot(3, 2, 2)
sns.kdeplot(appli_data['AMT_REQ_CREDIT_BUREAU_QRT'], shade = True)

plt.subplot(3, 2, 3)
sns.kdeplot(appli_data['AMT_REQ_CREDIT_BUREAU_MON'], shade = True)

plt.subplot(3, 2, 4)
sns.kdeplot(appli_data['AMT_REQ_CREDIT_BUREAU_WEEK'], shade = True)

plt.subplot(3, 2, 5)
sns.kdeplot(appli_data['AMT_REQ_CREDIT_BUREAU_DAY'], shade = True)

plt.subplot(3, 2, 6)
sns.kdeplot(appli_data['AMT_REQ_CREDIT_BUREAU_HOUR'], shade = True)

plt.show()
```

```
   ♥ Click here to ask Blackbox to help you code faster
1  col = ['AMT_REQ_CREDIT_BUREAU_YEAR', 'AMT_REQ_CREDIT_BUREAU_QRT',
2  'AMT_REQ_CREDIT_BUREAU_MON','AMT_REQ_CREDIT_BUREAU_WEEK',
3  'AMT_REQ_CREDIT_BUREAU_DAY','AMT_REQ_CREDIT_BUREAU_HOUR']
4
5  # Iterate through each column and print its value counts amd filling NaN values.
6  for column in col:
7      print(appli_data[column].value_counts())
8      print("\n")
9      appli_data[column].fillna(appli_data[column].mode()[0], inplace = True)
10     print(f'There are {appli_data[column].isna().sum()} NaN values present in {column}.')
```

-→ Previous Application:

The previous application dataset contains 49999 rows and 37 columns.

Checking the NaN values in the form of percentage to remove the columns who are having the NaN values greater than 35%.

| | column_name | null_count | null_percentage |
|---|---|---|---|
| 14 | RATE_INTEREST_PRIVILEGED | 49834 | 99.67 |
| 13 | RATE_INTEREST_PRIMARY | 49834 | 99.67 |
| 12 | RATE_DOWN_PAYMENT | 25198 | 50.40 |
| 6 | AMT_DOWN_PAYMENT | 25198 | 50.40 |
| 20 | NAME_TYPE_SUITE | 24243 | 48.49 |
| 36 | NFLAG_INSURED_ON_APPROVAL | 19160 | 38.32 |
| 31 | DAYS_FIRST_DRAWING | 19160 | 38.32 |
| 32 | DAYS_FIRST_DUE | 19160 | 38.32 |
| 33 | DAYS_LAST_DUE_1ST_VERSION | 19160 | 38.32 |
| 34 | DAYS_LAST_DUE | 19160 | 38.32 |
| 35 | DAYS_TERMINATION | 19160 | 38.32 |
| 7 | AMT_GOODS_PRICE | 10744 | 21.49 |
| 3 | AMT_ANNUITY | 10592 | 21.18 |
| 28 | CNT_PAYMENT | 10592 | 21.18 |
| 30 | PRODUCT_COMBINATION | 8 | 0.02 |
| 25 | CHANNEL_TYPE | 0 | 0.00 |
| 24 | NAME_PRODUCT_TYPE | 0 | 0.00 |
| 29 | NAME_YIELD_GROUP | 0 | 0.00 |
| 26 | SELLERPLACE_AREA | 0 | 0.00 |
| 27 | NAME_SELLER_INDUSTRY | 0 | 0.00 |
| 22 | NAME_GOODS_CATEGORY | 0 | 0.00 |
| 23 | NAME_PORTFOLIO | 0 | 0.00 |
| 0 | SK_ID_PREV | 0 | 0.00 |
| 21 | NAME_CLIENT_TYPE | 0 | 0.00 |
| 19 | CODE_REJECT_REASON | 0 | 0.00 |
| 1 | SK_ID_CURR | 0 | 0.00 |

```
    ♀ Click here to ask Blackbox to help you code faster
1 # displaying all the collumn who are having more than 40% NaN values in it.
2 null_value_gre40 = null_percentage[null_percentage > 35].index
3 null_value_gre40

✓ 0.0s

Index(['AMT_DOWN_PAYMENT', 'RATE_DOWN_PAYMENT', 'RATE_INTEREST_PRIMARY',
       'RATE_INTEREST_PRIVILEGED', 'NAME_TYPE_SUITE', 'DAYS_FIRST_DRAWING',
       'DAYS_FIRST_DUE', 'DAYS_LAST_DUE_1ST_VERSION', 'DAYS_LAST_DUE',
       'DAYS_TERMINATION', 'NFLAG_INSURED_ON_APPROVAL'],
      dtype='object')
```

```
    ♀ Click here to ask Blackbox to help you code faster
1 col_drop = ['NFLAG_LAST_APPL_IN_DAY', 'FLAG_LAST_APPL_PER_CONTRACT', 'HOUR_APPR_PROCESS_START', 'WEEKDAY_APPR_PROCESS_START']
2
3 previ_data.drop(columns = col_drop, inplace = True)

✓ 0.0s
```

```
    ♀ Click here to ask Blackbox to help you code faster
1 # Display the count of missing values for each column, sorted in descending order,
2 # and show the top 10 columns with the highest missing values
3 missing_values_top10 = previ_data.isna().sum().sort_values(ascending=False).nlargest(10)
4
5 # Print the results
6 print(missing_values_top10)

✓ 0.3s

AMT_GOODS_PRICE         10744
AMT_ANNUITY             10592
CNT_PAYMENT             10592
PRODUCT_COMBINATION         8
NAME_CLIENT_TYPE            0
NAME_YIELD_GROUP            0
NAME_SELLER_INDUSTRY        0
SELLERPLACE_AREA            0
CHANNEL_TYPE               0
NAME_PRODUCT_TYPE          0
dtype: int64
```

➔ Filling NaN Values with top frequency:

```
  Click here to ask Blackbox to help you code faster
1 # Count the occurrences of each unique value in the 'PRODUCT_COMBINATION' column
2 val_count = previ_data.PRODUCT_COMBINATION.value_counts()
3
4 # Create a bar plot to visualize the counts of each product combination
5 plt.figure(figsize=(10, 5))
6 fig = sns.barplot(x=val_count.values, y=val_count.index, estimator='sum', errorbar=None)
7
8 # Add labels to the bars indicating the sum of counts
9 fig.bar_label(fig.containers[0], fontsize=10)
10
11 # Display the bar plot
12 plt.show()
13
14 # Fill missing values in the 'PRODUCT_COMBINATION' column with a specified default value
15 previ_data.PRODUCT_COMBINATION.fillna("POS household without interest", inplace=True)
```



```
  Click here to ask Blackbox to help you code faster
1 # Create a box plot to visualize the distribution of 'CNT_PAYMENT' column values
2 fig = px.box(previ_data, y = 'CNT_PAYMENT', title = "Box plot of CNT_PAYMENT")
3 fig.show()
4
5 # Fill missing values in the 'CNT_PAYMENT' column with the median value of the column
6 previ_data['CNT_PAYMENT'] = previ_data["CNT_PAYMENT"].fillna(previ_data["CNT_PAYMENT"].median())
```
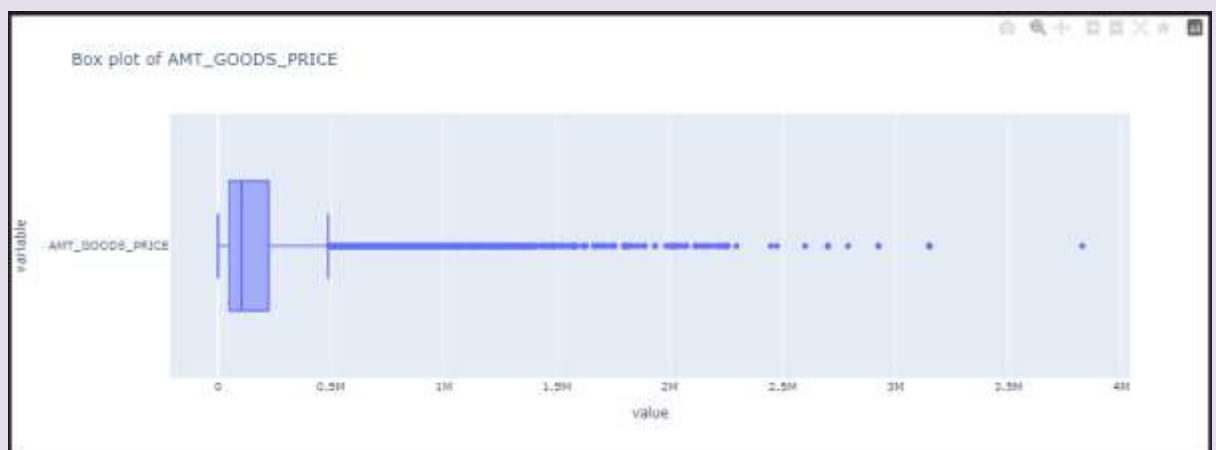


Box plot of CNT_PAYMENT

```
 Click here to ask Blackbox to help you code faster
1  # Create a box plot using Plotly Express for the 'AMT_ANNUITY' column in 'previ_data'
2  fig = px.box(previ_data, x="AMT_ANNUITY", title="Box plot of AMT_ANNUITY")
3  fig.show()
4
5  # Create a Kernel Density Estimate (KDE) plot using Seaborn for the 'AMT_ANNUITY' column
6  sns.kdeplot(previ_data['AMT_ANNUITY'], shade=True)
7
8  # Add a title to the KDE plot
9  plt.title("Density Plot of AMT_ANNUITY")
10 plt.show()
11
12 # Fill missing values in the 'AMT_ANNUITY' column with the median value
13 previ_data['AMT_ANNUITY'].fillna(previ_data['AMT_ANNUITY'].median(), inplace=True)
```

```
💡 Click here to ask Blackbox to help you code faster
1  # Create a box plot using Plotly Express for the 'AMT_GOODS_PRICE' column in 'previ_data'
2  fig = px.box(previ_data, x=['AMT_GOODS_PRICE'], title="Box plot of AMT_GOODS_PRICE")
3  fig.show()
4
5  # Create a Kernel Density Estimate (KDE) plot using Seaborn for the 'AMT_GOODS_PRICE' column
6  sns.kdeplot(previ_data['AMT_GOODS_PRICE'], shade=True)
7
8  # Display the plots
9  plt.show()
10
11  # Fill missing values in the 'AMT_GOODS_PRICE' column with the median value
12  previ_data['AMT_GOODS_PRICE'].fillna(previ_data['AMT_GOODS_PRICE'].median(), inplace=True)
✓ 14s
```



Box plot of AMT_GOODS_PRICE

➔ Checking all the colums who are having XNA and XAP in it to fill with top counts.

```python
# Click here to ask Blackbox to help you code faster
1  # List of column with XNA and XAP as unique values.
2  columns_with_xna_xap = []
3
4  # Iterating through each column in dataframe
5  for column in previ_data.columns:
6      # Checking if XNA and XAP is present in the unique values of the column
7      if "XNA" in previ_data[column].unique() or "XAP" in previ_data[column].unique():
8          columns_with_xna_xap.append(column)
```

```
['NAME_CONTRACT_TYPE',
 'NAME_CASH_LOAN_PURPOSE',
 'NAME_PAYMENT_TYPE',
 'CODE_REJECT_REASON',
 'NAME_CLIENT_TYPE',
 'NAME_GOODS_CATEGORY',
 'NAME_PORTFOLIO',
 'NAME_PRODUCT_TYPE',
 'NAME_SELLER_INDUSTRY',
 'NAME_YIELD_GROUP']
```

```python
# Click here to ask Blackbox to help you code faster
1   # Replace occurrences of "XNA" in the 'NAME_CONTRACT_TYPE' column with the mode (most frequent value)
2   # Mode is used as a replacement for "XNA" to handle missing or undefined values in a categorical column
3   previ_data['NAME_CONTRACT_TYPE'].replace("XNA",previ_data['NAME_CONTRACT_TYPE'].mode()[0], inplace = True)
4
5   # Replace occurrences of "XNA" in the 'NAME_PAYMENT_TYPE' column with the mode (most frequent value)
6   # Mode is used as a replacement for "XNA" to handle missing or undefined values in a categorical column
7   previ_data['NAME_PAYMENT_TYPE'].replace("XNA", previ_data['NAME_PAYMENT_TYPE'].mode()[0], inplace=True)
8
9   # Replace occurrences of "XNA" in the 'NAME_CLIENT_TYPE' column with the mode (most frequent value)
10  # Mode is used as a replacement for "XNA" to handle missing or undefined values in a categorical column
11  previ_data['NAME_CLIENT_TYPE'].replace("XNA", previ_data['NAME_CLIENT_TYPE'].mode()[0], inplace=True)
12
13  # Replace occurrences of "XNA" in the 'NAME_PORTFOLIO' column with the mode (most frequent value)
14  # Mode is used as a replacement for "XNA" to handle missing or undefined values in a categorical column
15  previ_data['NAME_PORTFOLIO'].replace("XNA", previ_data['NAME_PORTFOLIO'].mode()[0], inplace=True)
16
17  # Replace occurrences of "XNA" in the '"NAME_PRODUCT_TYPE"' column with the mode (most frequent value)
18  # Mode is used as a replacement for "XNA" to handle missing or undefined values in a categorical column
19  previ_data["NAME_PRODUCT_TYPE"].replace("XNA", "x-sell", inplace=True)
20
21  # Replace occurrences of "XNA" in the 'NAME_YIELD_GROUP' column with the mode (most frequent value)
22  # Mode is used as a replacement for "XNA" to handle missing or undefined values in a categorical column
23  previ_data['NAME_YIELD_GROUP'].replace("XNA", 'middle', inplace=True)
```

```python
1  # Display the counts of unique values in the 'NAME_CASH_LOAN_PURPOSE' column before modification
2  # Get the values excluding the top two
3  print(f"'NAME_CASH_LOAN_PURPOSE' column before modification:\n{previ_data['NAME_CASH_LOAN_PURPOSE'].value_counts()}")
4  print("——————————————————————————————————————————————————")
5
6  # Identify the 'XNA' values in the 'NAME_CASH_LOAN_PURPOSE' column
7  xna_values = previ_data['NAME_CASH_LOAN_PURPOSE'] == 'XNA'
8
9  # Divide 'XNA' values into four equal parts
10 division_size = len(previ_data[xna_values]) // 3
11 indices = np.where(xna_values)[0]  # Get the indices of 'XNA' values
12
13 # Assign parts to 'Repairs', 'Other', 'Urgent needs ' respectively
14 previ_data.loc[indices[:division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Repairs'
15 previ_data.loc[indices[division_size:2*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Other'
16 previ_data.loc[indices[2*division_size:3*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Urgent needs'
17
18 # Identify the 'XAP' values in the 'NAME_CASH_LOAN_PURPOSE' column
19 xap_values = previ_data['NAME_CASH_LOAN_PURPOSE'] == 'XAP'
20
21 # Divide 'XAP' values into four equal parts
22 division_size = len(previ_data[xap_values]) // 8
23 indices = np.where(xap_values)[0]  # Get the indices of 'XAP' values
24
25 # Assign parts to 'Repairs', 'Other', 'Urgent needs', 'Buying a used car','Building a house or an annex',
26 # 'Medicine', 'Payments on other loans', 'Everyday expenses', respectively
27 previ_data.loc[indices[:division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Repairs'
28 previ_data.loc[indices[division_size:2*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Other'
29 previ_data.loc[indices[2*division_size:3*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Urgent needs'
30 previ_data.loc[indices[3*division_size:4*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Buying a used car'
31 previ_data.loc[indices[4*division_size:5*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Building a house or an annex'
32 previ_data.loc[indices[5*division_size:6*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Medicine'
33 previ_data.loc[indices[6*division_size:7*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Payments on other loans'
34 previ_data.loc[indices[7*division_size:8*division_size], 'NAME_CASH_LOAN_PURPOSE'] = 'Everyday expenses'
35
36 # Renaming the remaining XNA and XAP to Refusal to name the goal
37 previ_data['NAME_CASH_LOAN_PURPOSE'].replace(['XNA', 'XAP'], "Refusal to name the goal", inplace = True)
38
39 # Display the counts of unique values in the 'NAME_CASH_LOAN_PURPOSE' column after modification
40 print(f"'NAME_CASH_LOAN_PURPOSE' column after modification:\n{previ_data['NAME_CASH_LOAN_PURPOSE'].value_counts()}")
```

```python
1  # Display the counts of unique values in the 'CODE_REJECT_REASON' column before modification
2  print(f"'CODE_REJECT_REASON' column before modification:\n{previ_data['CODE_REJECT_REASON'].value_counts()}")
3  print("——————————————————————————————————————————————————")
4
5  # Replace 'XNA' with 'HC' in the 'CODE_REJECT_REASON' column
6  previ_data['CODE_REJECT_REASON'].replace('XNA', 'HC', inplace=True)
7
8  # Identify the 'XAP' values in the 'CODE_REJECT_REASON' column
9  xap_values = previ_data['CODE_REJECT_REASON'] == 'XAP'
10
11 # Divide 'XAP' values into four equal parts
12 division_size = len(previ_data[xap_values]) // 4
13 indices = np.where(xap_values)[0]  # Get the indices of 'XAP' values
14
15 # Assign parts to 'HC', 'LIMIT', 'SCO', 'CLIENT' respectively
16 previ_data.loc[indices[:division_size], 'CODE_REJECT_REASON'] = 'HC'
17 previ_data.loc[indices[division_size:2*division_size], 'CODE_REJECT_REASON'] = 'LIMIT'
18 previ_data.loc[indices[2*division_size:3*division_size], 'CODE_REJECT_REASON'] = 'SCO'
19 previ_data.loc[indices[3*division_size:], 'CODE_REJECT_REASON'] = 'CLIENT'
20
21 # Display the counts of unique values in the 'CODE_REJECT_REASON' column after modification
22 print(f"'CODE_REJECT_REASON' column after modification:\n{previ_data['CODE_REJECT_REASON'].value_counts()}")
```
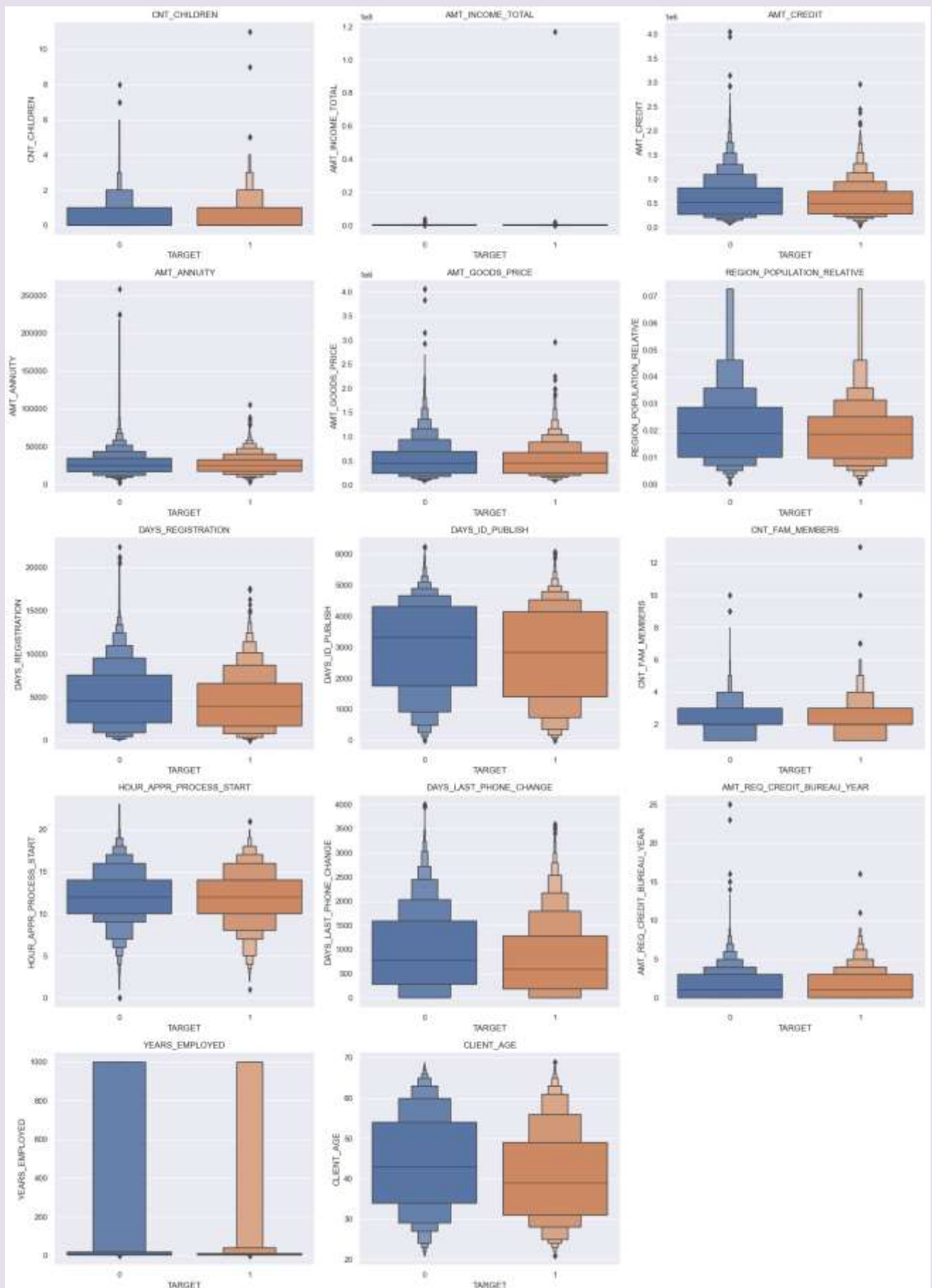
```python
1  # Display the counts of unique values in the 'NAME_GOODS_CATEGORY' column before modification
2  print(f"'NAME_GOODS_CATEGORY' column before modification:\n{previ_data['NAME_GOODS_CATEGORY'].value_counts()}")
3  print("—————————————————————————————————————————————————")
4
5  # Identify the 'XNA' values in the 'NAME_GOODS_CATEGORY' column
6  xna_values = previ_data['NAME_GOODS_CATEGORY'] == 'XNA'
7
8  # Divide 'XNA' values into four equal parts
9  division_size = len(previ_data[xna_values]) // 5
10 indices = np.where(xna_values)[0]  # Get the indices of 'XNA' values
11
12 # Assign parts to 'Mobile', 'Consumer Electronics', 'Computers', 'Audio/Video', 'Furniture' respectively
13 previ_data.loc[indices[:division_size], 'NAME_GOODS_CATEGORY'] = 'Mobile'
14 previ_data.loc[indices[division_size:2*division_size], 'NAME_GOODS_CATEGORY'] = 'Consumer Electronics'
15 previ_data.loc[indices[2*division_size:3*division_size], 'NAME_GOODS_CATEGORY'] = 'Computers'
16 previ_data.loc[indices[3*division_size:4*division_size], 'NAME_GOODS_CATEGORY'] = 'Audio/Video'
17 previ_data.loc[indices[4*division_size:], 'NAME_GOODS_CATEGORY'] = 'Furniture'
18
19 # Display the counts of unique values in the 'NAME_GOODS_CATEGORY' column after modification
20 print(f"'NAME_GOODS_CATEGORY' column after modification:\n{previ_data['NAME_GOODS_CATEGORY'].value_counts()}")
```

```python
1  # Display the counts of unique values in the 'NAME_SELLER_INDUSTRY' column before modification
2  print(f"'NAME_SELLER_INDUSTRY' column before modification:\n{previ_data['NAME_SELLER_INDUSTRY'].value_counts()}")
3  print("—————————————————————————————————————————————————")
4
5  # Identify the 'XNA' values in the 'NAME_SELLER_INDUSTRY' column
6  xna_values = previ_data['NAME_SELLER_INDUSTRY'] == 'XNA'
7
8  # Divide 'XNA' values into four equal parts
9  division_size = len(previ_data[xna_values]) // 2
10 indices = np.where(xna_values)[0]  # Get the indices of 'XNA' values
11
12 # Assign parts to 'Consumer electronics', 'Connectivity', 'Computers', 'Audio/Video', 'Furniture' respectively
13 previ_data.loc[indices[:division_size], 'NAME_SELLER_INDUSTRY'] = 'Consumer electronics'
14 previ_data.loc[indices[division_size:], 'NAME_SELLER_INDUSTRY'] = 'Connectivity'
15
16 # Display the counts of unique values in the 'NAME_SELLER_INDUSTRY' column after modification
17 print(f"'NAME_SELLER_INDUSTRY' column after modification:\n{previ_data['NAME_SELLER_INDUSTRY'].value_counts()}")
```
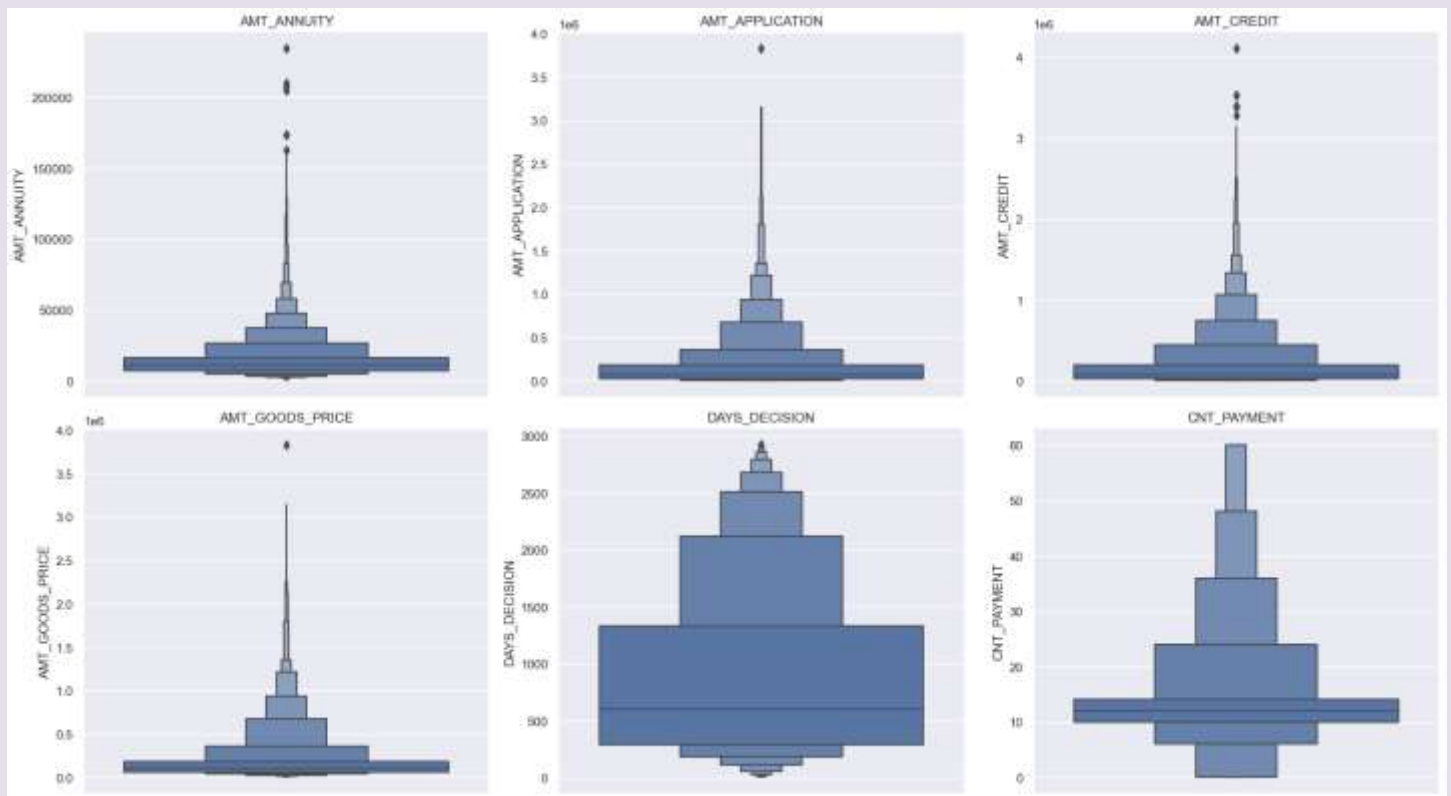
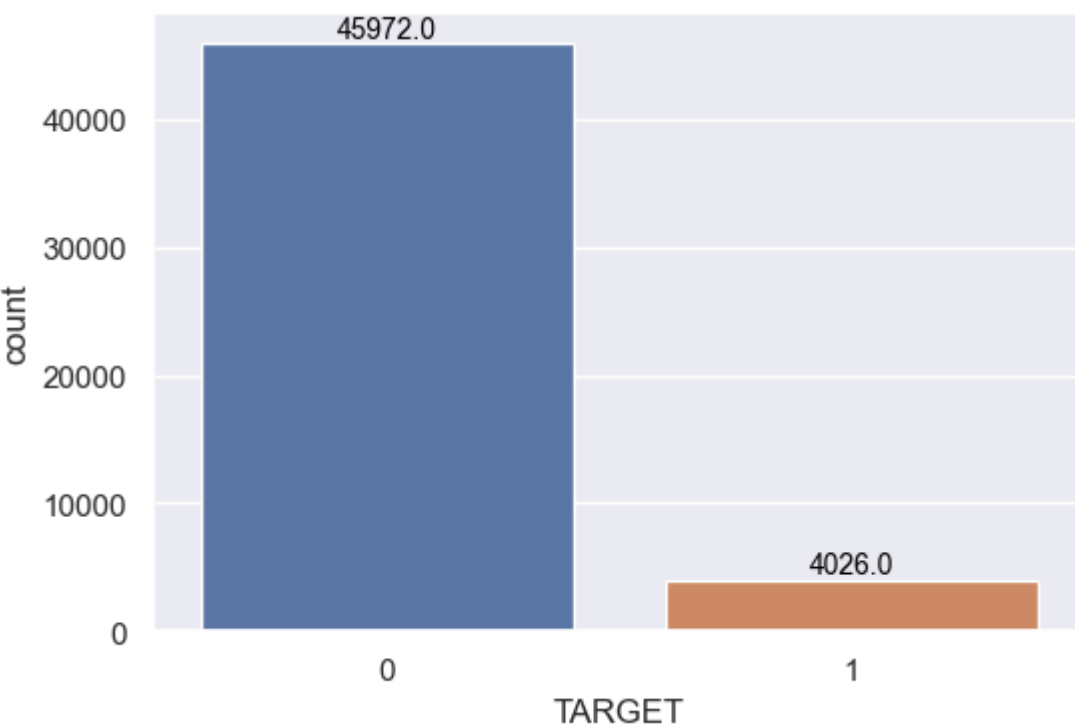✓ 0.0s

Detecting Outliers:

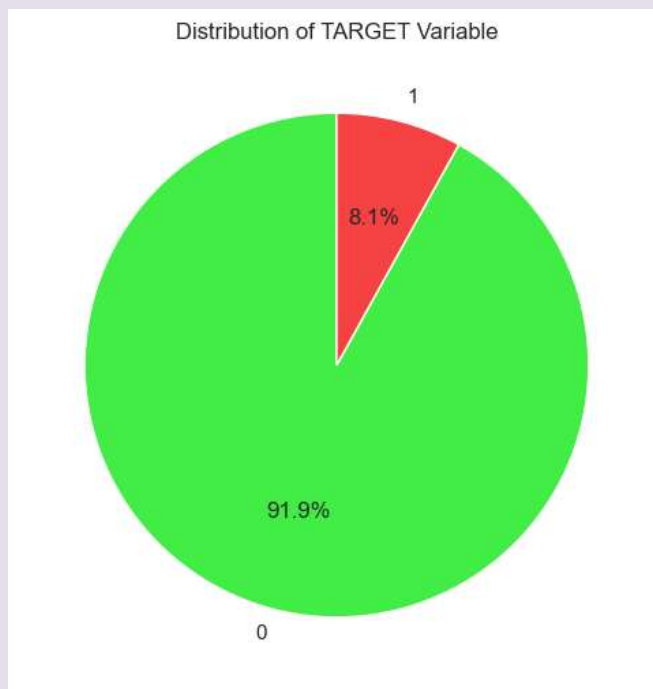1. Application data:

2. Previous Application Data:

Data Imbalance:

```python
# Assuming 'appli_data' is your DataFrame
sns.set(style="darkgrid")
plt.figure(figsize=(6, 4))

# Create a countplot
ax = sns.countplot(x="TARGET", data=appli_data)

# Add counts on top of each bar
for p in ax.patches:
    ax.annotate(f'{p.get_height()}', (p.get_x() + p.get_width() / 2., p.get_height()),
                ha='center', va='center', fontsize=10, color='black', xytext=(0, 5),
                textcoords='offset points')

plt.show()
```
✓ 0.2s



```python
# Assuming 'appli_data' is your DataFrame
target_counts = appli_data['TARGET'].value_counts()

# Create a pie chart
plt.figure(figsize=(6, 6))
plt.pie(target_counts, labels=target_counts.index, autopct='%1.1f%%', startangle=90, colors=['#42ed45', '#f54242'])
plt.title('Distribution of TARGET Variable')
plt.show()
```
✓ 0.1s

## Distribution of TARGET Variable



```
1   class_counts = appli_data['TARGET'].value_counts()
2   imbalance_ratio = class_counts[0] / class_counts[1]
3
4   "Imbalance Ratio:", round(imbalance_ratio,2)
```

✓ 0.0s

('Imbalance Ratio:', 11.42)

```
1   app_data_fil = appli_data['TARGET'].value_counts()
2   class_counts = app_data_fil
3   imbalance_ratio = class_counts[0] / class_counts[1]
4
5   "Imbalance Ratio: 1:{:.0f}".format(imbalance_ratio)
```

✓ 0.0s

'Imbalance Ratio: 1:11'

## Univariate Analysis:



The age group of applicant of age 30-40 and 40-50 and having highest numbers of applicants for loan process.



The employees who are having less than 5 years of work experience are applying for the loan and are not having problems while repaying loan.



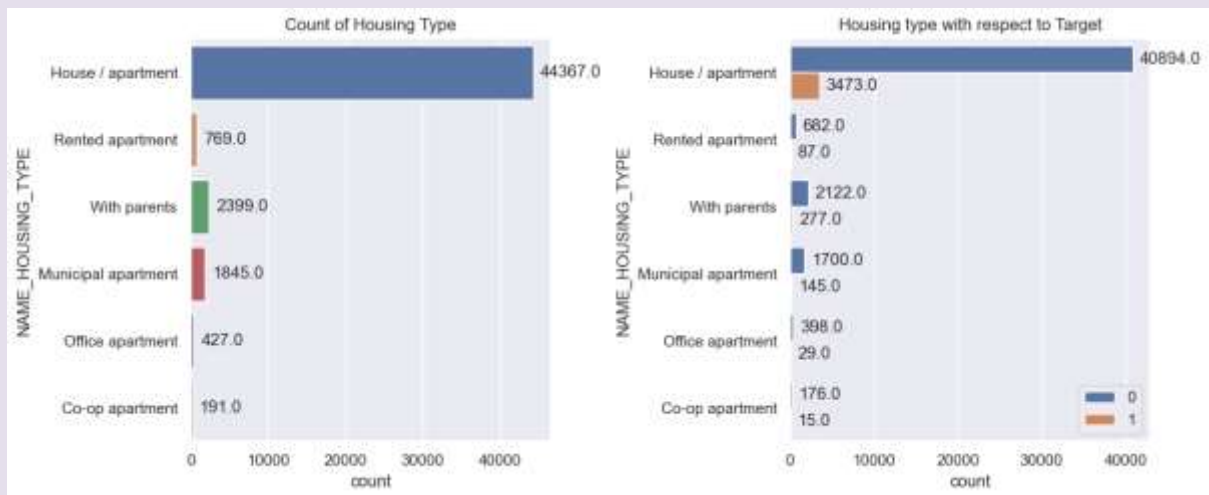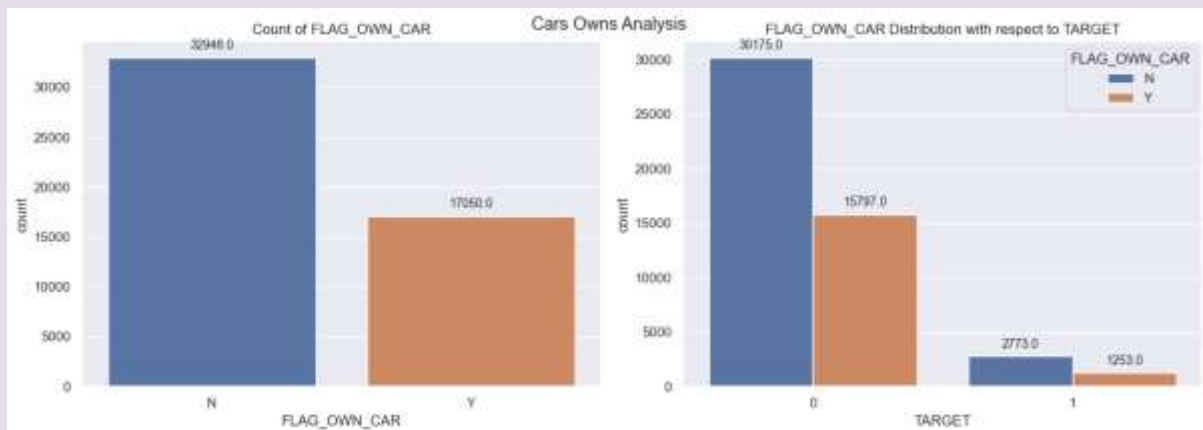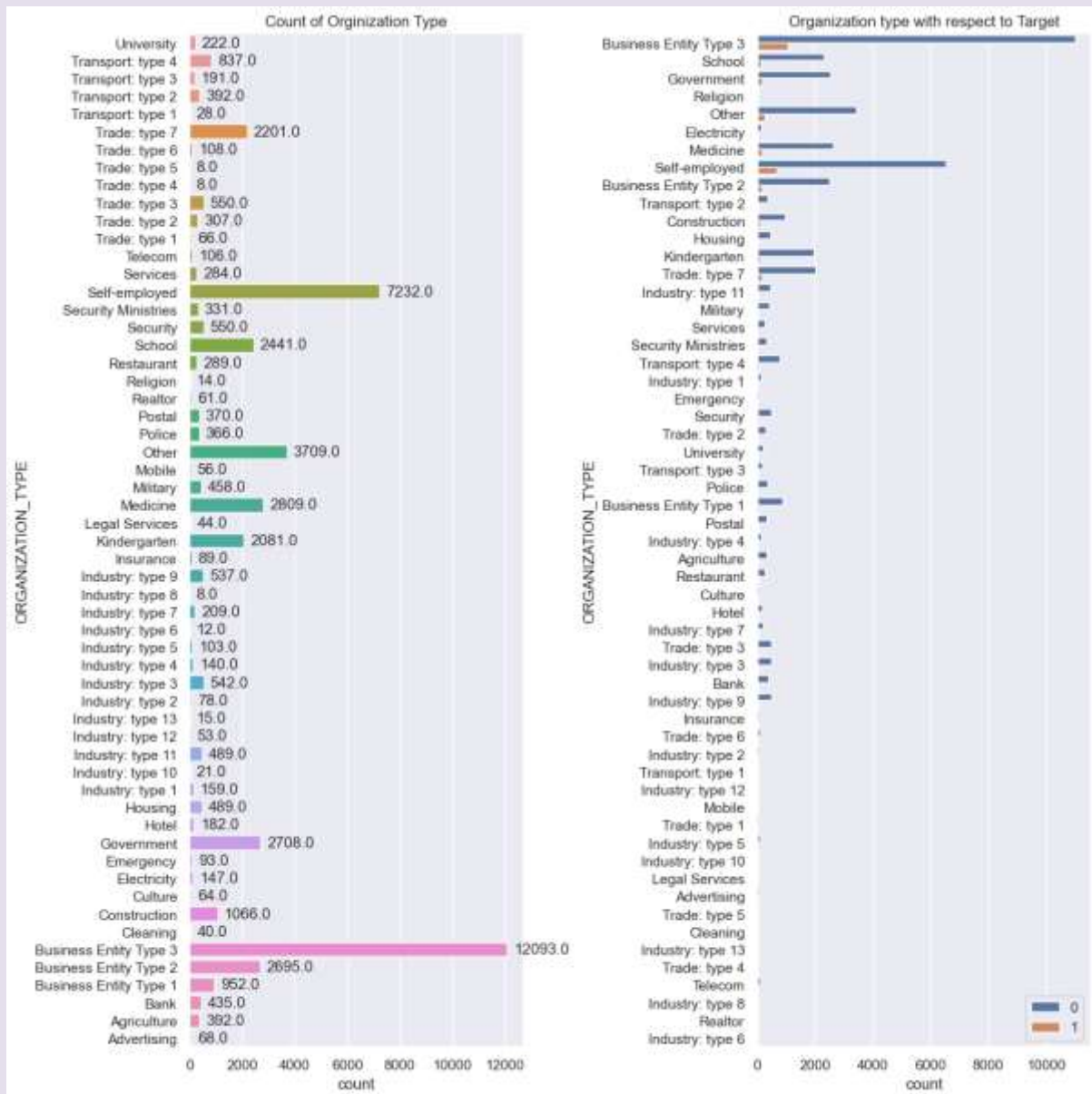It seems like Female client applied higher than male client for loan.

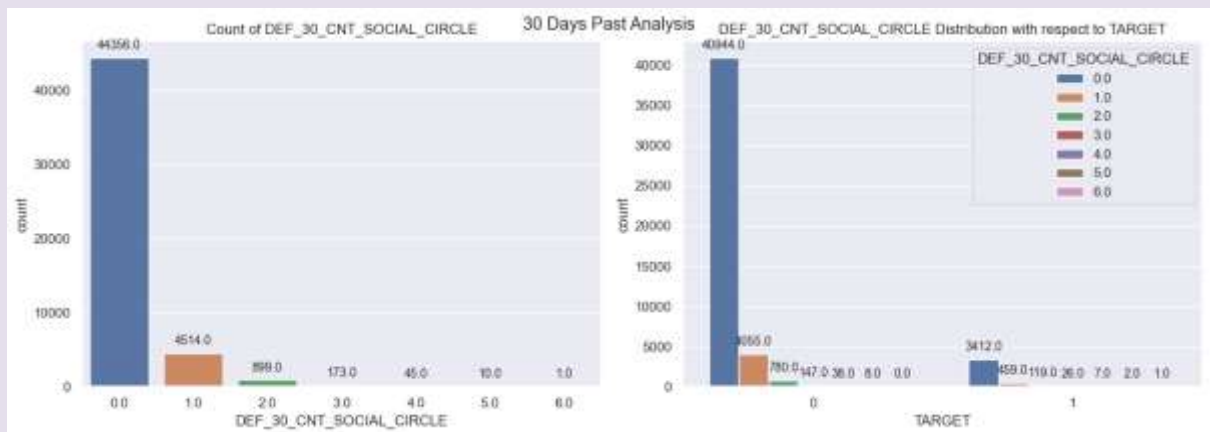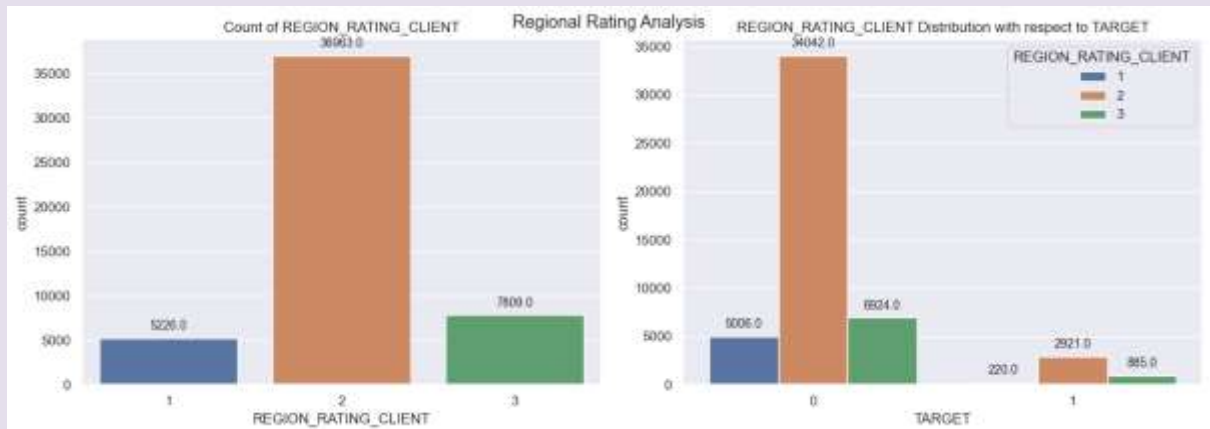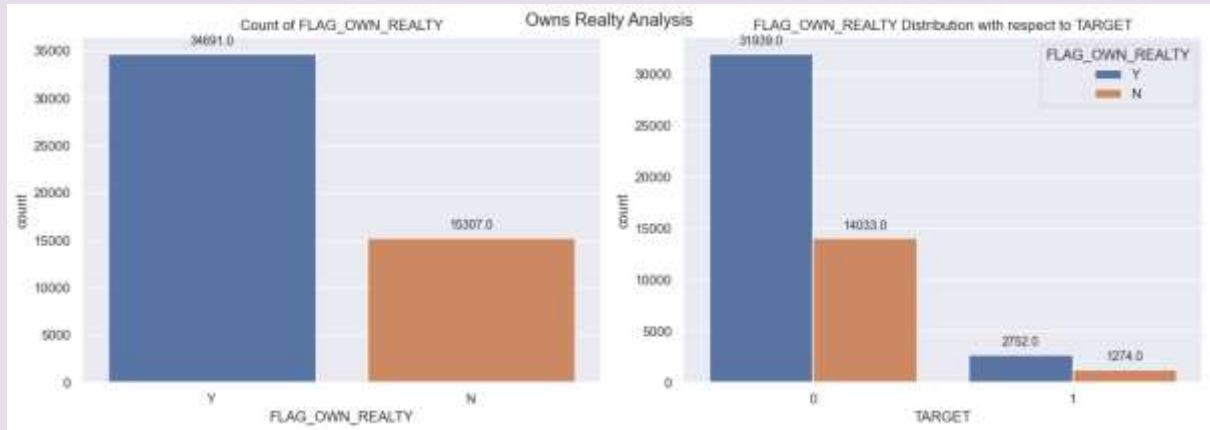We can see that 90% of peoples are taking Cash loans and 10% are applying for Revolving loans.
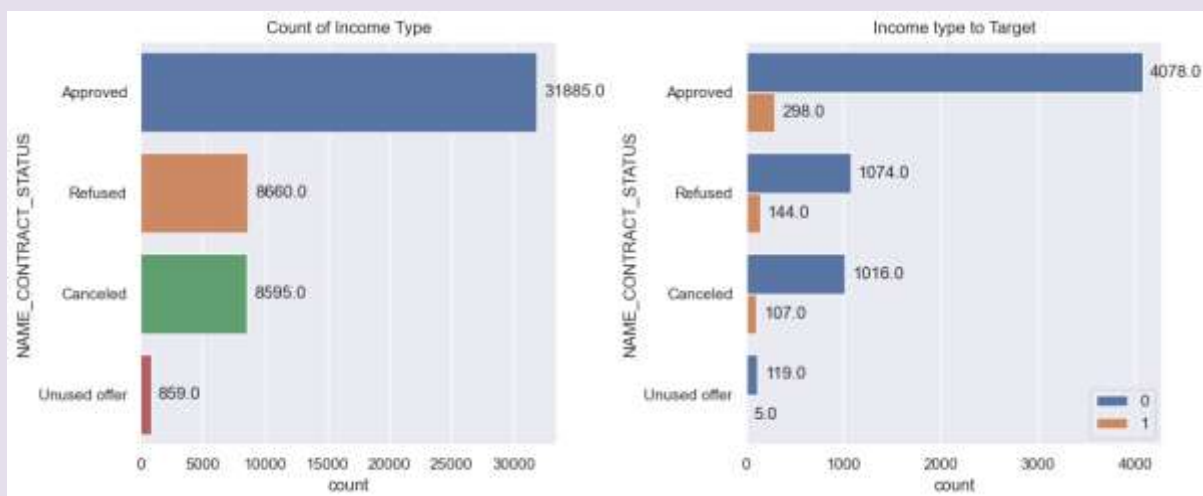
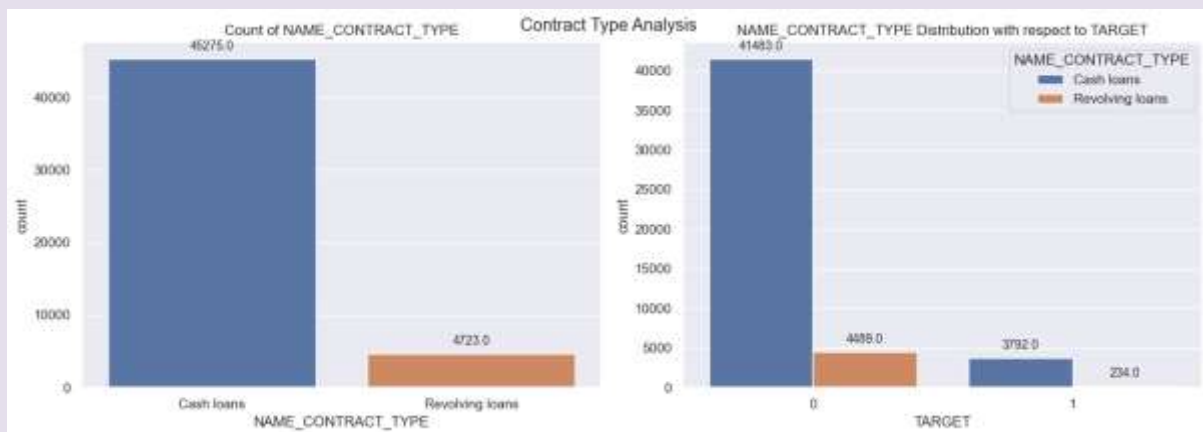Count of Family Status Type / Family Status with respect to Target

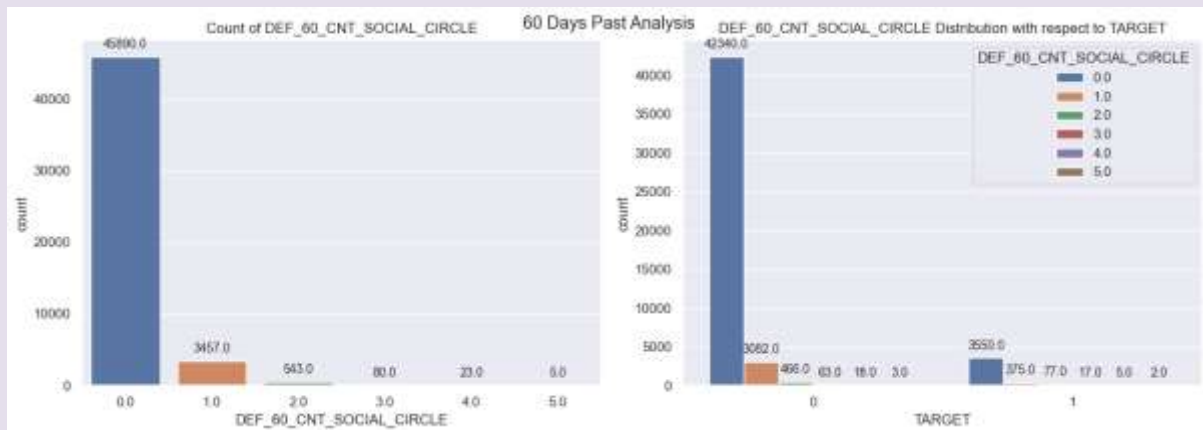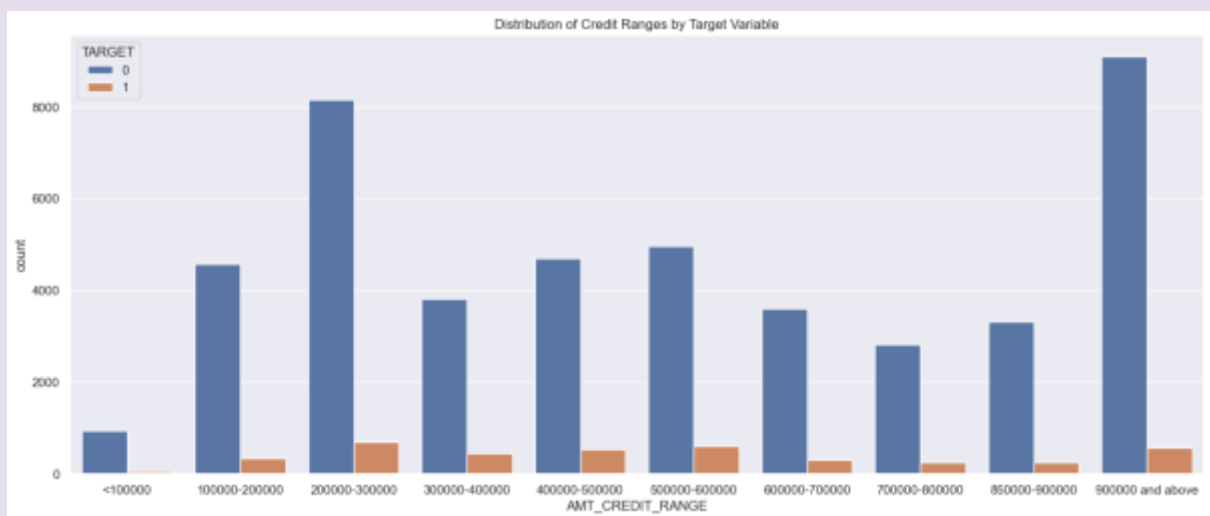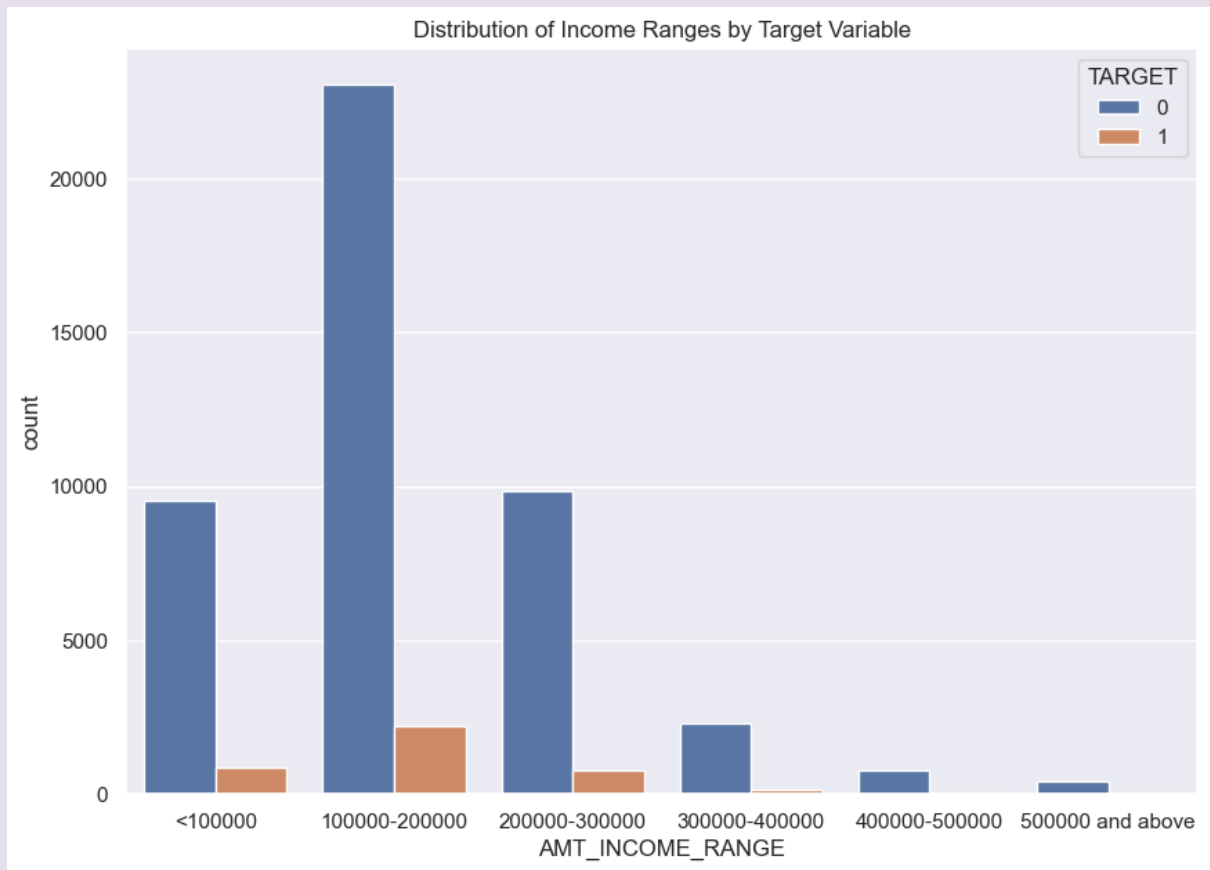We can see that married peoples are taking more loan than the single ,civil, window and separated people.



Count of Housing Type / Housing type with respect to Target



Count of Occupation Type / Occupation type with respect to Target

Count of Orginization Type

Organization type with respect to Target

Cars Owns Analysis

Count of FLAG_OWN_CAR

FLAG_OWN_CAR Distribution with respect to TARGET

Owns Realty Analysis

Count of FLAG_OWN_REALTY — 34691.0 (Y), 15307.0 (N)

FLAG_OWN_REALTY Distribution with respect to TARGET — TARGET 0: 31939.0 (Y), 14033.0 (N); TARGET 1: 2752.0 (Y), 1274.0 (N)



Regional Rating Analysis

Count of REGION_RATING_CLIENT — 5226.0 (1), 36963.0 (2), 7800.0 (3)

REGION_RATING_CLIENT Distribution with respect to TARGET — TARGET 0: 5006.0 (1), 34042.0 (2), 6924.0 (3); TARGET 1: 220.0 (1), 2921.0 (2), 885.0 (3)



30 Days Past Analysis

Count of DEF_30_CNT_SOCIAL_CIRCLE — 44356.0 (0.0), 4514.0 (1.0), 899.0 (2.0), 173.0 (3.0), 45.0 (4.0), 10.0 (5.0), 1.0 (6.0)

DEF_30_CNT_SOCIAL_CIRCLE Distribution with respect to TARGET — TARGET 0: 40944.0, 4055.0, 780.0, 147.0, 38.0, 8.0, 0.0; TARGET 1: 3412.0, 459.0, 119.0, 26.0, 7.0, 2.0, 1.0

60 Days Past Analysis

Count of DEF_60_CNT_SOCIAL_CIRCLE

DEF_60_CNT_SOCIAL_CIRCLE Distribution with respect to TARGET

Contract Type Analysis

Count of NAME_CONTRACT_TYPE

NAME_CONTRACT_TYPE Distribution with respect to TARGET

Count of Income Type

Income type to Target

Distribution of Income Ranges by Target Variable



Distribution of Credit Ranges by Target Variable

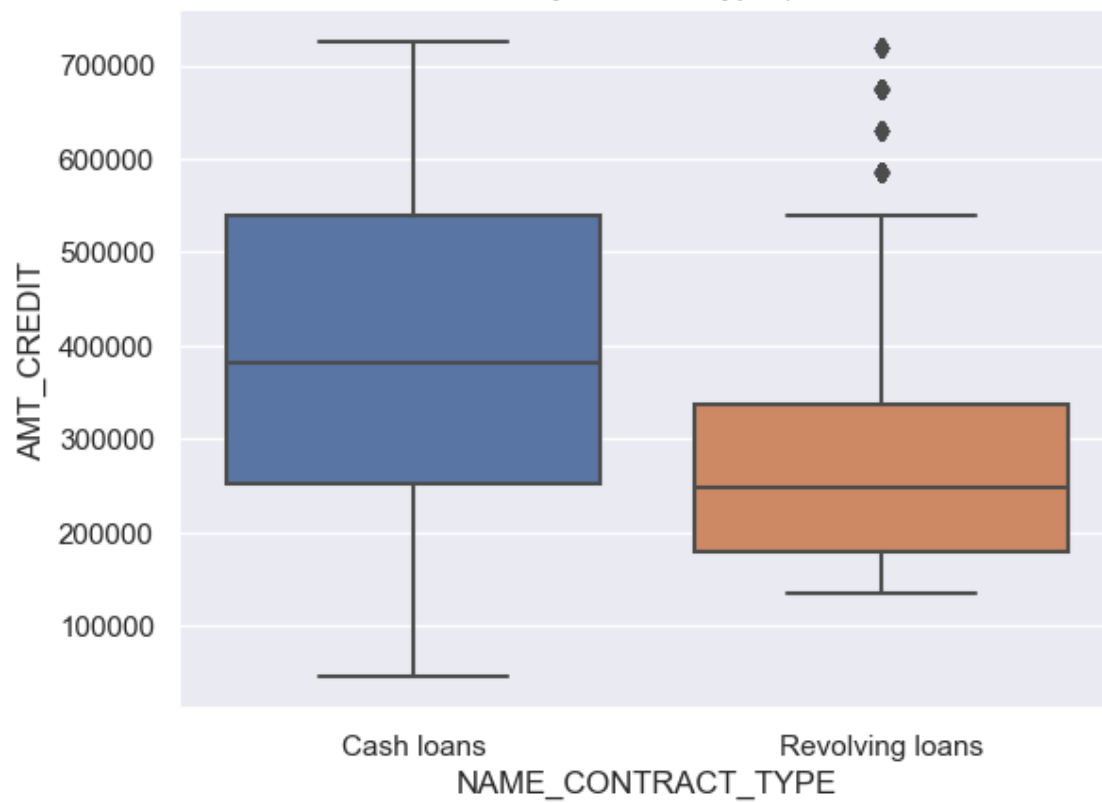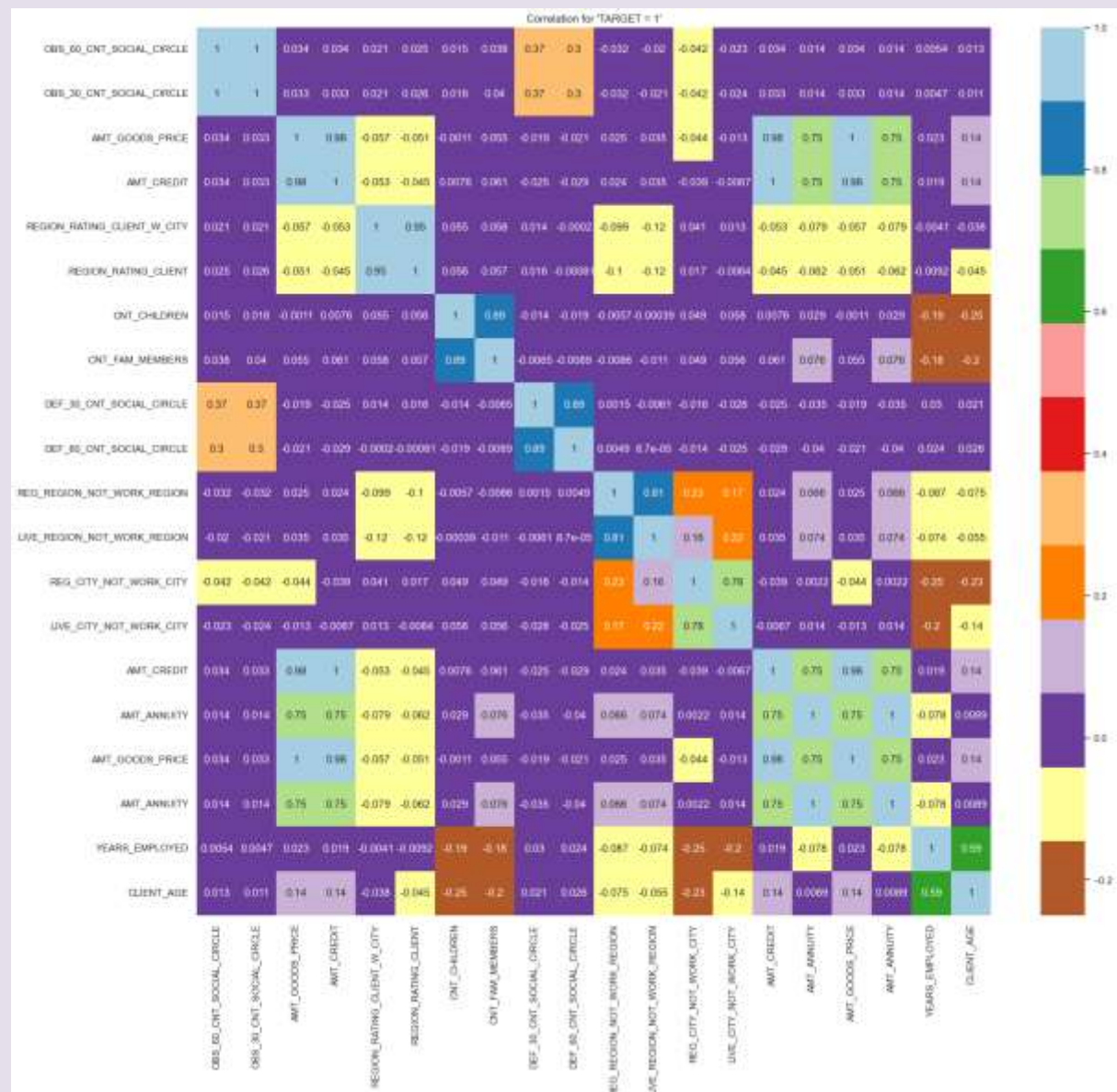Distribution of Income by Family Status (Income <= $135,000)

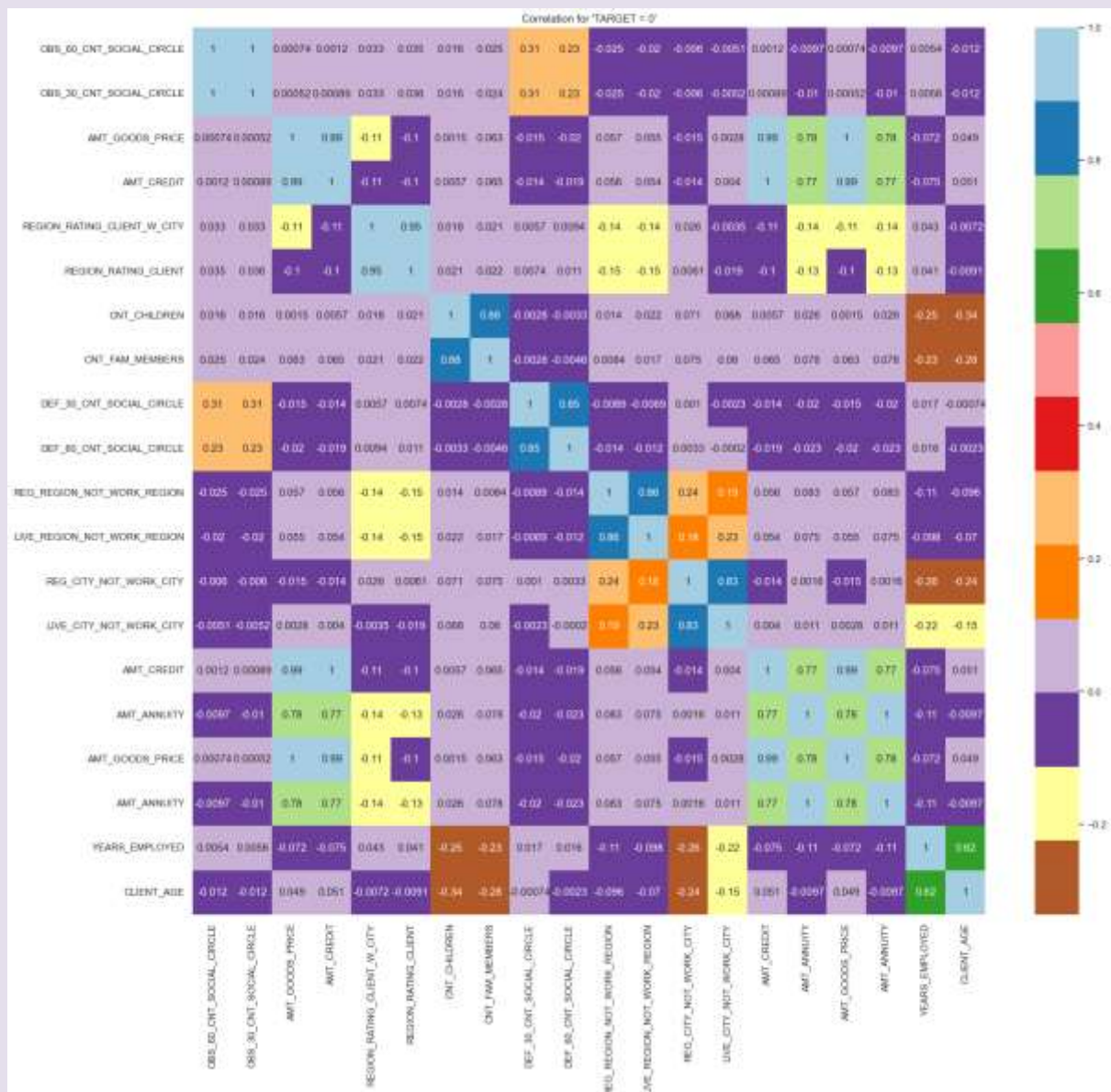Distribution of Credit Amounts by Contract Type (Credit Amount <= $725,067)

# Correlation Analysis:

Correlation analysis for Target 1


Correlation for 'TARGET = 1'

Correlation analysis for Target 0



Link of complete Python workbook:
https://drive.google.com/file/d/1BeSueJIIeSqkeTRgK0l28DjjZqMpIhCx/view?usp=drive_link

Video Link:
https://drive.google.com/file/d/1Aq_ILYDxTbOUepryI4ibx6ckQehAppD8/view?usp=drive_link