

# BANK LOAN CASE STUDY

-MRUNALI PETAKAR



## **DESCRIPTION**

The project aims to analyze the loan application dataset to identify patterns and factors that influence loan default. The goal is to use Exploratory Data Analysis (EDA) to understand how customer attributes and loan attributes impact the likelihood of default, so that the company can make informed decisions and reduce the risk of default. The dataset contains information about loan applications, including customers with payment difficulties and all other cases. The project will focus on identifying missing data, outliers, data imbalance, and performing univariate, segmented univariate, and bivariate analysis to gain insights into the driving factors of loan default. The top correlations for different scenarios will also be identified.

# SCOPE

The project will focus on analyzing the loan application dataset provided, which contains information about loan applicants and their loan applications. The dataset will be analyzed to identify patterns and trends that can help predict the likelihood of loan default.

## TECH - STACK USED

- Python 3.11
- Libraries: Pandas, Matplotlib, Seaborn, NumPy
- Google Drive
- PowerPoint

here I am giving the link to Excel sheet & github link

# APPROACH

## **Descriptive Statistics and Visualization**

to detect and  
identify  
outliers.

## **Model Building and Evaluation**

to evaluate the  
performance of the  
classification  
models.

## **Data Preparation and Cleaning**

Removing any irrelevant or  
redundant data to ensure  
that the remaining variables  
are relevant to the analysis.

## **Bivariate and Multivariate Analysis**

to identify any  
interactions or non-linear  
relationships between  
variables.

## **Insights and Recommendations**

to create visualizations and  
present the findings and  
recommendations in a clear  
and concise manner.





# UNDERSTANDING THE DATA

- a. `previous_application.csv`: Contains information about previous loan applications.
- b. `application_data.csv`: Provides details about the current loan applications.
- c. `columns_description.csv`: Describes the columns present in the other datasets, explaining what each column represents.

## Approach-

- Loaded three datasets: `application_data`, `previous_application`, and `columns_description`.
- Examined the structure and content of the datasets

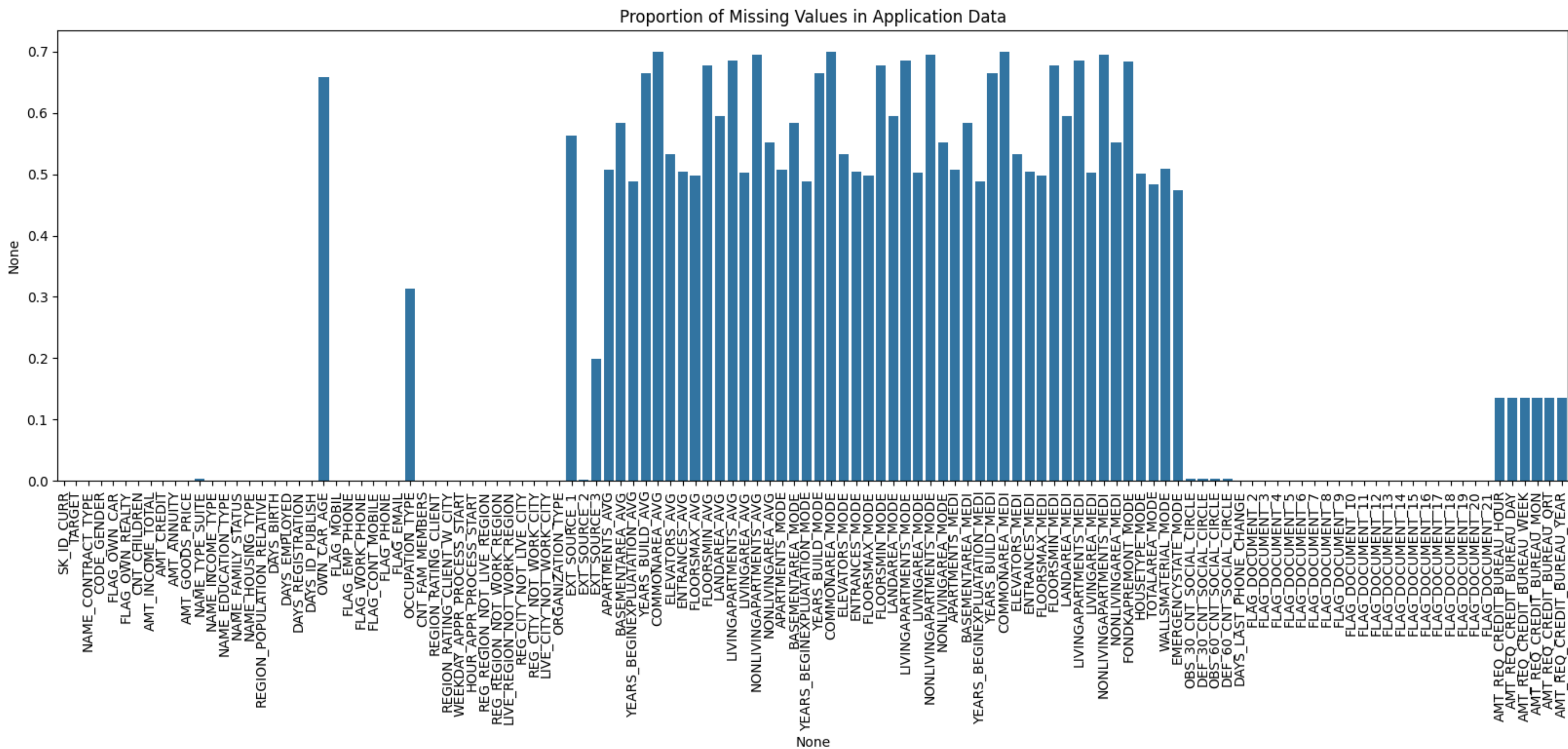
# 1] DATA CLEANING

1. Identified missing values in application\_data.
2. Visualized the proportion of missing values using a bar plot.
3. Handled missing values:
  - Imputed missing numerical values with the mean.
  - Imputed missing categorical values with the mode.

columns to drop -

HOUSETYPE\_MODE  
WALLSMATERIAL\_MODE  
BASEMENTAREA\_MEDI  
FLOORSMIN\_MEDI  
LIVINGAREA\_AVG  
ELEVATORS\_AVG  
LANDAREA\_AVG  
LIVINGAPARTMENTS\_AVG  
LIVINGAREA\_MODE  
ELEVATORS\_MODE  
LANDAREA\_MODE  
LIVINGAPARTMENTS\_MODE  
LIVINGAREA\_MEDI  
ELEVATORS\_MEDI  
LANDAREA\_MEDI  
LIVINGAPARTMENTS\_MEDI  
ENTRANCES\_AVG  
NONLIVINGAREA\_AVG  
OWN\_CAR\_AGE  
FONDKAPREMONT\_MODE  
ENTRANCES\_MODE  
NONLIVINGAREA\_MODE  
YEARS\_BUILD\_AVG  
NONLIVINGAPARTMENTS\_AVG  
ENTRANCES\_MEDI  
NONLIVINGAREA\_MEDI  
YEARS\_BUILD\_MODE  
NONLIVINGAPARTMENTS\_MODE  
APARTMENTS\_AVG  
EXT\_SOURCE\_1  
YEARS\_BUILD\_MEDI  
NONLIVINGAPARTMENTS\_MEDI  
APARTMENTS\_MODE  
BASEMENTAREA\_AVG  
FLOORSMIN\_AVG  
COMMONAREA\_AVG  
APARTMENTS\_MEDI  
BASEMENTAREA\_MODE  
FLOORSMIN\_MODE  
COMMONAREA\_MODE  
COMMONAREA\_MEDI

# data cleaning



```
# Identify numerical columns
```

```
numerical_columns = application_data.select_dtypes(include=['float64', 'int64']).columns
```

```
# Replace missing values with the mean for numerical columns
```

```
application_data[numerical_columns] =
```

```
application_data[numerical_columns].fillna(application_data[numerical_columns].mean())
```

```
# Identify categorical columns
```

```
categorical_columns = application_data.select_dtypes(include=['object']).columns
```

```
# Replace missing values with the median for categorical columns
```

```
application_data[categorical_columns] =
```

```
application_data[categorical_columns].fillna(application_data[categorical_columns].mode().iloc[0])
```

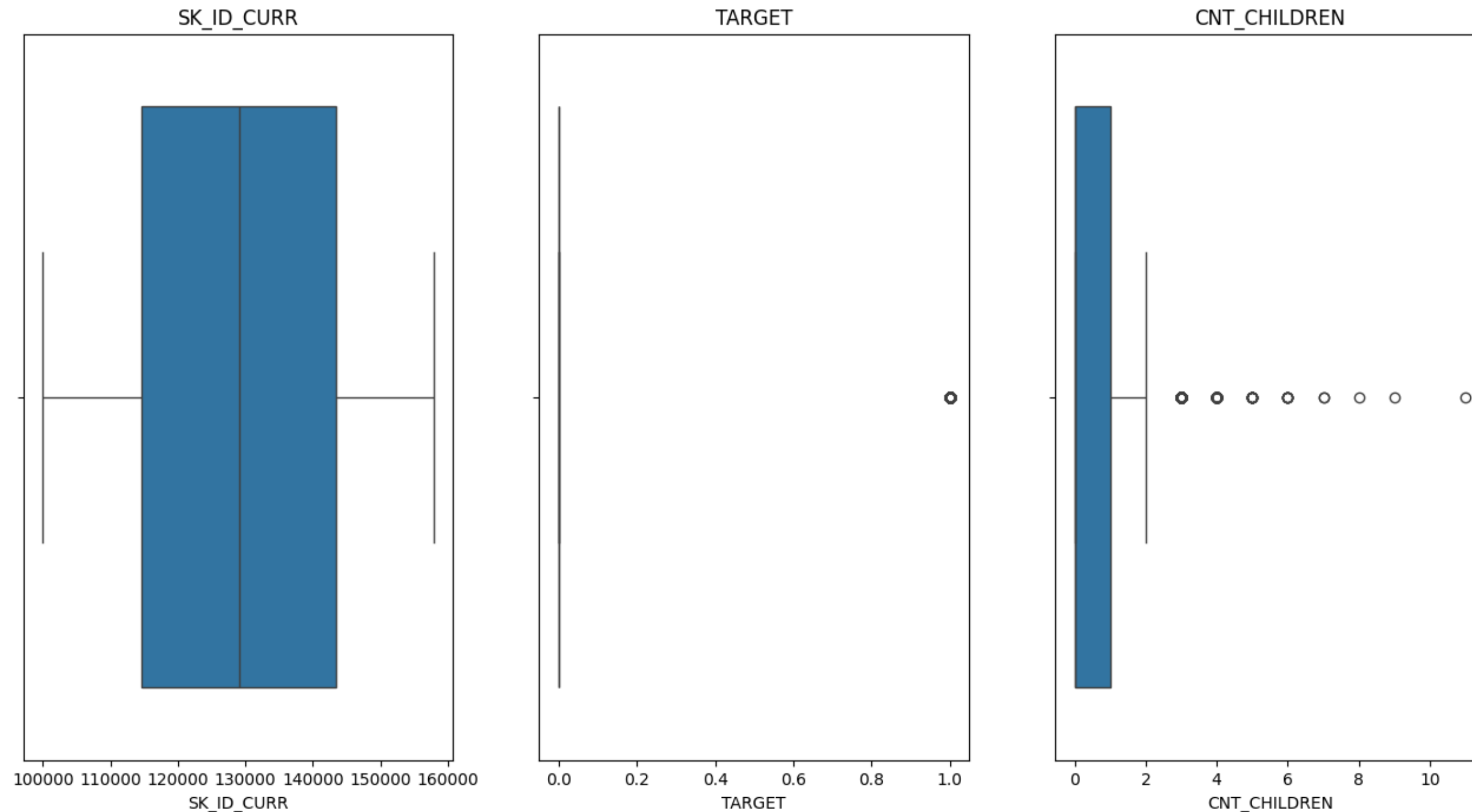


# TASK 2]

Detect and identify outliers in the dataset using Excel statistical functions and features, focusing on numerical variables.

- Outliers can only be identified on Numeric variables.

Identified outliers in numerical columns of application\_data using box plots shown as below



# TASK C] Analyze Data Imbalance

- Determining if there is data imbalance in the loan application dataset and calculating the ratio of data imbalance  
Visualized the distribution of the target variable (TARGET) using a pie chart and bar chart.

```
# Assuming 'TARGET' is your target variable column
```

```
target_counts = application_data['TARGET'].value_counts()
```

```
# Pie Chart
```

```
plt.figure(figsize=(4,4))
```

```
plt.pie(target_counts, labels=target_counts.index, autopct='%1.1f%%', startangle=90, colors=['skyblue', 'lightcoral'])
```

```
plt.title('Distribution of Target Variable (TARGET)')
```

```
plt.show()
```

```
# Bar Chart
```

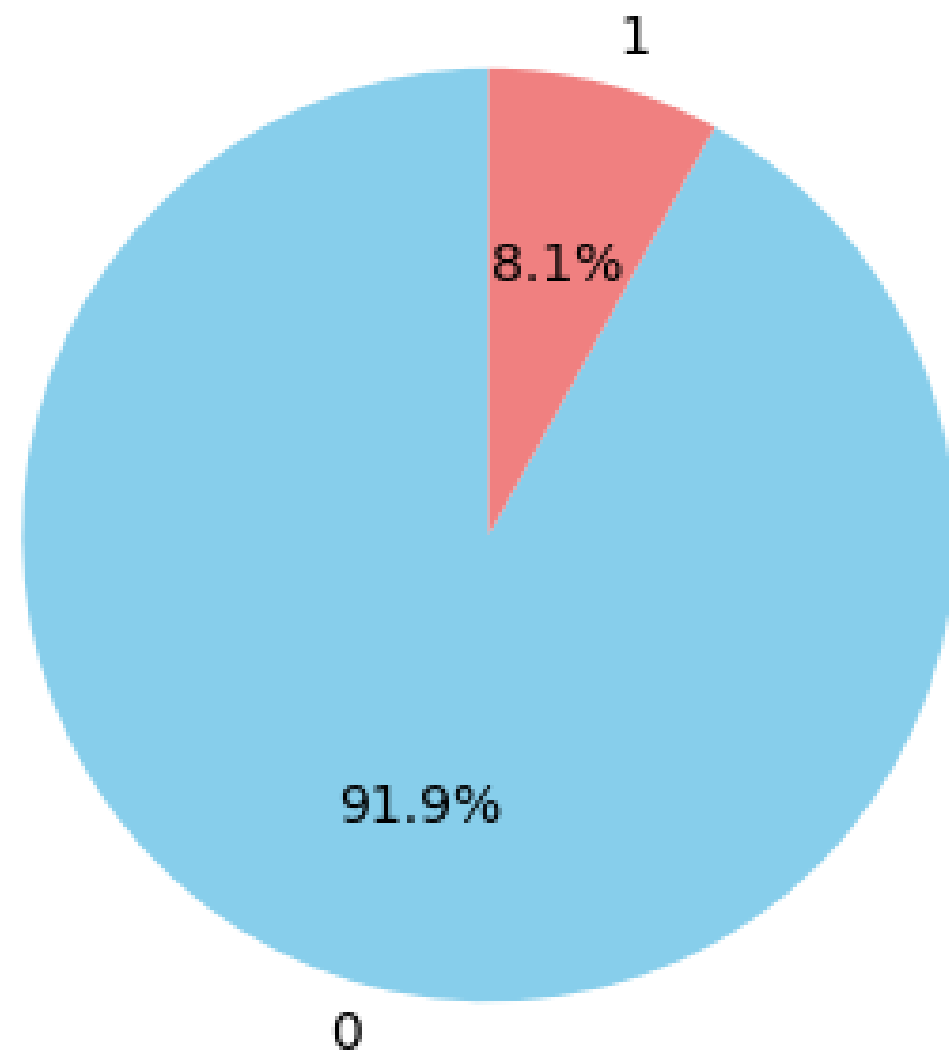
```
plt.figure(figsize=(6, 6))
```

```
sns.countplot(x='TARGET', data=application_data, palette='pastel')
```

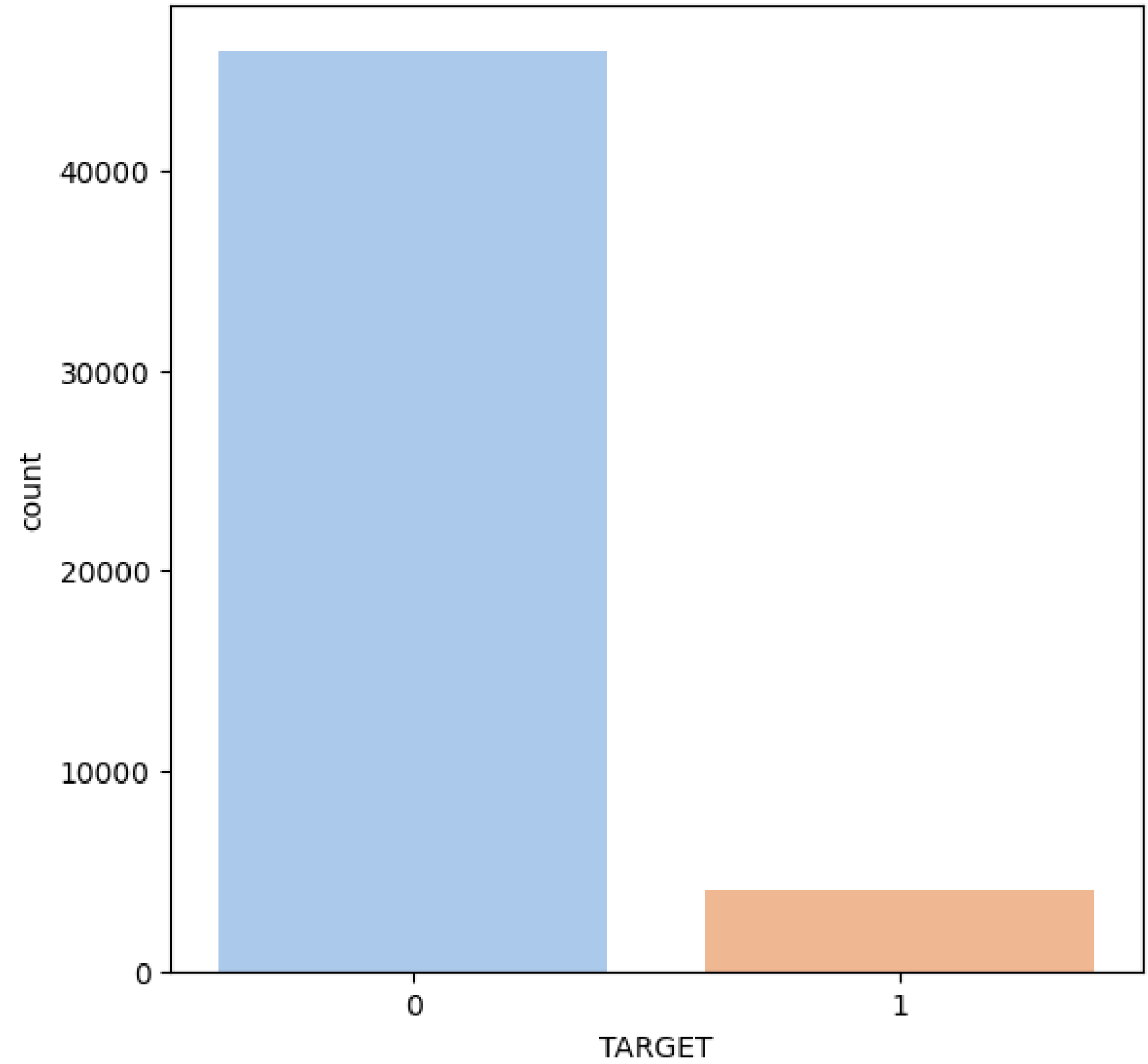
```
plt.title('Distribution of Target Variable (TARGET)')
```

```
plt.show()
```

Distribution of Target Variable (TARGET)



Distribution of Target Variable (TARGET)



## **TASK D] Perform Univariate, Segmented Univariate, and Bivariate Analysis:**

### Univariate Analysis:

Plotted histograms for numerical variables.

Plotted count plots for categorical variables.

### Segmented Univariate Analysis:

Utilized grouped bar charts for categorical variables across scenarios.

### Bivariate Analysis:

Created scatter plots for numerical-target variable relationships.

Utilized box plots for analyzing the income distribution by education and target.

Performing univ# Example for a numerical variable 'income'

```
plt.figure(figsize=(8, 6))
```

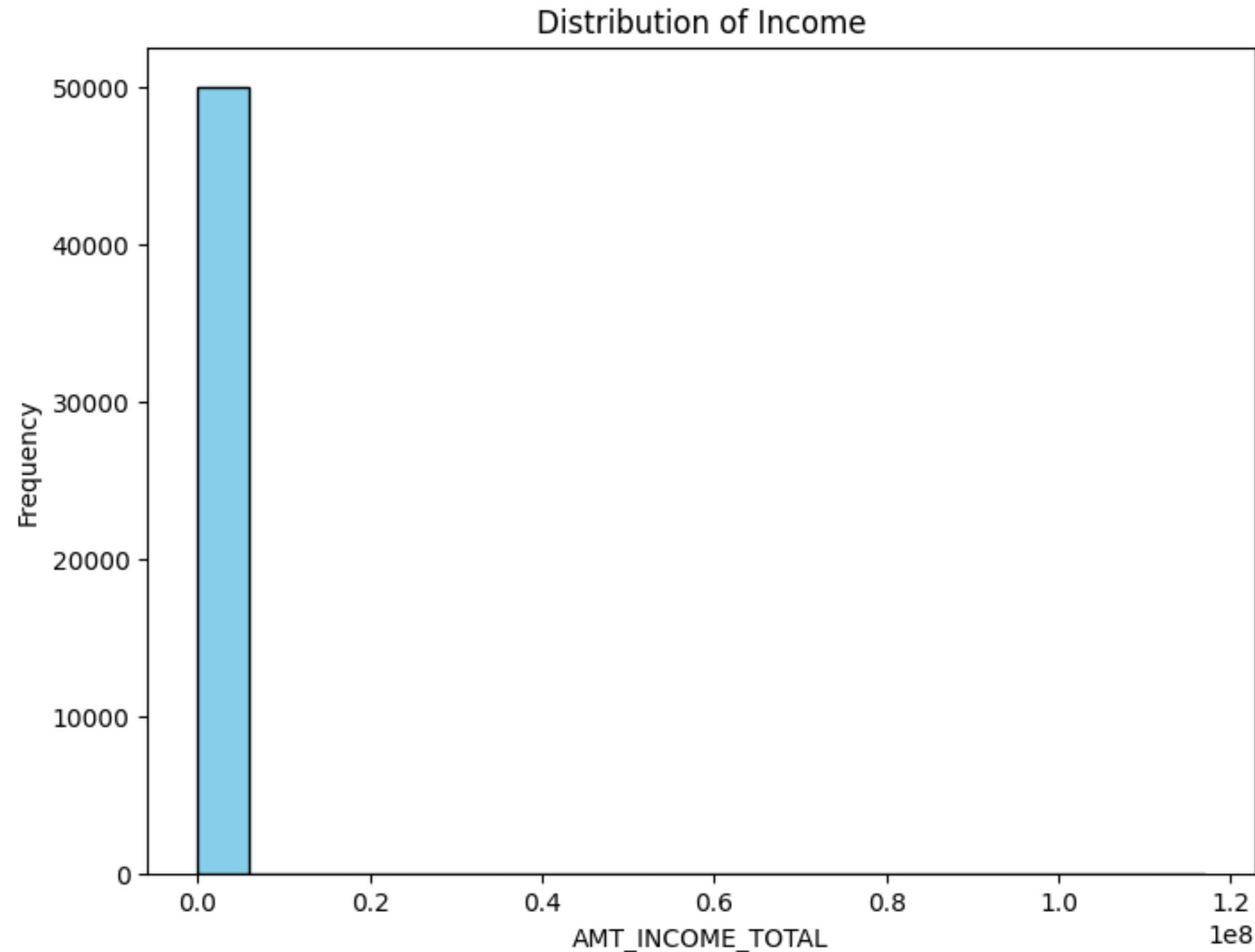
```
plt.hist(application_data['AMT_INCOME_TOTAL'], bins=20, color='skyblue', edgecolor='black')
```

```
plt.title('Distribution of Income')
```

```
plt.xlabel('AMT_INCOME_TOTAL')
```

```
plt.ylabel('Frequency')
```

```
plt.show()
```





## #Histograms for Numerical Variables:

```
# Selecting numerical columns
```

```
numerical_columns =  
application_data.select_dtypes(include=  
['int64', 'float64']).columns
```

```
# Plotting histograms for all numerical  
columns
```

```
plt.figure(figsize=(16, 12))
```

```
for i, col in enumerate(numerical_columns,  
1):
```

```
    plt.subplot(1, 3, i)
```

```
    plt.hist(application_data[col], bins=20,  
color='skyblue', edgecolor='black')
```

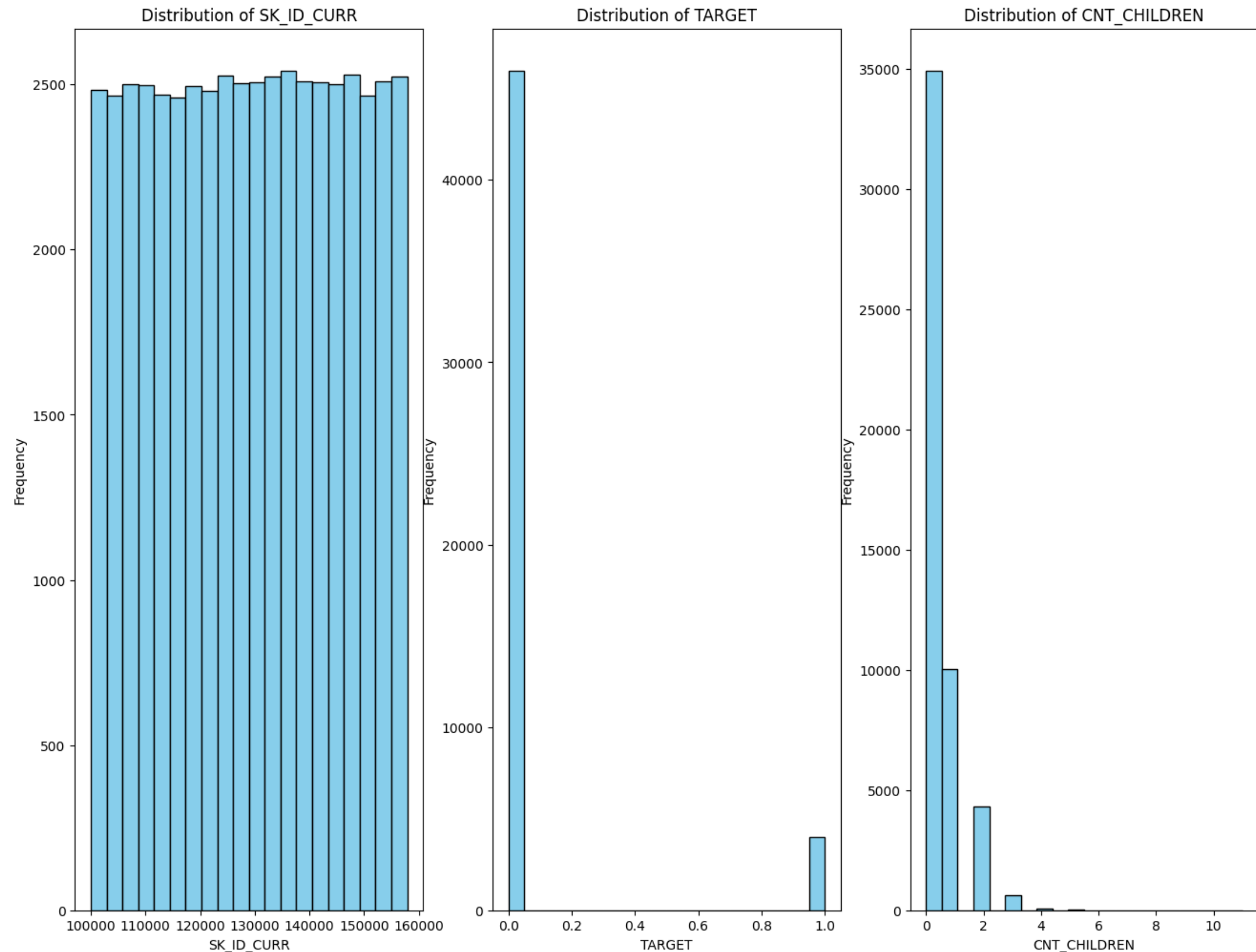
```
    plt.title(f'Distribution of {col}')
```

```
    plt.xlabel(col)
```

```
    plt.ylabel('Frequency')
```

```
plt.tight_layout()
```

```
plt.show()
```



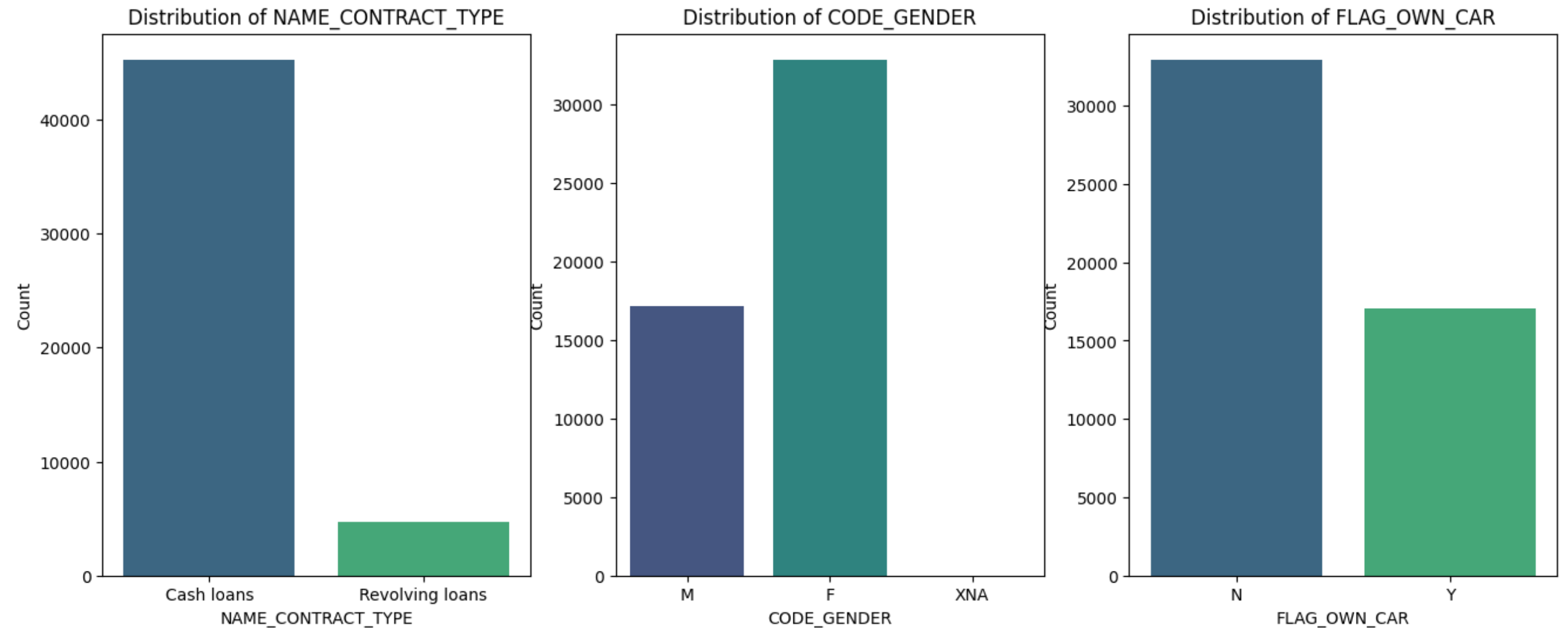
#Bar Charts for Categorical Variables:

```
import matplotlib.pyplot as plt
import seaborn as sns
```

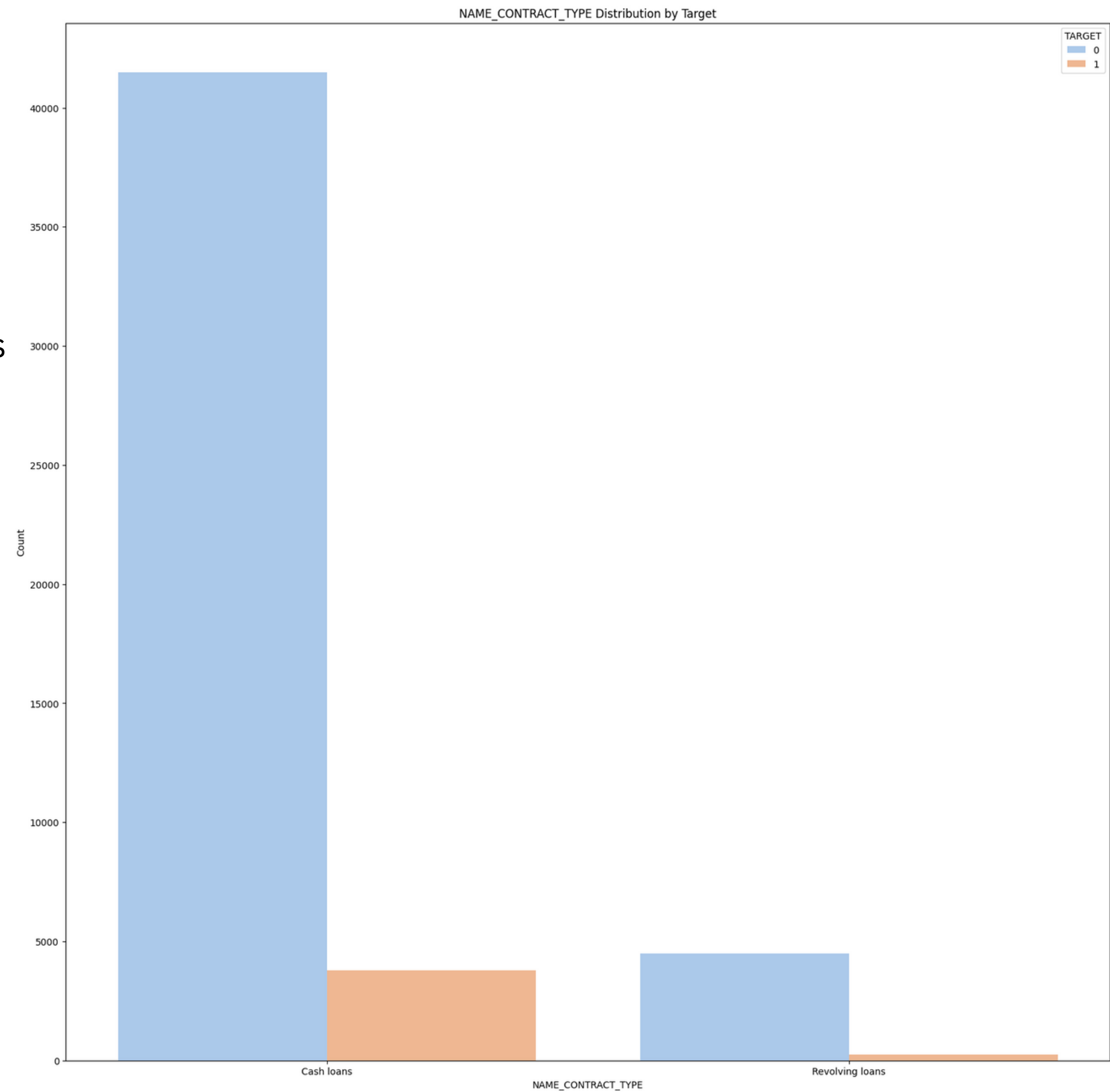
```
# Selecting categorical columns
categorical_columns =
application_data.select_dtypes(include=['object']).columns
```

```
# Plotting count plots for all
categorical columns
plt.figure(figsize=(16, 6))
for i, col in
enumerate(categorical_columns, 1):
    plt.subplot(1, 3, i)
    sns.countplot(x=col,
data=application_data,
palette='viridis')
    plt.title(f'Distribution of {col}')
    plt.xlabel(col)
    plt.ylabel('Count')
```

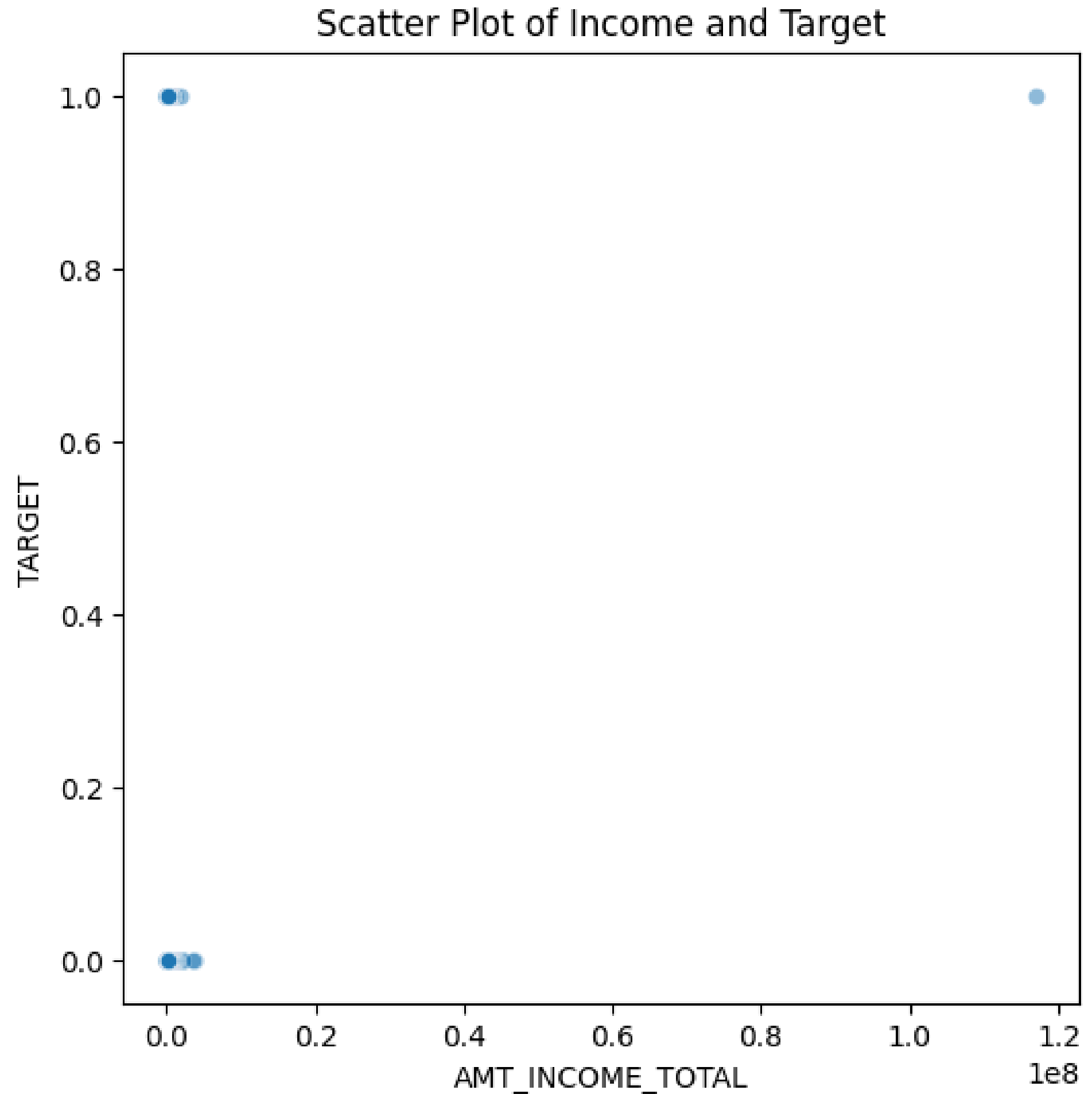
```
plt.tight_layout()
plt.show()
```



```
#Segmented Univariate Analysis:
#Grouped Bar Charts for Categorical Variables Across
Scenarios:
plt.figure(figsize=(20, 20))
sns.countplot(x='NAME_CONTRACT_TYPE',
hue='TARGET', data=application_data,
palette='pastel')
plt.title('NAME_CONTRACT_TYPE Distribution by
Target')
plt.xlabel('NAME_CONTRACT_TYPE')
plt.ylabel('Count')
plt.show()
```



```
#Bivariate Analysis:  
#Scatter Plots for Numerical-  
Target Variable Relationships:  
plt.figure(figsize=(6,6))  
sns.scatterplot(x='AMT_INCOME_  
TOTAL', y='TARGET',  
data=application_data,  
alpha=0.5)  
plt.title('Scatter Plot of Income  
and Target')  
plt.xlabel('AMT_INCOME_TOTAL')  
plt.ylabel('TARGET')  
plt.show()
```



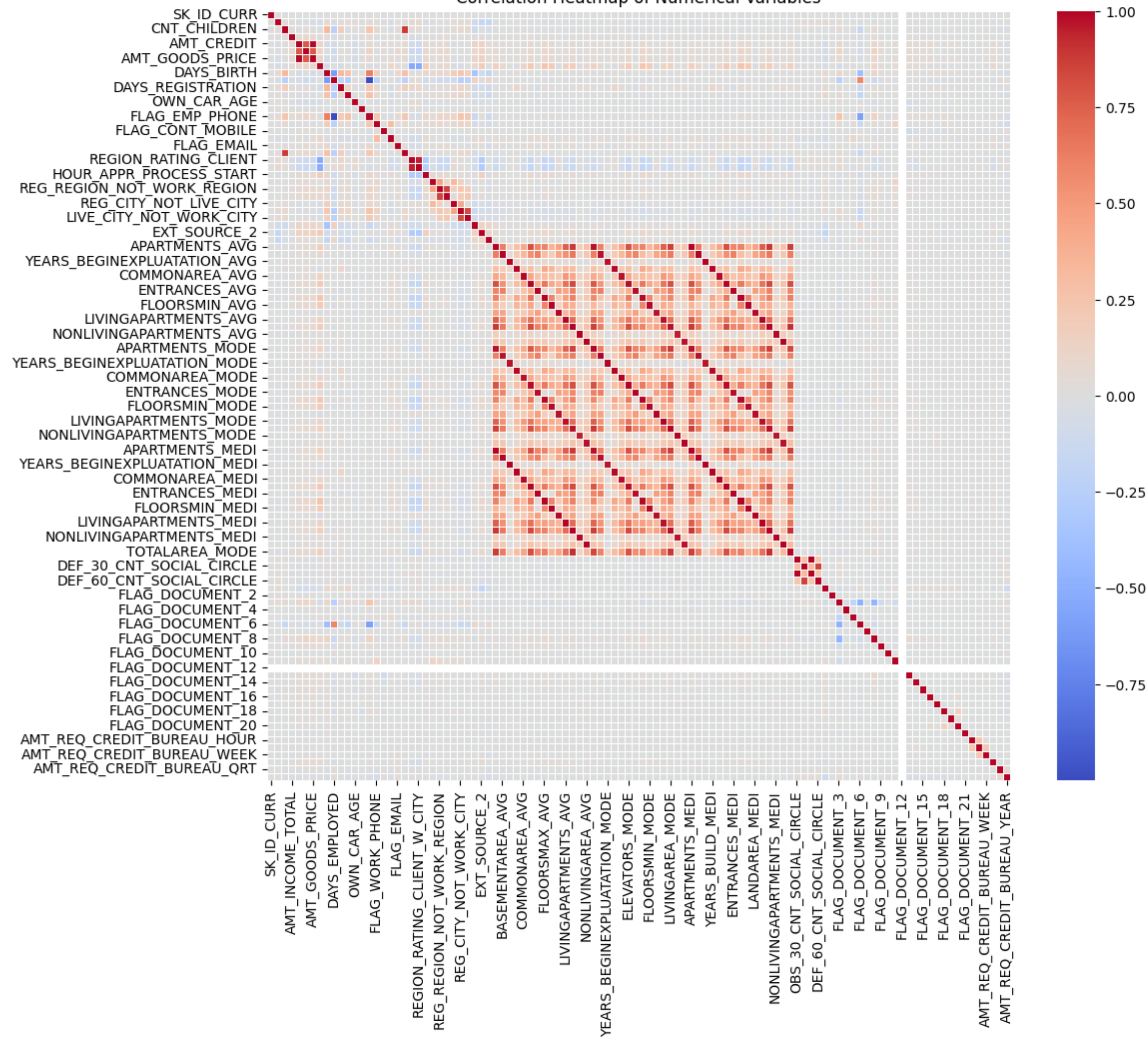
## Task E]

Identify Top Correlations for Different Scenarios: Understanding the correlation between variables and the target variable can provide insights into strong indicators of loan default. Created a function (analyze\_correlations) to analyze and visualize correlation matrices for different segments.

CORRELATION FOR APPLICANTS WITH PAYMENT MADE ON TIME								
CNT_CHILDREN	1	0.027	0.003	-0.024	-0.337	-0.245	0.029	0.023
AMT_INCOME_TOTAL	0.027	1	0.343	0.168	-0.063	-0.140	-0.023	-0.187
AMT_CREDIT	0.003	0.343	1	0.101	0.047	-0.070	0.001	-0.103
REGION_POPULATION_RELATIVE	-0.024	0.168	0.101	1	0.025	-0.007	0.001	-0.539
DAYS_BIRTH(Years)	-0.337	-0.063	0.047	0.025	1	0.626	0.271	-0.002
DAYS_EMPLOYED (Years)	-0.245	-0.140	-0.070	-0.007	0.626	1	0.277	0.038
DAYS_ID_PUBLISH(Years)	0.029	-0.023	0.001	0.001	0.271	0.277	1	0.009
REGION_RATING_CLIENT	0.023	-0.187	-0.103	-0.539	-0.002	0.038	0.009	1
	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	REGION_POPULATION_RELATIVE	DAYS_BIRTH(Years)	DAYS_EMPLOYED (Years)	DAYS_ID_PUBLISH(Years)	REGION_RATING_CLIENT



Correlation Heatmap of Numerical Variables



```
application_data = pd.read_csv('application_data.csv')

def analyze_correlations(data, target_column, top_n=5):
    # Selecting only numeric columns for correlation analysis
    numeric_data = data.select_dtypes(include=[np.number])

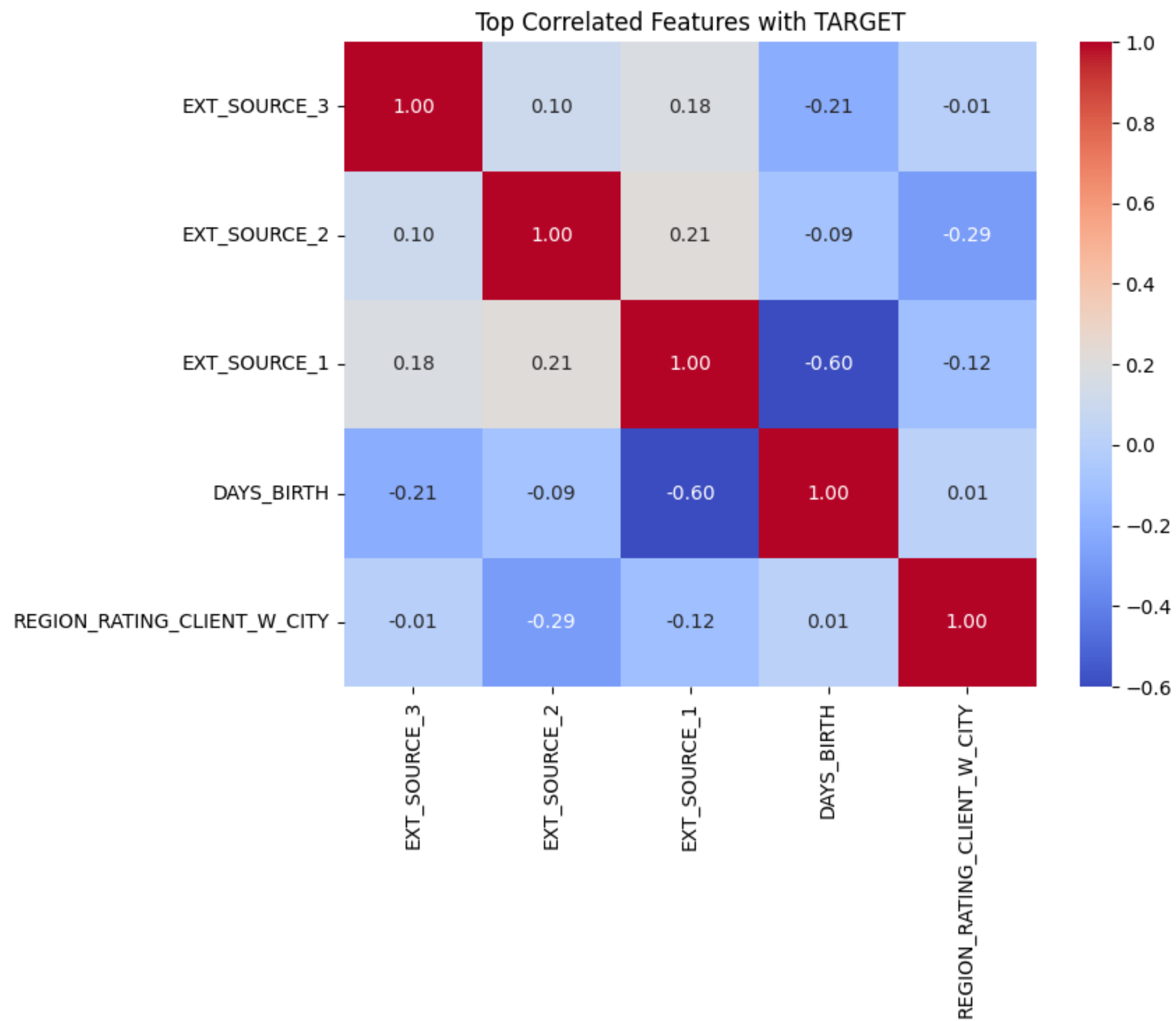
    # Calculating correlation matrix
    correlation_matrix = numeric_data.corr()

    # Getting correlations with the target variable
    target_correlations = correlation_matrix[target_column].drop(target_column)

    # Getting the top correlated features
    top_correlations = target_correlations.abs().sort_values(ascending=False).head(top_n)

    # Plotting heatmap
    plt.figure(figsize=(8, 6))
    sns.heatmap(correlation_matrix.loc[top_correlations.index, top_correlations.index], annot=True, cmap='coolwarm', fmt=".2f")
    plt.title(f'Top Correlated Features with {target_column}')
    plt.show()

# Assuming 'TARGET' is the column indicating payment difficulties
analyze_correlations(application_data, target_column='TARGET', top_n=5)
```



*Through this project, I accomplished several key objectives:*

- Handling Missing Data:
  - Successfully addressed missing data by imputing values based on data types.
  - Achieved a dataset with minimal missing values.
- Outlier Identification:
  - Identified outliers in numerical columns, providing insights into potential data quality issues.
- Data Imbalance:
  - Explored the distribution of the target variable, revealing an imbalanced dataset.
- Exploratory Data Analysis:
  - Conducted a comprehensive EDA, revealing patterns and trends in the data.
  - Utilized various visualization techniques for effective data exploration.
- Correlation Analysis:
  - Examined correlations between numerical variables in different scenarios.
  - Identified top correlated features with the target variable.

Here are some key insights that I obtained from the Bank Loan Case Study project:

1. Income and Loan Requests:
  - Individuals with higher incomes are less likely to apply for loans.
  - The credit amount of a bank loan typically falls within the range of 45,000 to 1,045,000.
  - The majority of loan applications come from people between the ages of 35 and 50.
  - Those with 0 to 8 years of work experience are more likely to seek loans.
2. Homeownership and Marital Status:
  - Individuals who own homes are more likely to apply for loans.
  - Married individuals are more inclined to take out loans compared to singles or those with other marital statuses.
3. Employment Status:
  - Individuals with jobs are more likely to request loans.
  - Unaccompanied minors have requested additional loans.
4. Loan Outcomes:
  - Customers who live in low-rating areas are more likely to have loan defaults.
  - Individuals with lower incomes are more likely to default.
  - Younger individuals are more likely to default, with the trend of defaulters decreasing with age.
  - Females are less likely than males to have defaults.



- Family Size and Education:
  - Customers with more than five family members are more likely to default on their bank loans.
  - Customers with fewer educational qualifications are more likely to fail to repay their loans.
  - Clients with little work experience are more likely to have defaults.
- Loan Types and Approval Rates:
  - Consumer loans have a significantly lower rate of cancellations and the highest approval rate.
  - Loans requested for the first Selling Place Area group experienced a higher rate of cancellations.
- Previous Loan History:
  - Clients who have applied for previous loans tend to have no defaults in their current loans.
- Top Correlations with Loan Default:
  - The top factors correlated with loan default include income type, family size, children count, external source, region rating of the client, age, months employed, amount credit, amount goods price, and amount total income.

## Result

- The project successfully achieved its objectives, providing valuable insights into the bank loan dataset. By addressing missing data, identifying outliers, and conducting thorough exploratory data analysis, the team gained a better understanding of the dataset's characteristics. The correlation analysis helped in identifying features strongly correlated with the target variable.

These insights are valuable for the finance company to make informed decisions about loan approvals, risk assessments, and strategies to reduce loan defaults. They highlight the importance of considering various customer attributes and loan attributes when evaluating loan applications.