

2023 | By: Pappu D. Kapgate



Credit EDA Case Study

Unlock the power of EDA on Loan Application Assessment

Table of Contents

Problem Statement I	2
Business Understanding:.....	2
Business Objectives:.....	3
Data Understanding:.....	4
Problem Statement II	5
Analysis Result	8
Introduction	8
Processing Current Application	9
Data cleansing	9
Handling Outliers	12
Float Outliers	12
Integer outliers	13
Derived Columns.....	14
Performing Data Analysis and Identifying Insights	15
Categorical unordered univariate analysis.....	15
Bivariate Analysis	17
Analysis between Numeric and Categorical variables	19
Processing Previous Application	22
Performing Data Analysis and Identifying Insights	26
Merging Current and Previous Application	30

Problem Statement I

This case study is designed to illustrate the practical application of Exploratory Data Analysis (EDA) in a business context. Through the exploration of EDA techniques, participants will gain insights into risk analytics within the banking and financial services sector. The study emphasizes the utilization of data to mitigate the risk associated with lending, providing a comprehensive understanding of strategies employed to minimize financial losses in customer transactions.

Business Understanding:

Loan providers face challenges in extending loans to individuals with insufficient or non-existent credit histories, as some consumers exploit this situation by defaulting on loans. As an employee of a consumer finance company specializing in lending various types of loans to urban customers, your role involves utilizing Exploratory Data Analysis (EDA) to analyze patterns in the data. The objective is to ensure that loan applicants capable of repaying are not unjustly rejected.

When evaluating loan applications, the company must make decisions based on the applicant's profile, considering two associated risks:

- Approving a loan for a capable applicant ensures business for the company.
- Approving a loan for a high-risk applicant may lead to financial losses if the applicant defaults.

The provided data includes information about loan applications at the time of submission and encompasses two scenarios:

- **Clients with payment difficulties:** Those who had late payments exceeding X days on at least one of the first Y installments.
- **All other cases:** Instances where payments were made on time.

Four types of decisions can be made by the client or company during the loan application process:

1. **Approved:** The company approves the loan application.
2. **Cancelled:** The client cancels the application, either due to a change of mind or, in some cases, due to receiving unfavorable terms.
3. **Refused:** The company rejects the loan application based on non-compliance with requirements.
4. **Unused offer:** The client cancels the loan at various stages of the process.

This case study employs EDA to discern how consumer attributes and loan characteristics influence the likelihood of default.

Business Objectives:

The primary goal of this case study is to identify patterns that serve as indicators of a client's difficulty in paying installments. The insights derived from these patterns can inform strategic actions, such as denying the loan, reducing loan amounts, or lending to riskier applicants at higher interest rates. The overarching objective is to ensure that consumers capable of repaying the loan are not unfairly rejected. The focus of this case study is on using Exploratory Data Analysis (EDA) to identify such applicants.

In essence, the company aims to comprehend the driving factors, or driver variables, behind loan defaults—those variables that strongly indicate a likelihood of default. The knowledge gained from this analysis can be leveraged for portfolio management and risk assessment within the company.

As part of developing a comprehensive understanding of the domain, participants are encouraged to conduct independent research on risk analytics. This research should encompass an exploration of the types of variables involved and their significance in the context of risk assessment.

Data Understanding:

The dataset for this analysis comprises three files:

- **application_data.csv**: This file encompasses comprehensive information about the client at the time of loan application. It specifically focuses on whether the client faces payment difficulties.
- **previous_application.csv**: This file provides insights into the client's previous loan data, detailing whether the prior application was approved, cancelled, refused, or if it resulted in an unused offer.
- **columns_description.csv**: Serving as a data dictionary, this file elucidates the meaning of the variables present in the dataset, offering a reference guide for understanding the information contained in the other files.

Problem Statement II

In the pursuit of extracting meaningful insights from our dataset, the anticipated results encapsulate a comprehensive understanding of various facets crucial to our analysis. As we embark on this exploration, we aim to present a clear overview of our analytical approach, addressing key elements such as handling missing data, identifying outliers, and analyzing potential data imbalances.

Our objective is to not only provide a structured presentation of our findings but to also translate these results into actionable business insights. Through a meticulous examination of the data, we anticipate unveiling patterns, correlations, and anomalies that will contribute to a robust comprehension of the underlying dynamics within the dataset.

This section will guide you through the expected outcomes of our analysis, encompassing the identification and resolution of missing data, insights into potential outliers, a nuanced understanding of data imbalances, and a comprehensive interpretation of analysis results in terms of their business implications. Each expected result is aligned with the overarching goal of extracting knowledge that is not only statistically significant but also strategically valuable in the context of our problem statement.

1. Overall Approach Presentation:

- a.** Provide a brief overview of the problem statement and the analysis approach.

2. Handling Missing Data:

- a.** Identify missing data in the dataset.
- b.** Utilize an appropriate method to address missing values, whether through removal of columns or replacement with suitable values.

- c. Clearly articulate the rationale behind the chosen approach.

3. Outlier Identification:

- a. Identify outliers in the dataset.
- b. Provide an explanation for why each identified point is considered an outlier.
- c. Emphasize that, for this exercise, removal of data points is not necessary.

4. Data Imbalance Analysis:

- a. Assess data imbalance in the 'Target variable' (clients with payment difficulties and all other cases).
- b. Calculate the ratio of data imbalance.
- c. Employ a mix of univariate and bivariate analysis, considering different scales for graphs if necessary.

5. Explaining Analysis Results:

- a. Interpret the results of univariate, segmented univariate, and bivariate analyses in business terms.
- b. Utilize loops for appropriate columns to extract insights.

6. Top 10 Correlations:

- a. Identify the top 10 correlations for clients with payment difficulties and all other cases.
- b. Segment the data frame with respect to the target variable.
- c. Correlate variables within each segment, excluding the target variable as it is categorical.

7. Visualizations and Summarization:

- a. Include relevant visualizations to explain numerical/categorical variables.

- b.** Summarize the most important results, providing insights into why specific variables are crucial for distinguishing clients with payment difficulties from all other cases.

Analysis Result

The analysis undertaken provides a comprehensive exploration into the intricacies of our dataset, delving into critical aspects that shape the landscape of consumer loan applications. Through a meticulous examination, we have unearthed valuable insights that not only address missing data but also identify outliers and assess the balance within our dataset.

Our journey begins with a strategic approach to handling missing data, ensuring a robust foundation for subsequent analyses. The identification and rationale behind addressing outliers are presented, with a notable emphasis on the fact that their removal is not imperative for the scope of this exercise.

Furthermore, our exploration extends to understanding the balance within our data, particularly focusing on the 'Target variable.' A nuanced analysis, blending univariate and bivariate perspectives, offers a profound understanding of the dynamics at play.

This introduction sets the stage for a detailed presentation of our analysis results, offering a glimpse into the multifaceted insights gleaned from our data exploration journey.

Introduction

As part of this case study exercise we were given two datasets of a bank namely

1. application_data.csv
2. previous_application_data.csv

The objective for our analysis was to understand this dataset, clean and identify patterns and behaviors which can help determine risk in banking sector while lending money to customers. Key factors that needed to be considered were

1. Lending money to customers that are incapable of paying back
2. Not lending money to customer who can pay back resulting in revenue loss.

Processing Current Application

We started looking at the application data csv file and explored data inconsistencies and null value columns.

1. Exploring the Data in Application Data file

```
[1]: appl_data.head()

[4]: #Checking number of rows and columns in the file
      appl_data.shape

[4]: (307511, 122)

[5]: #Check the column data size and data types
      appl_data.info()

      #This is a lot of data!!!
      <class 'pandas.core.frame.DataFrame'>
      RangeIndex: 307511 entries, 0 to 307510
      Columns: 122 entries, SK_ID_CURR to AMT_REQ_CREDIT_BUREAU_YEAR
      dtypes: float64(65), int64(41), object(16)
      memory usage: 286.2+ MB

[6]: #checking statistical data
      appl_data.describe()
```

We have cleared output to save space

Data cleansing

1. Before analyzing the data we identified percentage of null values in data columns.

```
[10]: # Identify the % of missing values and list greater than 25% null value columns
      null_cols=((appl_data.isnull().sum()*100)/appl_data.shape[0]).round(2)
      null_cols

[10]: SK_ID_CURR          0.0
      TARGET            0.0
      NAME_CONTRACT_TYPE 0.0
      CODE_GENDER        0.0
      FLAG_OWN_CAR       0.0
      ...
      AMT_REQ_CREDIT_BUREAU_DAY 13.5
      AMT_REQ_CREDIT_BUREAU_WEEK 13.5
      AMT_REQ_CREDIT_BUREAU_MON 13.5
      AMT_REQ_CREDIT_BUREAU_QRT 13.5
      AMT_REQ_CREDIT_BUREAU_YEAR 13.5
      Length: 122, dtype: float64
```

2. As there were a lot of columns having higher percentage of missing values. We decided to drop any column which has more than 25% of null values.

```
[9]: #Drop the columns
appl_data.drop(appl_data.loc[:,appl_data.isnull().mean()>=.25],axis=1,inplace=True)

[10]: #Once dropped the total number of columns are 72
appl_data.shape

[10]: (307511, 72)
```

3. Impute missing values with mean, median, mode, 0 where applicable.

Columns Name	Impute Strategy
AMT_GOODS_PRICE	Mean
NAME_TYPE_SUITE	Mode
CNT_FAM_MEMBERS	0
OBS_60_CNT_SOCIAL_CIRCLE	Mean
DEF_60_CNT_SOCIAL_CIRCLE	Mean
OBS_30_CNT_SOCIAL_CIRCLE	Mean
DEF_30_CNT_SOCIAL_CIRCLE	Mean
AMT_REQ_CREDIT_BUREAU_HOUR	0
AMT_REQ_CREDIT_BUREAU_DAY	0
AMT_REQ_CREDIT_BUREAU_WEEK	0
AMT_REQ_CREDIT_BUREAU_MON	0
AMT_REQ_CREDIT_BUREAU_QRT	0
AMT_REQ_CREDIT_BUREAU_YEAR	Median
AMT_ANNUITY	Mean
DAYS_LAST_PHONE_CHANGE	0

4. Checking and converting data type to appropriate type based on columns data.

a Object

b Float

We converted the below columns from float to integers as it had only whole numbers.

Counts and days are whole numbers in real life, so converting them into integers.

```
Memory usage: 171.51 Mb
[55]: #Check if any Objects need to be converted. Looks good!
appl_data.select_dtypes('object')
```

```
[56]: #Check if any float columns need to be converted. Looks like some can be converted into integers!
appl_data.select_dtypes('float')
```

Columns Name	Old Data type	New Data Type
CNT_FAM_MEMBERS	Float64	Int64
DAYS_LAST_PHONE_CHANGE	Float64	Int64
DAYS_REGISTRATION	Float64	Int64
AMT_REQ_CREDIT_BUREAU_HOUR	Float64	Int64
AMT_REQ_CREDIT_BUREAU_DAY	Float64	Int64
AMT_REQ_CREDIT_BUREAU_WEEK	Float64	Int64
AMT_REQ_CREDIT_BUREAU_MON	Float64	Int64
AMT_REQ_CREDIT_BUREAU_QRT	Float64	Int64
AMT_REQ_CREDIT_BUREAU_YEAR	Float64	Int64

5. Converting all flag variables values [1,0] to a Binary variable form (Y or N)

```
[67]: #converting all flag columns to Y or N instead of 1 and 0
appl_data[flag_columns]=appl_data[flag_columns].replace((0, 1), ('N', 'Y'))

[68]: #All the flag have been converted to Y and N
appl_data[list(appl_data.filter(regex='FLAG'))]

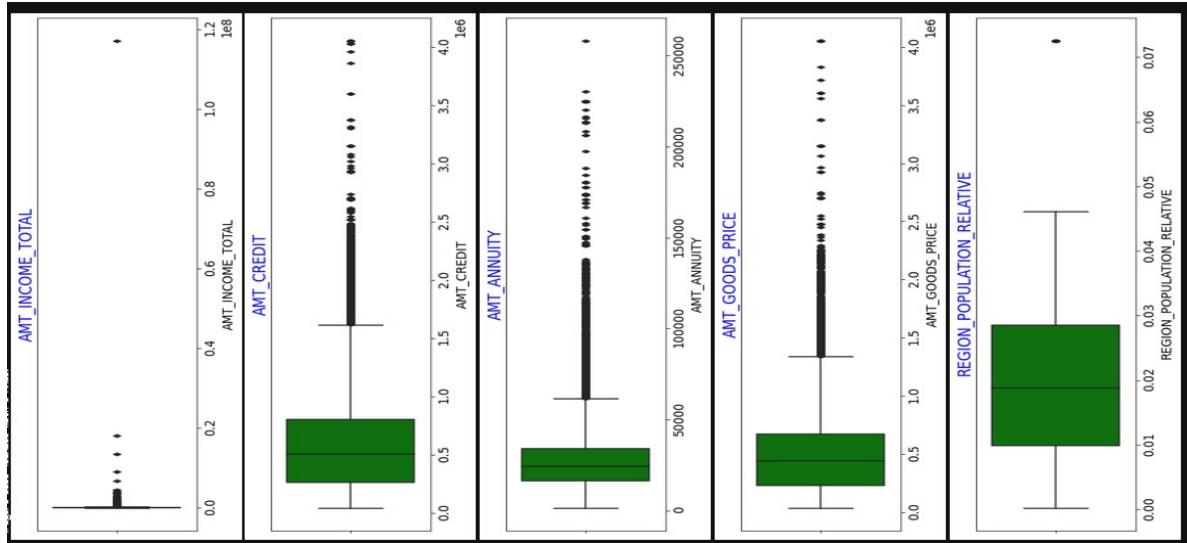
[68]:
   FLAG_own_car FLAG_own_realty FLAG_mobil FLAG_emp_phone FLAG_work_phone FLAG_cont_mobile
0           N             Y             Y             Y             N             Y
1           N             N             Y             Y             N             Y
2           Y             Y             Y             Y             Y             Y
3           N             Y             Y             Y             N             Y
4           N             Y             Y             Y             N             Y
...         ...
307506      N             N             Y             Y             N             Y
307507      N             Y             Y             N             N             Y
```

Handling Outliers

We identified outliers for both floating and integer columns.

Float Outliers

We can see the outliers depicted in the boxplot below are far away from the 75% quantile.

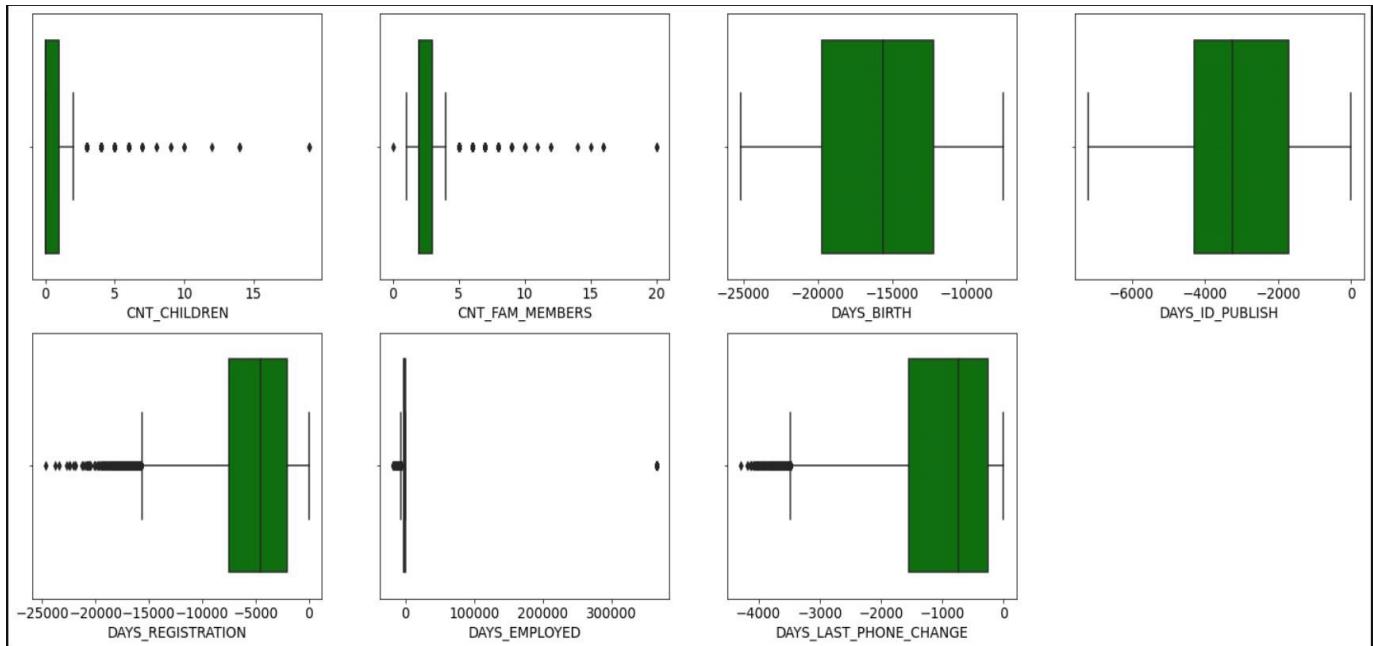


By looking at the data columns individually we decided on what percentile of data should be retained in the dataset. Which can make analysis accurate and does not result in skewed patterns. For determining the percentile (%) we looked at max values and retained only those outliers which are relatively closer to the 75% quantile values. We made sure all high value anomalies are completely removed.

Following is the list of columns and the corresponding quantiles below which we have retain the data.

Columns Name	Retain % data
AMT_INCOME_TOTAL	95%
AMT_CREDIT	99%
AMT_ANNUITY	99%
AMT_GOODS_PRICE	90%

Integer outliers



Along with outliers we observed that there were a huge population of negative values. Just like float attributes we handled the outliers by looking at the data and figuring out the anomalies.

Columns Name	Retain % data
CNT_CHILDREN	99%
CNT_FAM_MEMBERS	99%
DAYS_REGISTRATION	90%
DAYS_EMPLOYED	75%
DAYS_LAST_PHONE_CHANGE	90%

Derived Columns

We organized the applicable data columns into ranges and summarized it for effective analysis. • [Converting DAYS into YEARS](#)

```
: # Converting days into years

appl_data["YEARS_BIRTH"] = appl_data.DAYS_BIRTH.apply(lambda x: x/365)
appl_data["YEARS_ID_PUBLISH"] = appl_data.DAYS_ID_PUBLISH.apply(lambda x: x/365)
appl_data["YEARS_REGISTRATION"] = appl_data.DAYS_REGISTRATION.apply(lambda x: x/365)
appl_data["YEARS_LAST_PHONE_CHANGE"] = appl_data.DAYS_LAST_PHONE_CHANGE.apply(lambda x: x/365)
```

- [Categorical Variables](#)

[INCOME RANGE](#)

A) AMT_INCOME

```
8]: #creating categories for customer incomes.

label = ['Very Low', 'Low', 'Moderate', 'High', 'Very High']
appl_data["AMT_INCOME_CATEG"] = pd.qcut(appl_data.AMT_INCOME_TOTAL,q=[0, .2, .4, .6, .8, 1],labels=label)
appl_data.AMT_INCOME_CATEG.value_counts()

8]: Low      48663
Very Low   31909
High       28609
Very High  27430
Moderate   18632
Name: AMT_INCOME_CATEG, dtype: int64
```

- [AGE GROUP RANGE](#)

```
[109]: # Creating age groups of customer
bins = [18, 30, 40, 50, 60, 70, 120]
age_group_labels = ['18-29', '30-39', '40-49', '50-59', '60-69', '70+']
appl_data['AGE_GROUP_CUSTOMER'] = pd.cut(appl_data.AGE_CUSTOMER,bins,labels=age_group_labels,include_lowest=True)
appl_data['AGE_GROUP_CUSTOMER'].value_counts()

[109]: 30-39    57167
        18-29   39999
        40-49   38551
        50-59   17733
        60-69   1793
        70+      0
Name: AGE_GROUP_CUSTOMER, dtype: int64
```

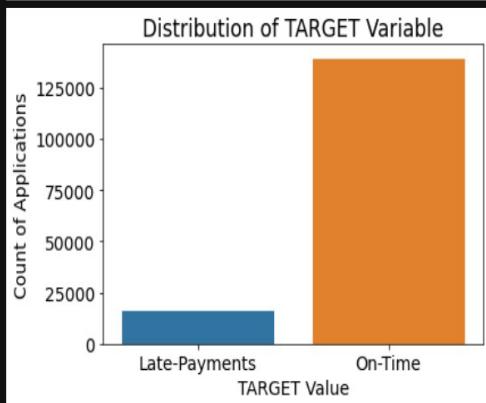
Performing Data Analysis and Identifying Insights

Categorical unordered univariate analysis

- **TARGET Variable**

- We created a new column and divided the dataset into on-time and late paying customers and assigned this category to paystatus column.

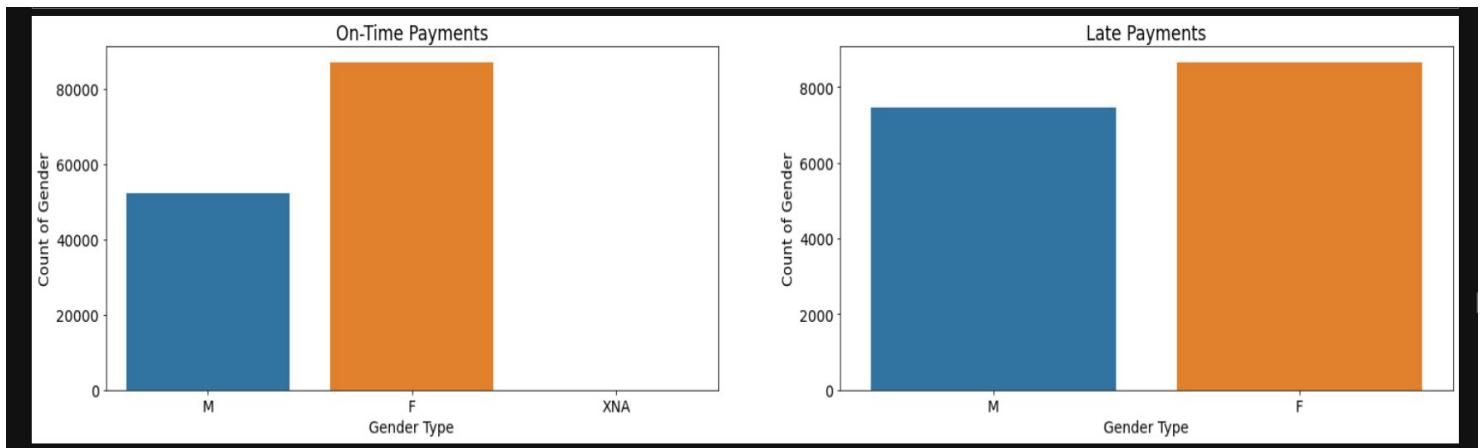
```
[112]: sns.countplot(appl_data.paystatus)
plt.rcParams.update({'font.size': 14})
plt.xlabel("TARGET Value")
plt.ylabel("Count of Applications")
plt.title("Distribution of TARGET Variable")
plt.show()
```



Observation: In the given population, the number of people having difficulties making payments is less than other applicants. The data been leaning towards ontime payments which may be due to random sampling.

⇒ Late payments are < 25K and others are >125K

- GENDER Variable



Observation: Females make more on-time payments than Male counterparts. But for late payments both genders are considerably similar.

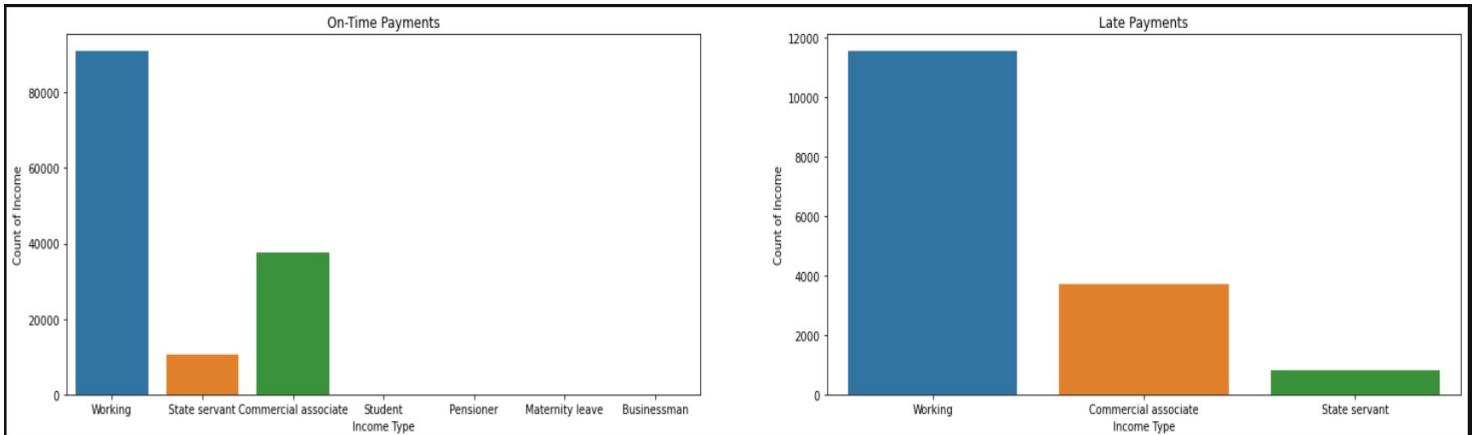
⇒ Ontime payments

- Males are 50-60K,
- Females are > 80K

⇒ Late payments

- Males are 7-8K,
- Female are >8K.

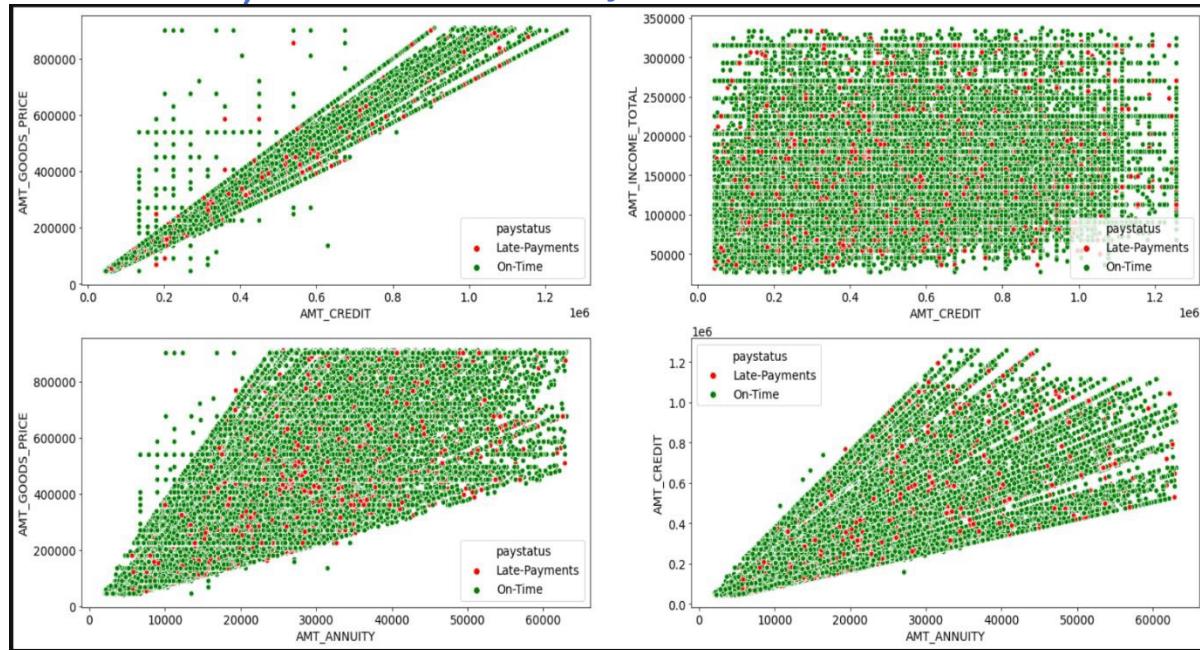
- INCOME TYPE Variable



Observation: Student, Pensioner, Maternity Leave and Businessman are not present in the late payments data group. Also On-time paying customers are higher than late paying customers accross all income type groups.

Bivariate Analysis

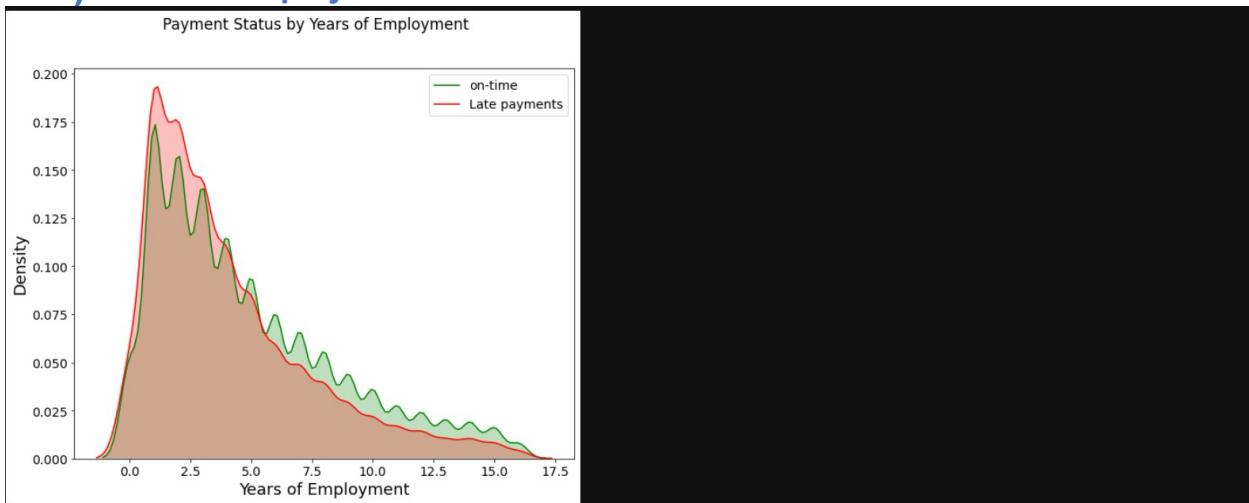
a) Numeric Variable Analysis



Observation:

- We can see a linear progression between credit amount, annuity amount and goods price. Even customers with payment difficulties have linearly progressed accross all the above plot.
- We didnt find strong linear corelation between income and credit amount. We were expecting to see people with higher income get more credit.

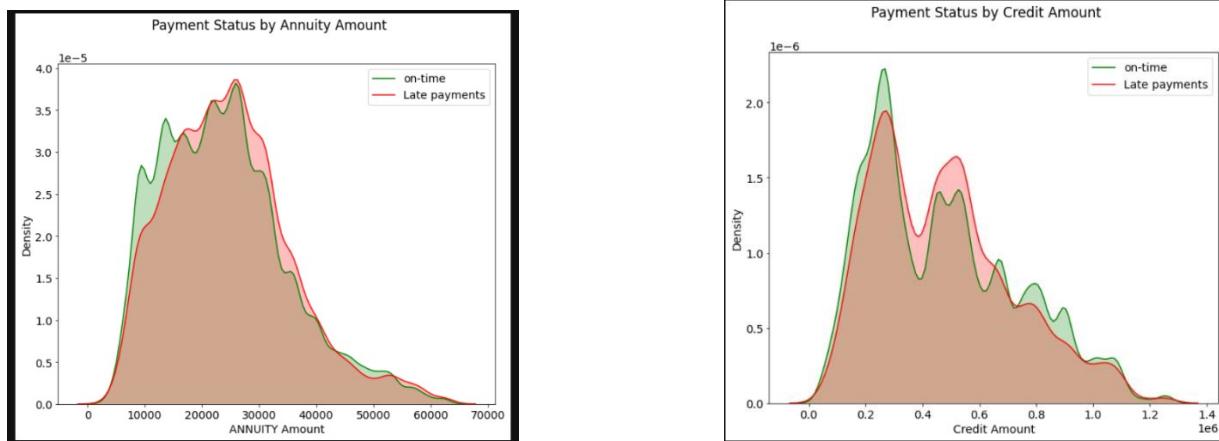
b) Years of Employment



Observation:

1. We observe that customers with less than 5 years of employment have a higher possibility to default on their loans.
2. Customers with greater than 5 years of employment have less difficulties in paying their loans, hence they are a better group to attract on loan products.

c) ANNUITY Amount and CREDIT Amount



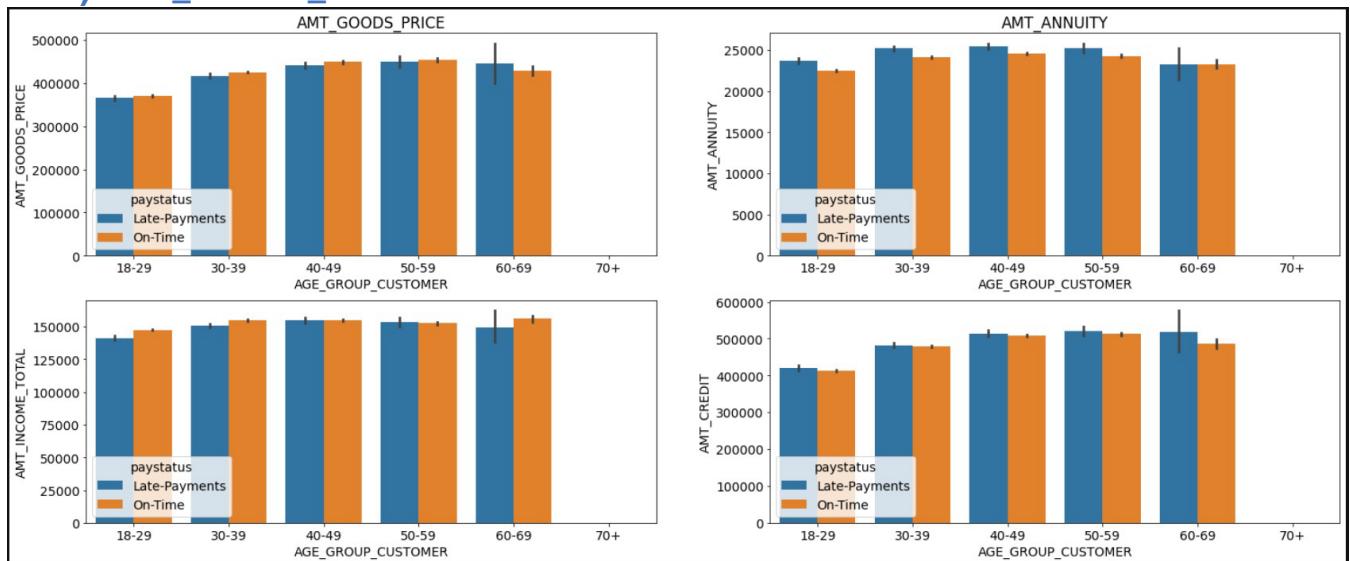
Observation:

- ⇒ We observe that annuity amount below 20,000 have higher on-time payments rates.
- ⇒ Customer with loan credits of 2.5 lakhs have a higher on-time payment than rest of the credit ranges.

⇒ So we can conclude there is a low risk in giving out loans for 2.5 Lakhs or repayment annuity of below 20K.

Analysis between Numeric and Categorical variables

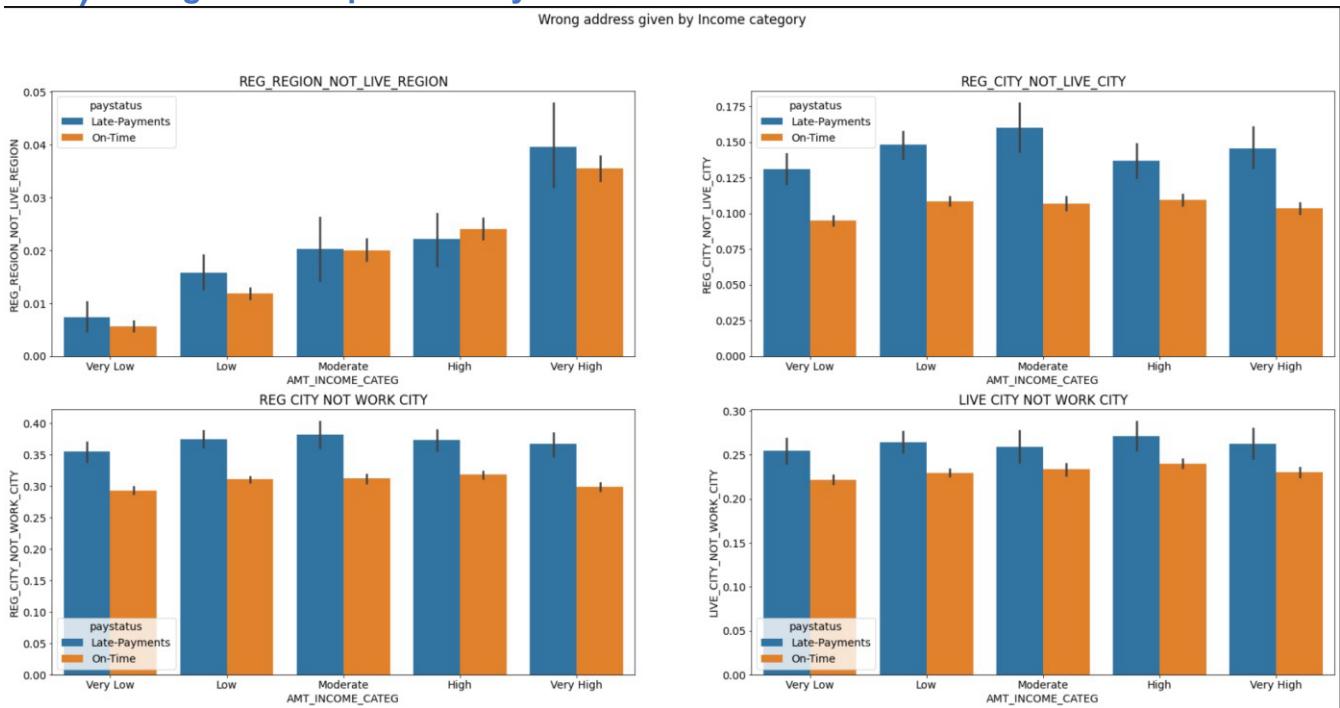
a) AGE_GROUP_CUSTOMER



Observation:

- ⇒ We observe within age group 60-69 if their total income is more than 1.25 Lakhs and credited amount is more than 5 Lakhs, then there may be some difficulty to repay. So, we can say that if we offer those customers loans below 5 Lakhs then there is a higher probability to pay installments on time.
- ⇒ While considering annuity in range 20-25 thousand, all age groups apart from 60-69 have a higher possibility to default.
- ⇒ So we can conclude there is a low risk in giving out loans to age group 60-69 if the annuity is around 20-25 thousand and income is more than 1.25 Lakhs and credited amount is below 5 lakhs.

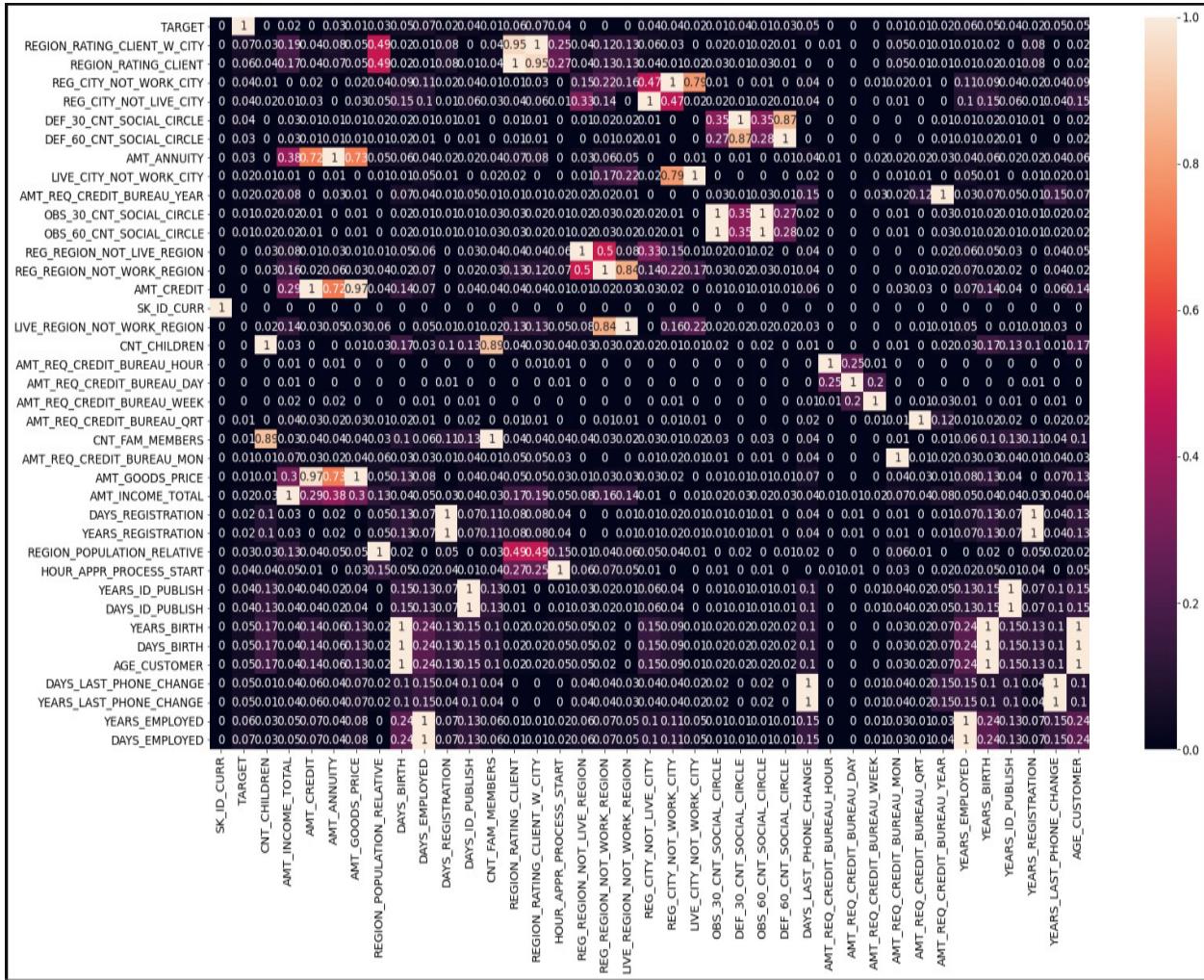
b) Wrong Address provided by Customers



Observation:

- ⇒ We observed that no matter what income category the customer belongs to, if they provide wrong contact address then they have a higher probability to default.
- ⇒ As all the bars indicate higher late payment values when the address is not provided.

C) Multivariate Analysis



Based on the above heatmap, we have identified the highly correlated attributes:

[351]:	AMT_GOODS_PRICE	AMT_CREDIT	0.97
	AMT_CREDIT	AMT_GOODS_PRICE	0.97
	REGION_RATING_CLIENT_W_CITY	REGION_RATING_CLIENT	0.95
	REGION_RATING_CLIENT	REGION_RATING_CLIENT_W_CITY	0.95
	CNT_CHILDREN	CNT_FAM_MEMBERS	0.89
	CNT_FAM_MEMBERS	CNT_CHILDREN	0.89
	DEF_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	0.87
	DEF_30_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE	0.87
	REG_REGION_NOT_WORK_REGION	LIVE_REGION_NOT_WORK_REGION	0.84
	LIVE_REGION_NOT_WORK_REGION	REG_REGION_NOT_WORK_REGION	0.84
	LIVE_CITY_NOT_WORK_CITY	REG_CITY_NOT_WORK_CITY	0.79
	REG_CITY_NOT_WORK_CITY	LIVE_CITY_NOT_WORK_CITY	0.79
	AMT_GOODS_PRICE	AMT_ANNUITY	0.73
	AMT_ANNUITY	AMT_GOODS_PRICE	0.73
	AMT_CREDIT	AMT_ANNUITY	0.72
	AMT_ANNUITY	AMT_CREDIT	0.72
	dtype: float64		

Processing Previous Application

a) Data Exploration

```
[135]: p_appl_data.head()
[135]:
   SK_ID_PREV    SK_ID_CURR NAME_CONTRACT_TYPE AMT_ANNUITY AN
0      2030495      271877  Consumer loans        1730.430
1      2802425      108129  Cash loans          25188.615
2      2523466      122040  Cash loans          15060.735
3      2819243      176158  Cash loans          47041.335
4      1784265      202054  Cash loans          31924.395
5 rows × 37 columns
<
[136]: p_appl_data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1670214 entries, 0 to 1670213
Data columns (total 37 columns):
 #   Column           Non-Null Count   Dtype  
--- 
 0   SK_ID_PREV       1670214 non-null   int64  
 1   SK_ID_CURR       1670214 non-null   int64  
 2   NAME_CONTRACT_TYPE 1670214 non-null   object  
 3   AMT_ANNUITY      1297979 non-null   float64
 4   AMT_APPLICATION  1670214 non-null   float64
 5   AMT_CREDIT        1670213 non-null   float64
 6   AMT_DOWN_PAYMENT  774370 non-null   float64
 7   AMT_GOODS_PRICE   1284699 non-null   float64
 8   WEEKDAY_APPR_PROCESS_START 1670214 non-null   object  
 9   HOUR_APPR_PROCESS_START 1670214 non-null   int64  
10  FLAG_LAST_APPL_PER_CONTRACT 1670214 non-null   object  
11  NFLAG_LAST_APPL_IN_DAY     1670214 non-null   int64 

[137]: p_appl_data.shape
[137]: (1670214, 37)
```

b) Handling NULL values

```
[500]: null_cols=((p_appl_data.isnull().sum()*100)/p_appl_data.shape[0]).round(2)
        null_cols=null_cols>0
[500]:
    AMT_ANNUITY            22.29
    AMT_DOWN_PAYMENT        53.64
    AMT_GOODS_PRICE          23.08
    RATE_DOWN_PAYMENT        53.64
    RATE_INTEREST_PRIMARY    99.64
    RATE_INTEREST_PRIVILEGED 99.64
    NAME_TYPE_SUITE          49.12
    CNT_PAYMENT              22.29
    PRODUCT_COMBINATION       0.02
    DAYS_FIRST_DRAWING       40.30
    DAYS_FIRST_DUE             40.30
    DAYS_LAST_DUE_1ST_VERSION 40.30
    DAYS_LAST_DUE              40.30
    DAYS_TERMINATION           40.30
    NFLAG_INSURED_ON_APPROVAL 40.30
dtype: float64
```

We see high percentage of null values in the above columns, dropping all columns having more than 20% of null values.

```
[502]: p_appl_data.drop(p_appl_data.loc[:,p_appl_data.isnull().mean()>=.20],axis=1,inplace=True)

[503]: ((p_appl_data.isnull().sum()*100)/p_appl_data.shape[0]).round(2)

[503]:
SK_ID_PREV          0.00
SK_ID_CURR          0.00
NAME_CONTRACT_TYPE  0.00
AMT_APPLICATION     0.00
AMT_CREDIT          0.00
WEEKDAY_APPR_PROCESS_START 0.00
HOUR_APPR_PROCESS_START 0.00
FLAG_LAST_APPL_PER_CONTRACT 0.00
NFLAG_LAST_APPL_IN_DAY 0.00
NAME_CASH_LOAN_PURPOSE 0.00
NAME_CONTRACT_STATUS 0.00
DAYS_DECISION        0.00
NAME_PAYMENT_TYPE    0.00
CODE_REJECT_REASON   0.00
NAME_CLIENT_TYPE     0.00
NAME_GOODS_CATEGORY  0.00
NAME_PORTFOLIO        0.00
NAME_PRODUCT_TYPE    0.00
CHANNEL_TYPE          0.00
SELLERPLACE_AREA      0.00
NAME_SELLER_INDUSTRY 0.00
NAME_YIELD_GROUP      0.00
PRODUCT_COMBINATION   0.02
dtype: float64

[504]: p_appl_data.shape

[504]: (1670214, 23)
```

After dropping high null value columns, Impute Product_combination with mode, it is the only column which has null values left.

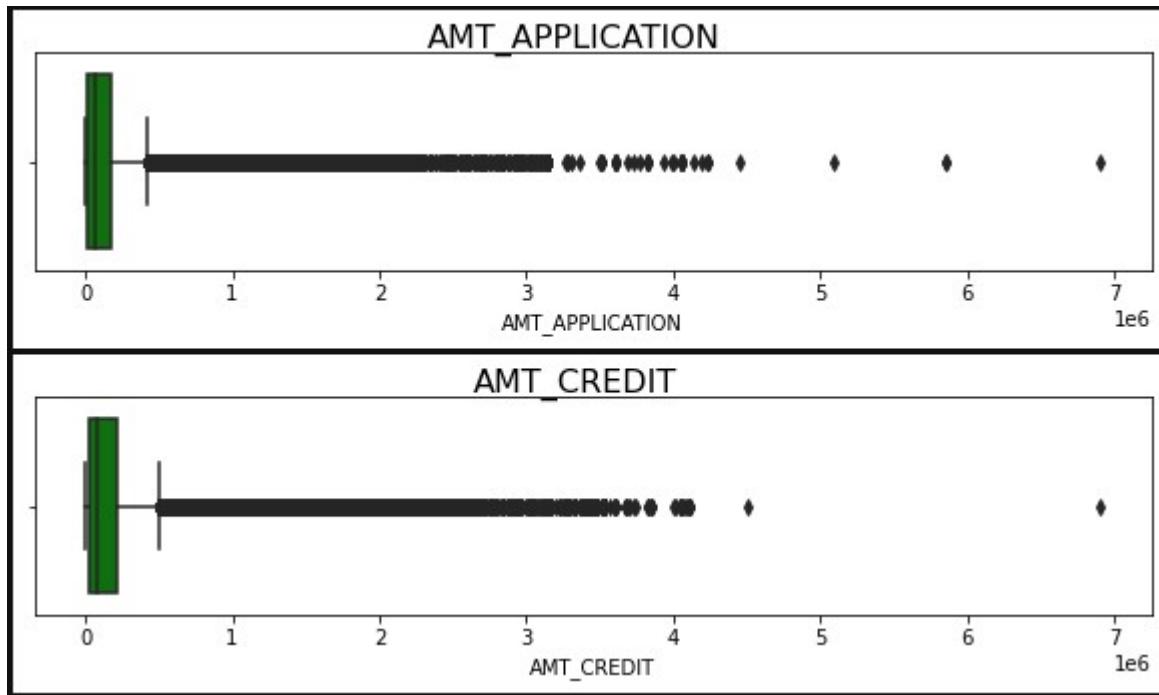
```
[505]: p_appl_data.PRODUCT_COMBINATION.mode()

[505]: 0    Cash
dtype: object

[506]: p_appl_data['PRODUCT_COMBINATION'].fillna(p_appl_data.PRODUCT_COMBINATION.mode()[0], inplace=True)
p_appl_data.PRODUCT_COMBINATION.isnull().sum()

[506]: 0
```

C) Identifying & Handling Outliers

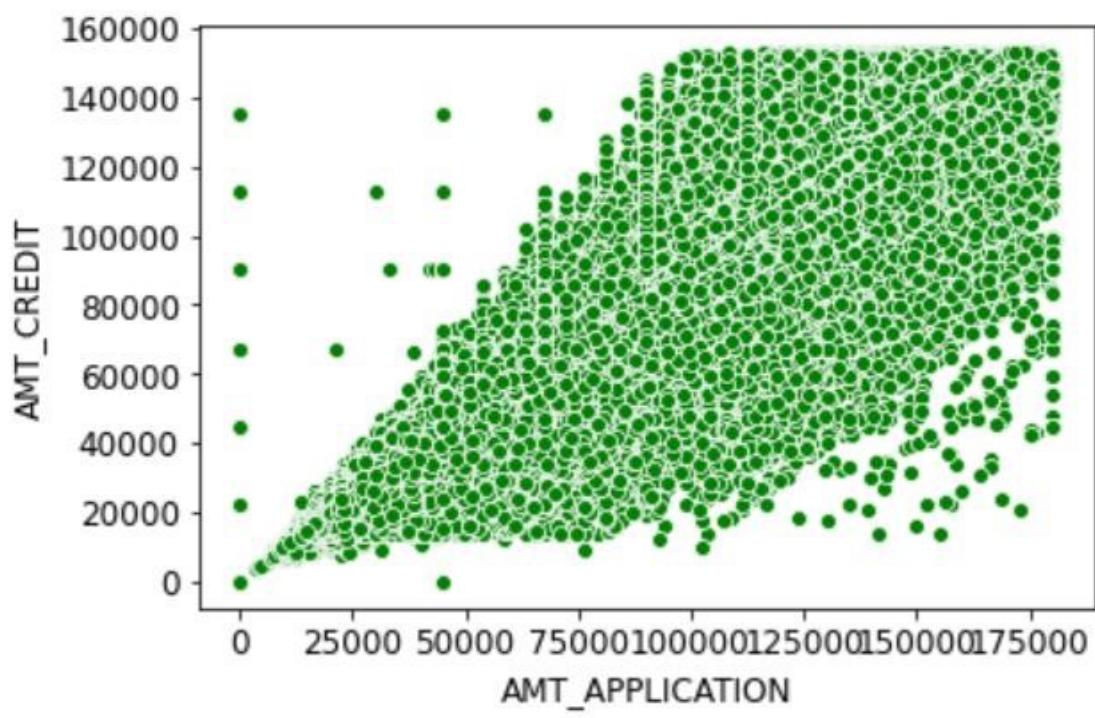


By looking at the data columns individually we decided on what percentile of data should be retained in the dataset. Which can make analysis accurate and does not result in skewed patterns. For determining the percentile (%) we looked at max values and retained only those outliers which are relatively closer to the 75% quantile values. We made sure all high value anomalies are completely removed.

Following is the list of columns and the corresponding quantiles below which we have retain the data.

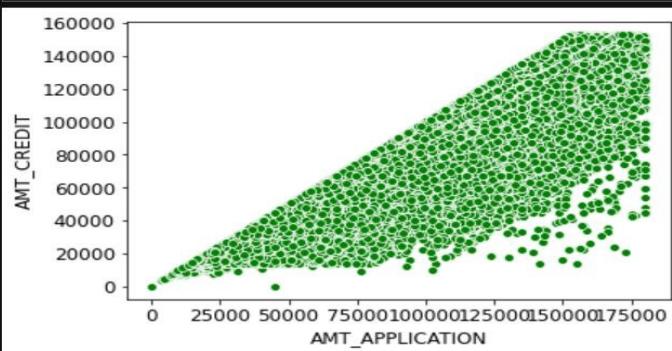
Columns Name	Retain % data
AMT_APPLICATION	75%
AMT_CREDIT	90%

d) Fixing AMT_CREDIT and AMT_APPLICATION Anomalies



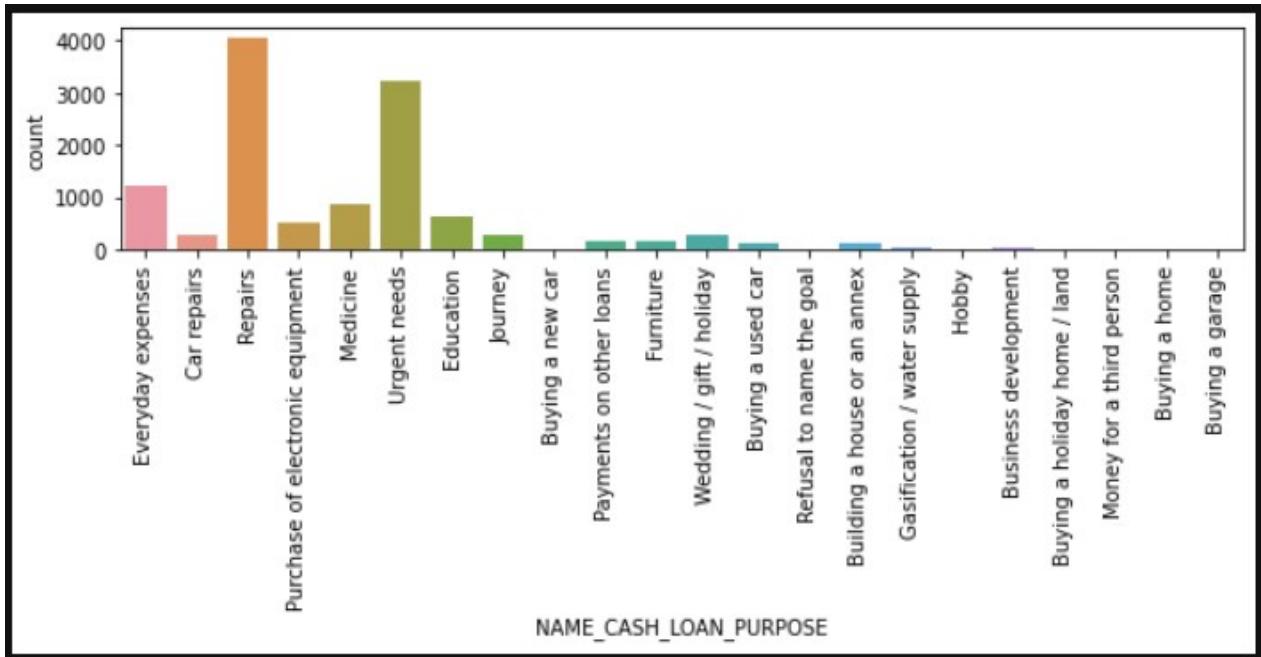
Ideally credit amount should not be greater than application amount. Hence, we should clean these values.

```
p_appl_data=p_appl_data[~(p_appl_data.AMT_CREDIT > p_appl_data.AMT_APPLICATION)]  
sns.scatterplot(x = 'AMT_APPLICATION', y="AMT_CREDIT",data=p_appl_data,color='green')  
plt.show()
```

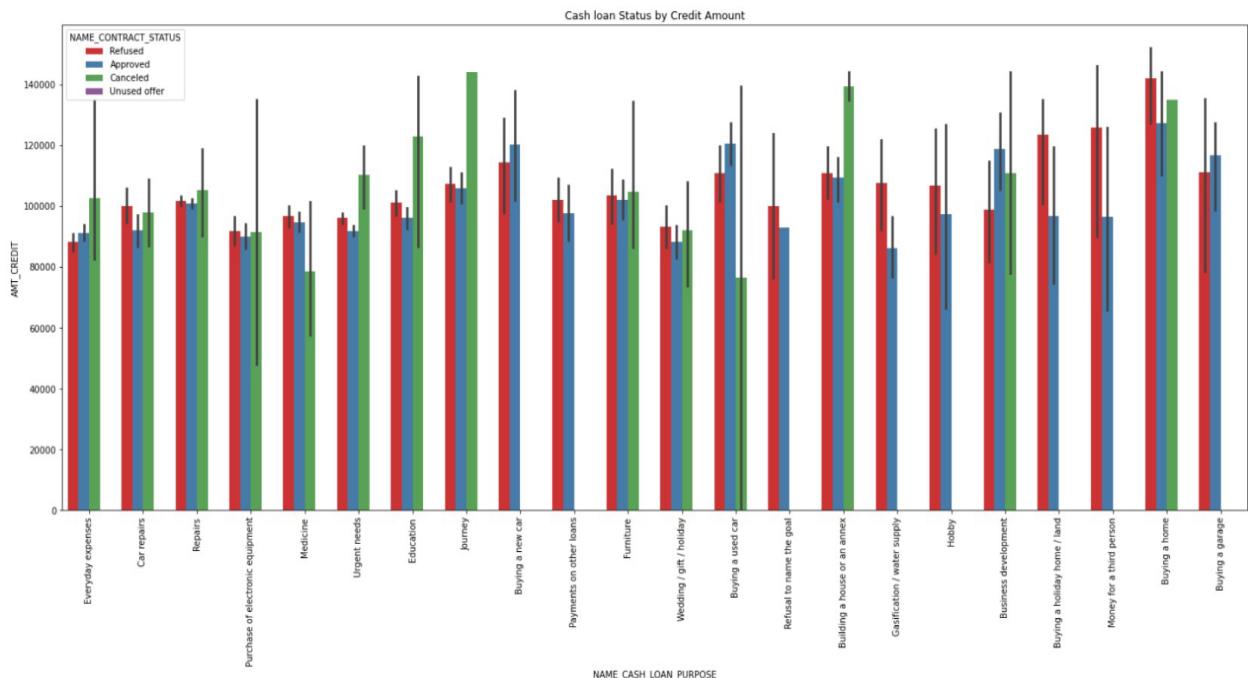


Performing Data Analysis and Identifying Insights

a) Cash Loans



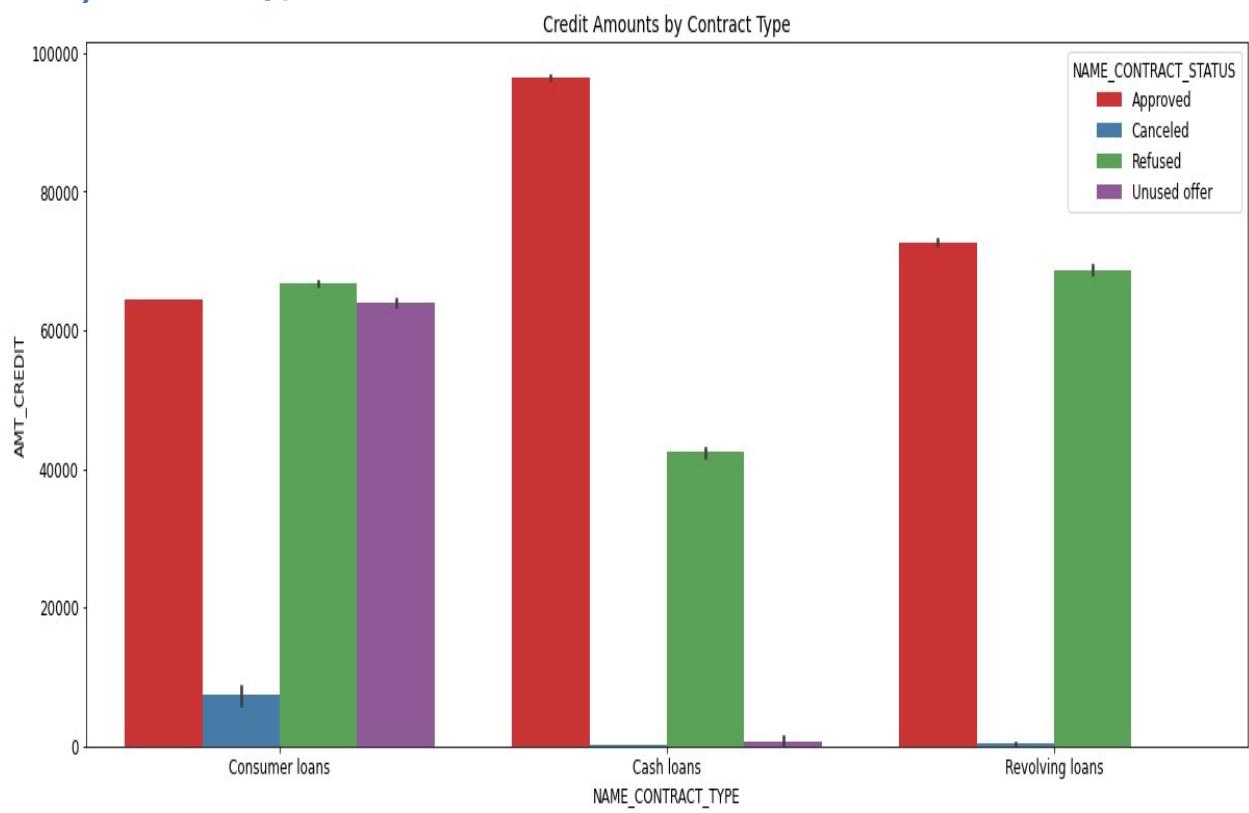
Most number of cash loan applications are for Repairs.



Observation:

- ⇒ Buying a home has the most cash loan approvals in terms of credit amounts.
- ⇒ Journey and building a house has the greatest number of cash loan cancellations.
- ⇒ Business development has the highest ratio of credit amounts approved vs refused.
- ⇒ Hobby, Gasification, buying new car, buying a garage has no cancellations.

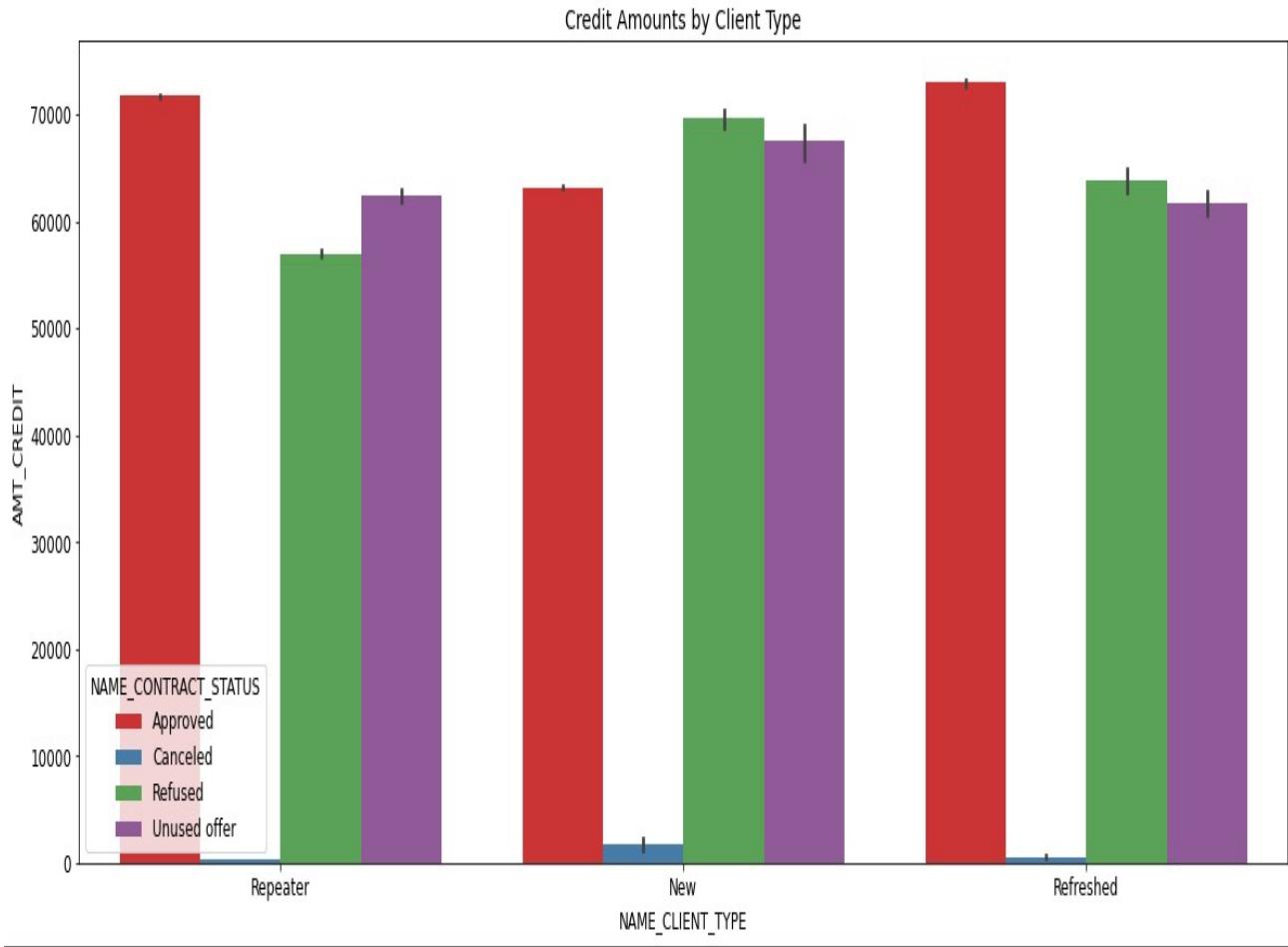
b) Contract Type



Observation:

- ⇒ Cash loans are approved for higher credit amounts than consumer and revolving loans and refused the least.
- ⇒ Consumer loans have higher cancellations as compared to cash loans and revolving loans.

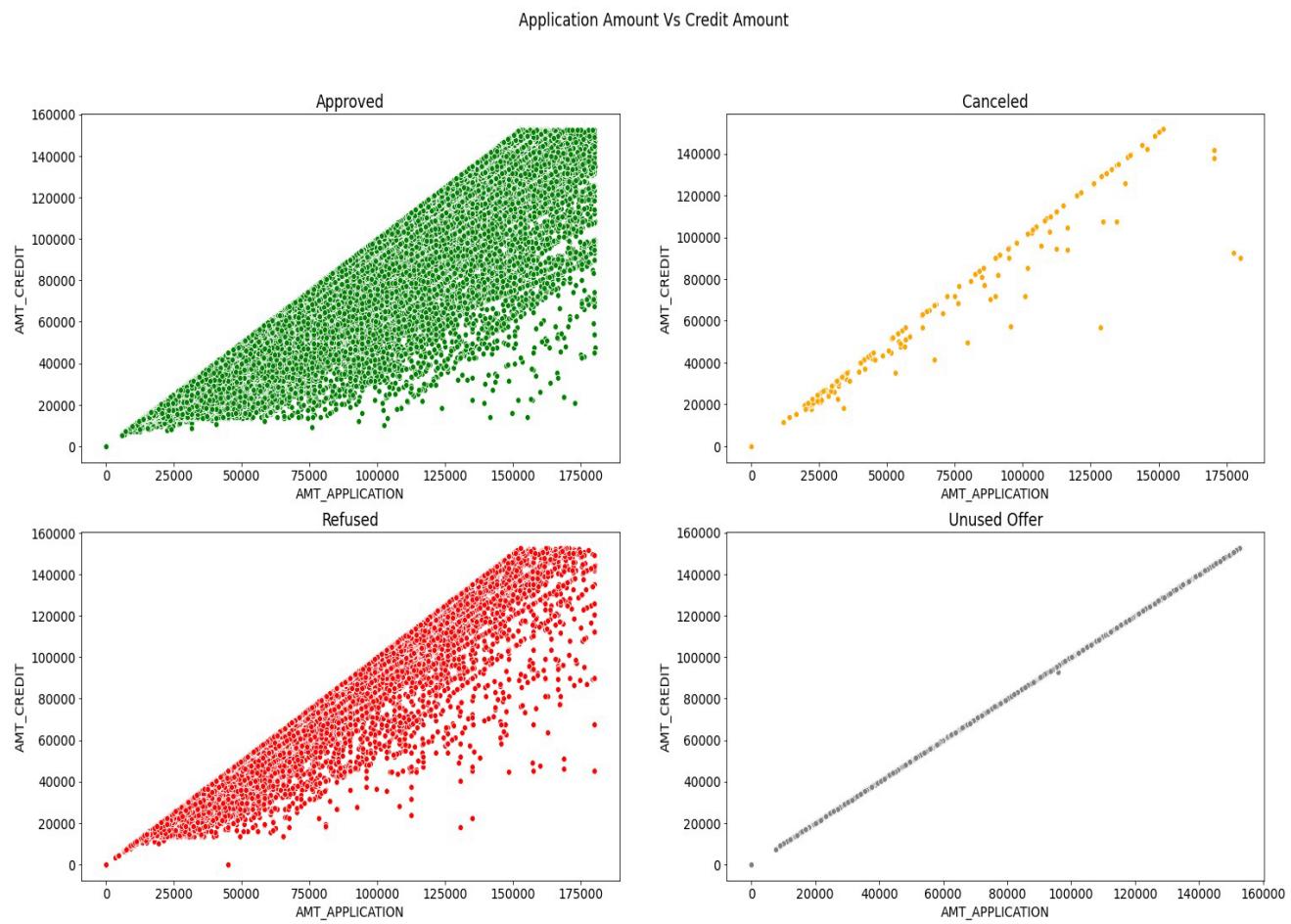
c) Client Type:



Observation:

- ⇒ Loans for greater than 70K+ are getting approved for Repeaters and Refreshed client as compared to New clients.
- ⇒ New clients have higher cancellations as compared to repeaters and refreshed clients.
- ⇒ New clients are refused on high credit loan applications.

Application Amount Vs Credit amount:



Observation:

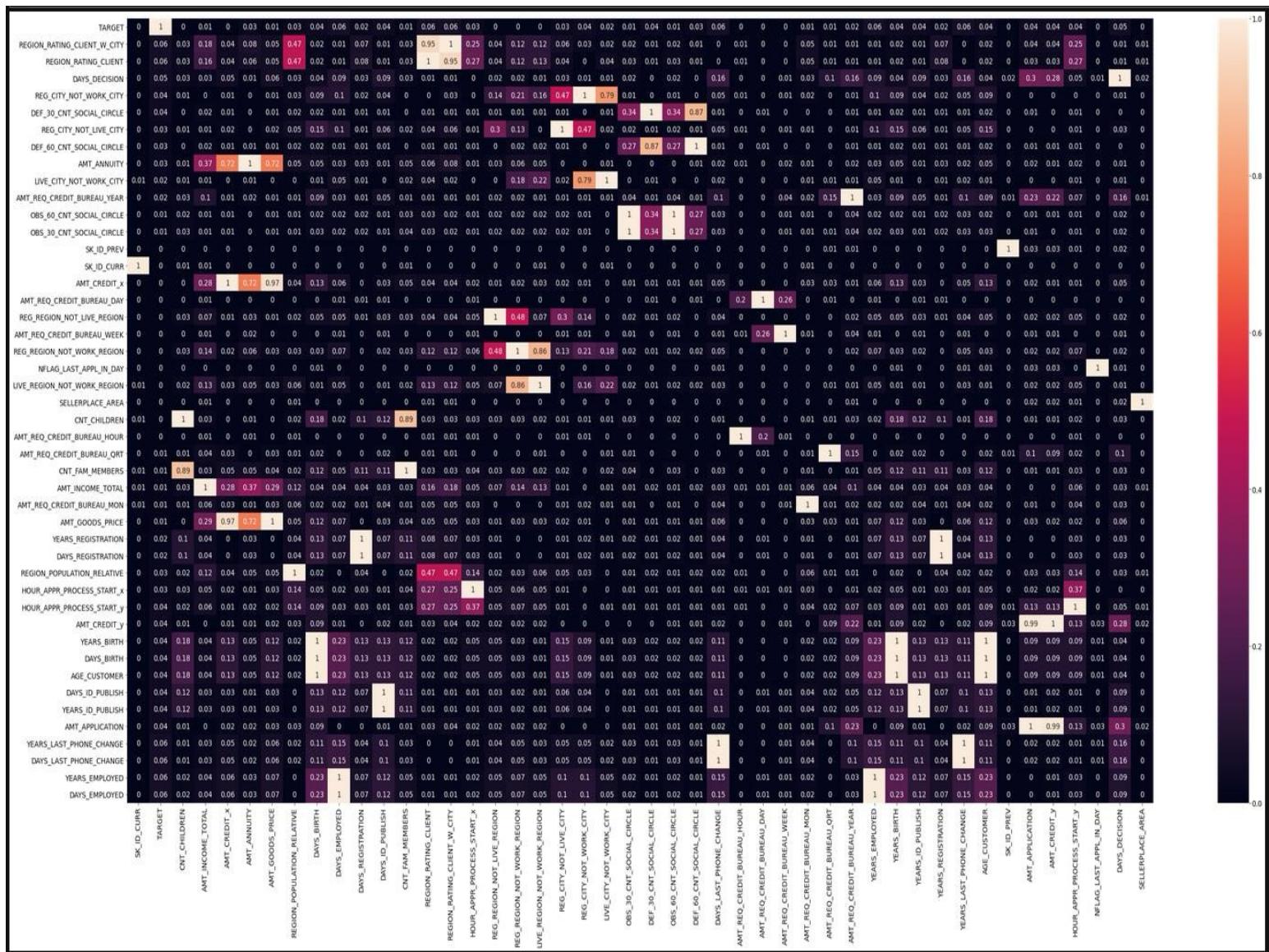
- ⇒ We can see a linear progression between credit amount and application amount across all decision's types.
- ⇒ We observed on cancelled applications the amount credit and amount application differ by some degree. This can be a reason why applications may have been cancelled.

Merging Current and Previous Application

9. Data Merging of Application data

```
combined_data = pd.merge(appl_data, p_appl_data, how='left', on=['SK_ID_CURR'])
combined_data
```

a) Identifying combined correlation of contributing factors

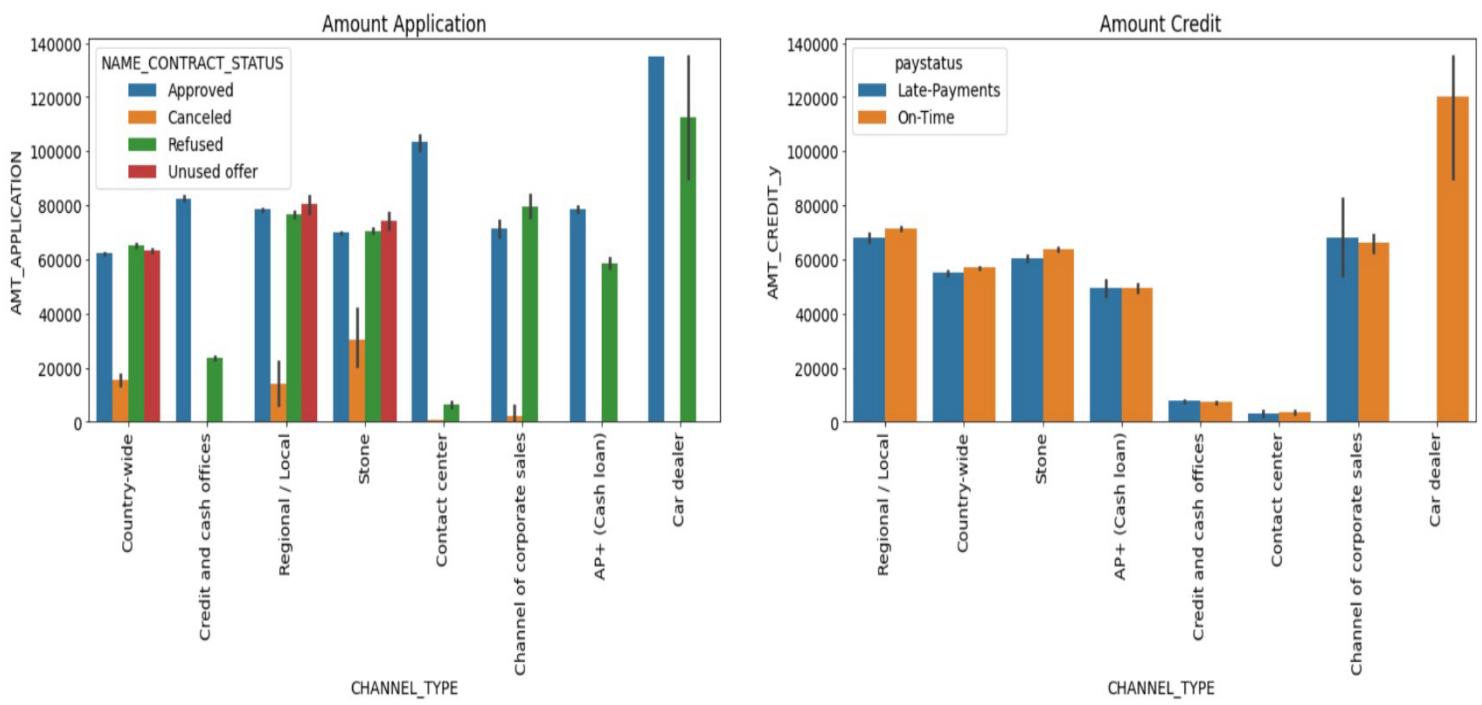


Based on the above heatmap, we have identified the highly correlated attributes:

AMT_APPLICATION	AMT_CREDIT_y	0.99
AMT_CREDIT_y	AMT_APPLICATION	0.99
AMT_GOODS_PRICE	AMT_CREDIT_x	0.97
AMT_CREDIT_x	AMT_GOODS_PRICE	0.97
REGION_RATING_CLIENT	REGION_RATING_CLIENT_W_CITY	0.95
REGION_RATING_CLIENT_W_CITY	REGION_RATING_CLIENT	0.95
CNT_CHILDREN	CNT_FAM_MEMBERS	0.89
CNT_FAM_MEMBERS	CNT_CHILDREN	0.89
DEF_30_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE	0.87
DEF_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	0.87
LIVE_REGION_NOT_WORK_REGION	REG_REGION_NOT_WORK_REGION	0.86
REG_REGION_NOT_WORK_REGION	LIVE_REGION_NOT_WORK_REGION	0.86
REG_CITY_NOT_WORK_CITY	LIVE_CITY_NOT_WORK_CITY	0.79
LIVE_CITY_NOT_WORK_CITY	REG_CITY_NOT_WORK_CITY	0.79
AMT_ANNUITY	AMT_GOODS_PRICE	0.72
AMT_CREDIT_x	AMT_ANNUITY	0.72
AMT_GOODS_PRICE	AMT_ANNUITY	0.72

b) Loan Decisions based on Channel Type:

Loan Decisions based on Channel Type



Observation:

- ⇒ Car dealers are good customers as they have less late payments and are approved more than they are refused.
- ⇒ Contact center are highly approved and have almost equal on-time vs late payment ratio.
- ⇒ Channel of Corporate sales are risky as they have higher late payment ratios and are refused loans more than approved.