# CASE STUDY – Credit EDA

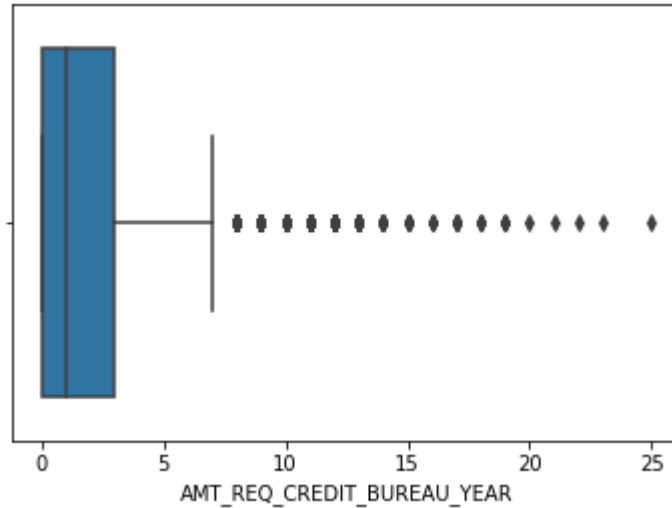By **Syed Saifullah Tarique and Aman Jha**

# Business Understanding

- The loan providing companies find it hard to give loans to the people due to their insufficient or non-existent credit history. Because of that, some consumers use it as their advantage by becoming a defaulter.

- When the company receives a loan application, the company has to decide for loan approval based on the applicant's profile. Two types of risks are associated with the bank's decision:

- If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company

- If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company.

- The data given below contains the information about the loan application at the time of applying for the loan. It contains two types of scenarios:

- **The client with payment difficulties:** he/she had late payment more than X days on at least one of the first Y instalments of the loan in our sample,

- **All other cases:** All other cases when the payment is paid on time.

- When a client applies for a loan, there are four types of decisions that could be taken by the client/company):

- **Approved:** The Company has approved loan Application

- **Cancelled:** The client cancelled the application sometime during approval. Either the client changed her/his mind about the loan or in some cases due to a higher risk of the client he received worse pricing which he did not want.

- **Refused:** The company had rejected the loan (because the client does not meet their requirements etc.).

- **Unused offer:** Loan has been cancelled by the client but on different stages of the process.
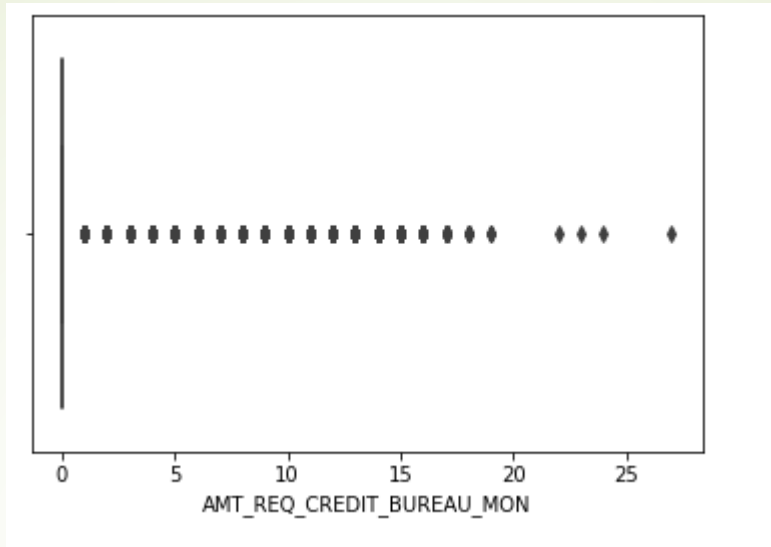
# Analysis – Application Data

- Missing Values and approaches

- Median of the column is the 50th percentile value of that column in the above table

- In the case of missing categorical column values we ignore them mostly.

- For the 'EXT_SOURCE_3' column we should fill the missing values with the Median because both are close so we prefer the Median because it is much unbiased approach.

- For the 'AMT_REQ_CREDIT_BUREAU_YEAR' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.

- For the 'AMT_REQ_CREDIT_BUREAU_MON' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.

- For the 'AMT_REQ_CREDIT_BUREAU_WEEK' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.

- For the 'AMT_REQ_CREDIT_BUREAU_DAY' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.

- For the 'AMT_REQ_CREDIT_BUREAU_HOUR' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.
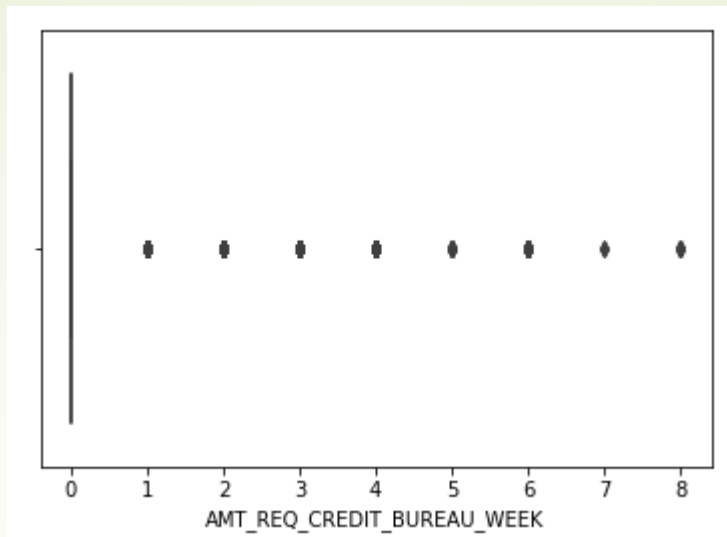
- For the 'AMT_REQ_CREDIT_BUREAU_QRT' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.

- For the 'OBS_30_CNT_SOCIAL_CIRCLE' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.

- - For the 'DEEF_30_CNT_SOCIAL_CIRCLE' column we should fill the missing values with the Median because it clearly have an outlier by looking at the maximum value and the mean value. And since the mean value gets influenced by the outlier and median don't so we would use median value here.
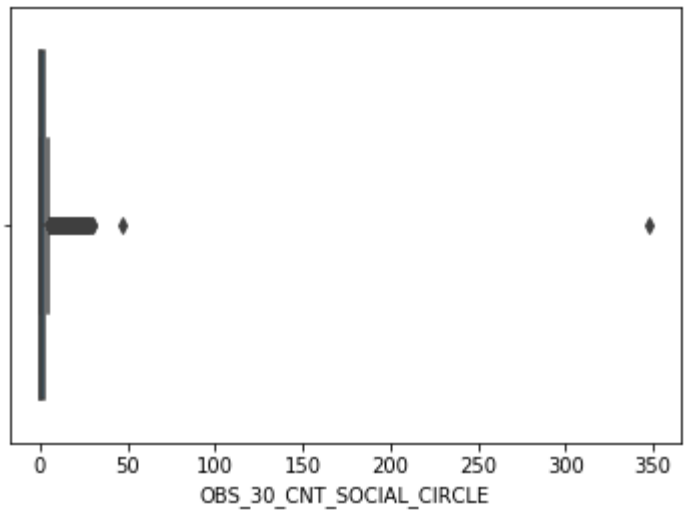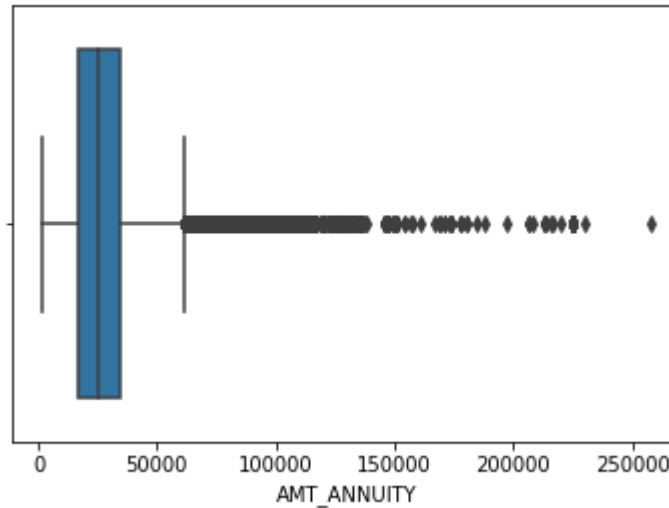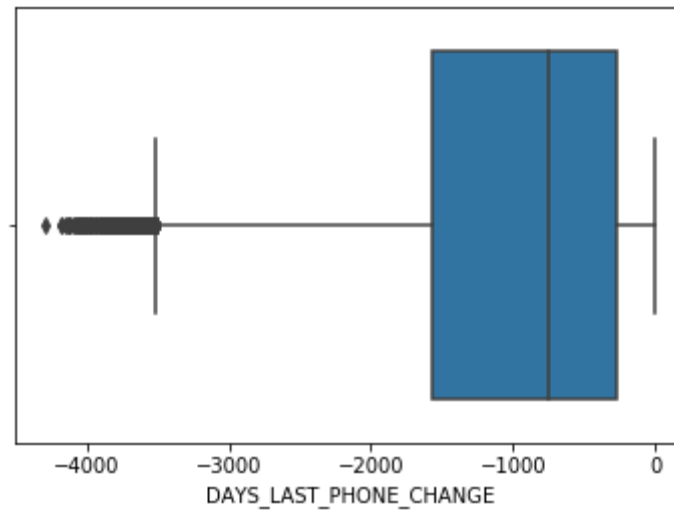
# Various Plots



AMT_REQ_CREDIT_BUREAU_YEAR

- Here we could easily see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

- This column represent the number of enquiries to credit bureau about the client one year before the application.

- So the outliers could be a mistake in which we have to rectify it(if we could) or remove it. In other case outlier could be a variance in the data which we could confirm by asking the company or by checking the other similar data for that employee.
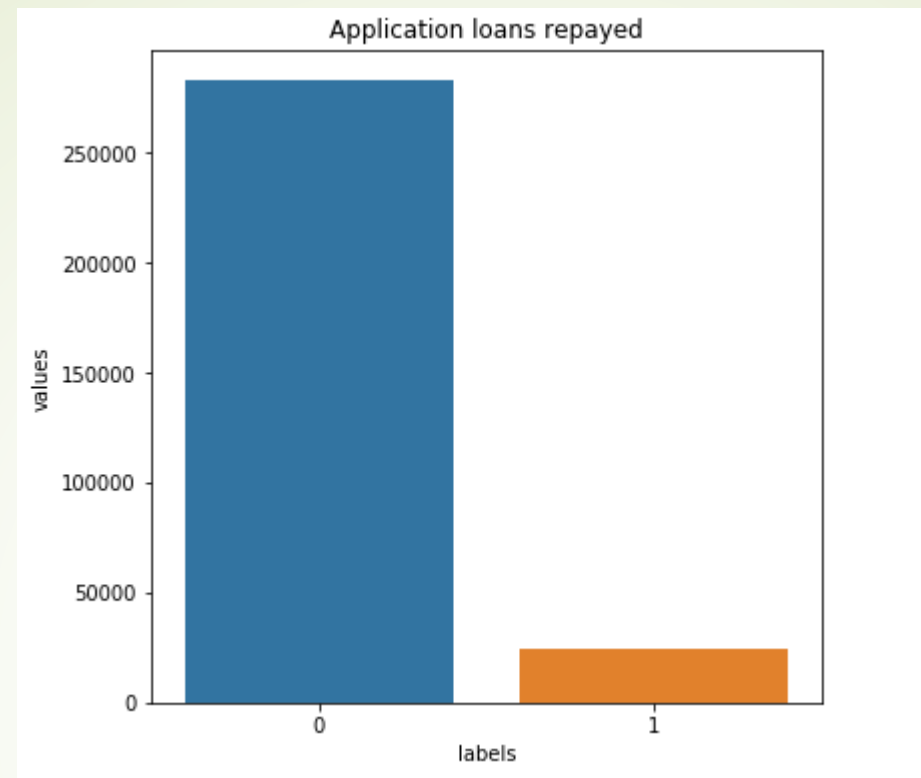
- Here we could easily see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

- This column represent the number of enquiries to credit bureau about the client a month before the application.

- So the outliers could be a mistake in which we have to rectify it(if we could) or remove it. In other case outlier could be a variance in the data which we could confirm by asking the company or by checking the other similar data for that employee.
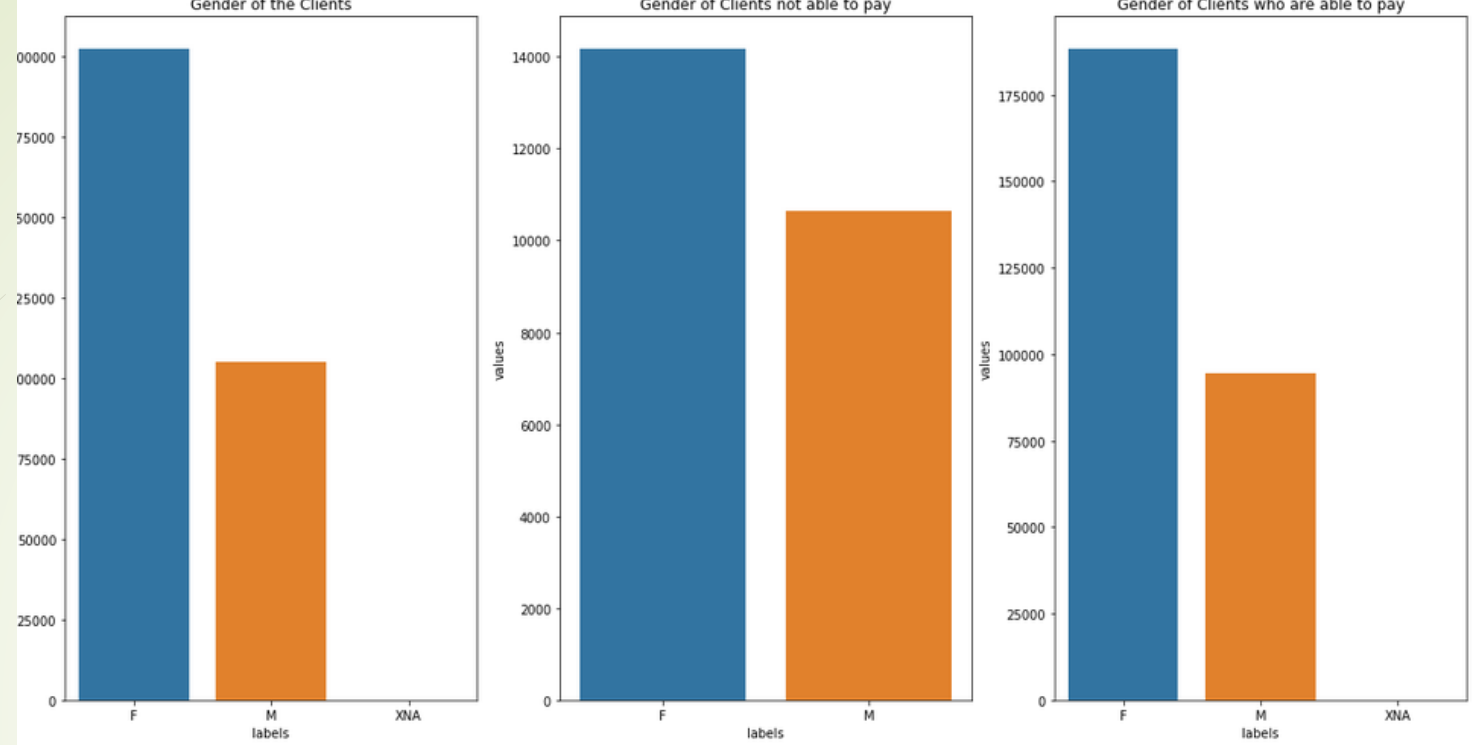
AMT_REQ_CREDIT_BUREAU_WEEK

- Here we could easily see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

- - This column represent the number of enquiries to credit bureau about the client a week before the application.

- - So the outliers could be a mistake in which we have to rectify it(if we could) or remove it. In other case outlier could be a variance in the data which we could confirm by asking the company or by checking the other similar data for that employee.
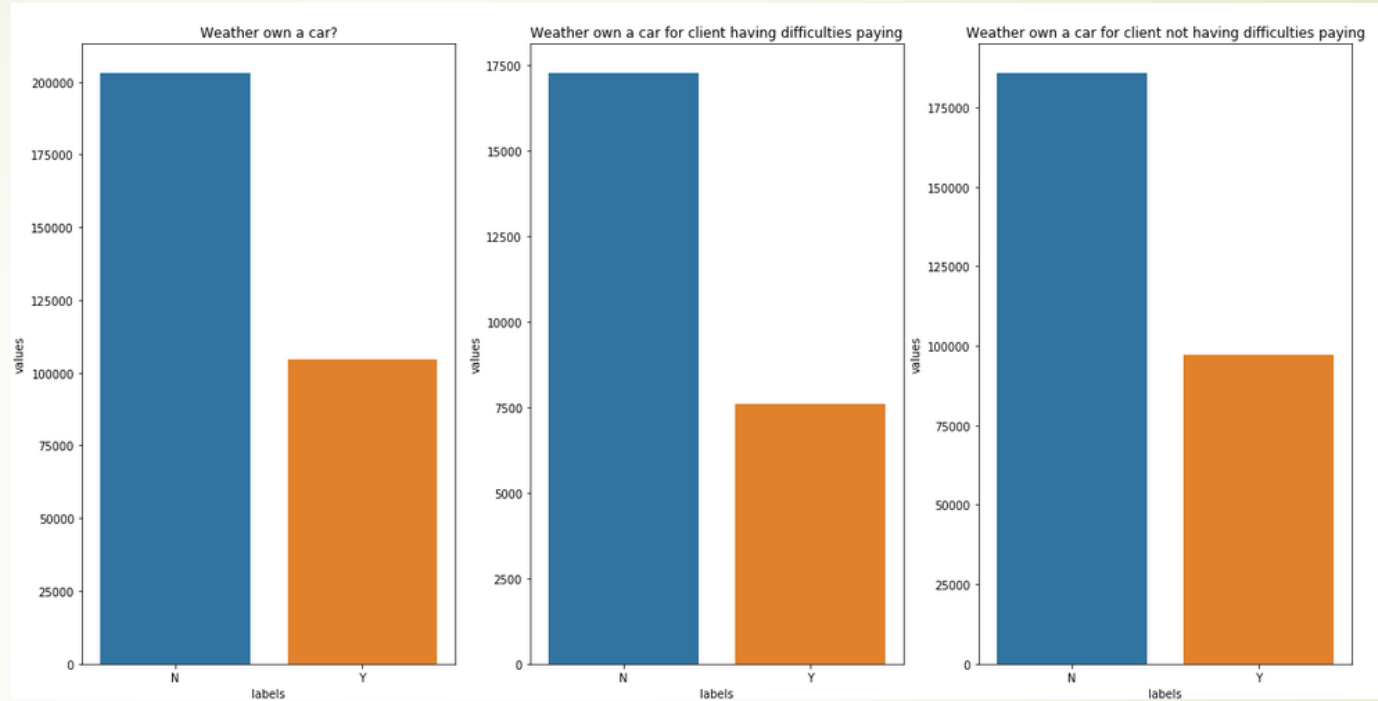
OBS_30_CNT_SOCIAL_CIRCLE

- Here we could easily see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

- - This column represent how many observation of the client social surrounding with observable 30 days past due.

- - Here the outliers is a mistake in which we have to rectify it (if we could) or remove it because the values can't be greater than 30 days.

AMT_ANNUITY

- Here we could easily see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

- This column represent the loan annuity.

- So the outliers could be a mistake in which we have to rectify it(if we could) or remove it. In other case outlier could be a variance in the data which we could confirm by asking the company or by checking the other similar data for that employee.

DAYS_LAST_PHONE_CHANGE

▬ - Here we could easly see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

▬ - This column represent how many days before the application did the client cahnge phone number.

▬ - Here the outliers is a mistake in which we have to rectify it (if we could) or remove it because he can't change the number very large number of days before because he won't be born to do so.

Application loans repayed

- In the above graph the 1 means Client has payment difficulties and 0 means he doesn't have payment difficulties

- We could see that there are more number of Female clients than Male Clients (Nearly twice as much).

- Here we could see in number wise there are more Female Client who has problem in paying the loan but If we see the above graph then we could see that there are much more number of Female who are taking loan. Therefore percentage wise there are more Male clients having problem in repaying loan than Female.

- Clients who don't own a car are nearly twice of those people who owns a car.
- Their are more client that don't have a car have difficulties paying than the ones who has a car.

- Revolving type of loans are very small fraction nearly 10% of the loans.

- There are large amount of cash loans are there which are not being paid but there are also very large amount of revolving loans are not being paid compare to its frequency.

- Clients that owns realty are nearly twice that who doesn't.
- Here weather the client would pay the loan is not much effected by weather they own the realty or not.

- Most of the clients don't have children.
- There are a large number of people who don't have children have problem paying the loan

- Most of clients are married, followed by Single/not married and civil marriage.

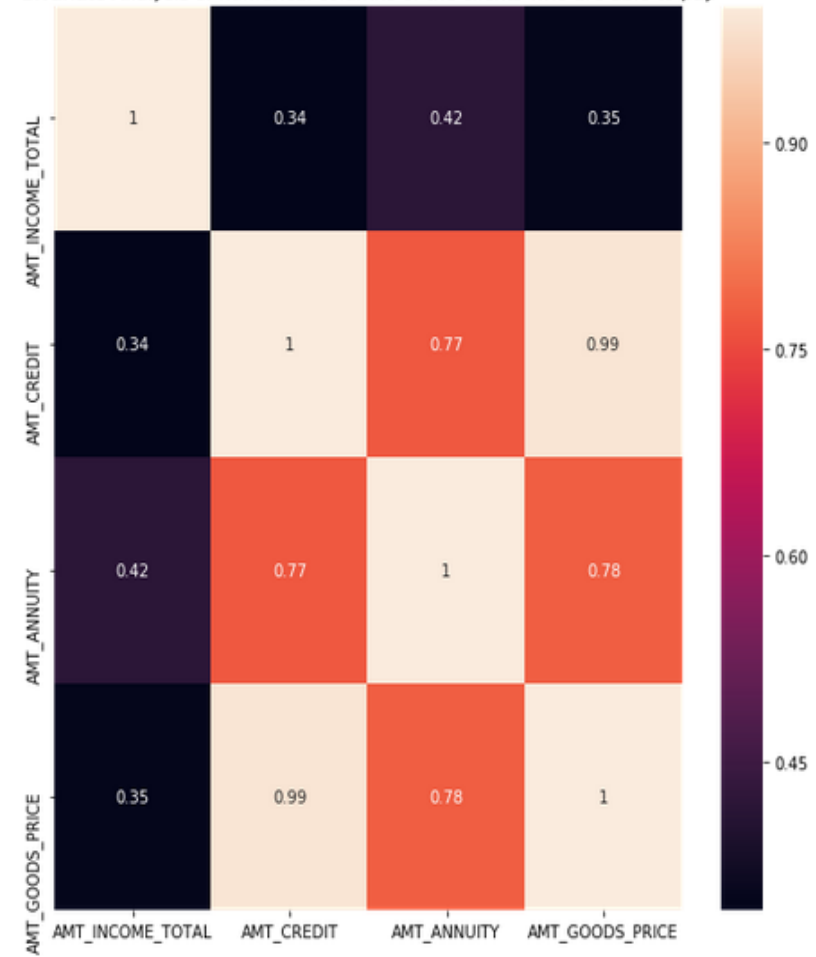- There are large number of client who are not able to pay are married and the least number of client who are unbale to pay are widowed.

Most of the loans are taken by Laborers,
followed by Sales staff.
IT staff take the lowest amount of loans.
Large numbers of unpaid loan by
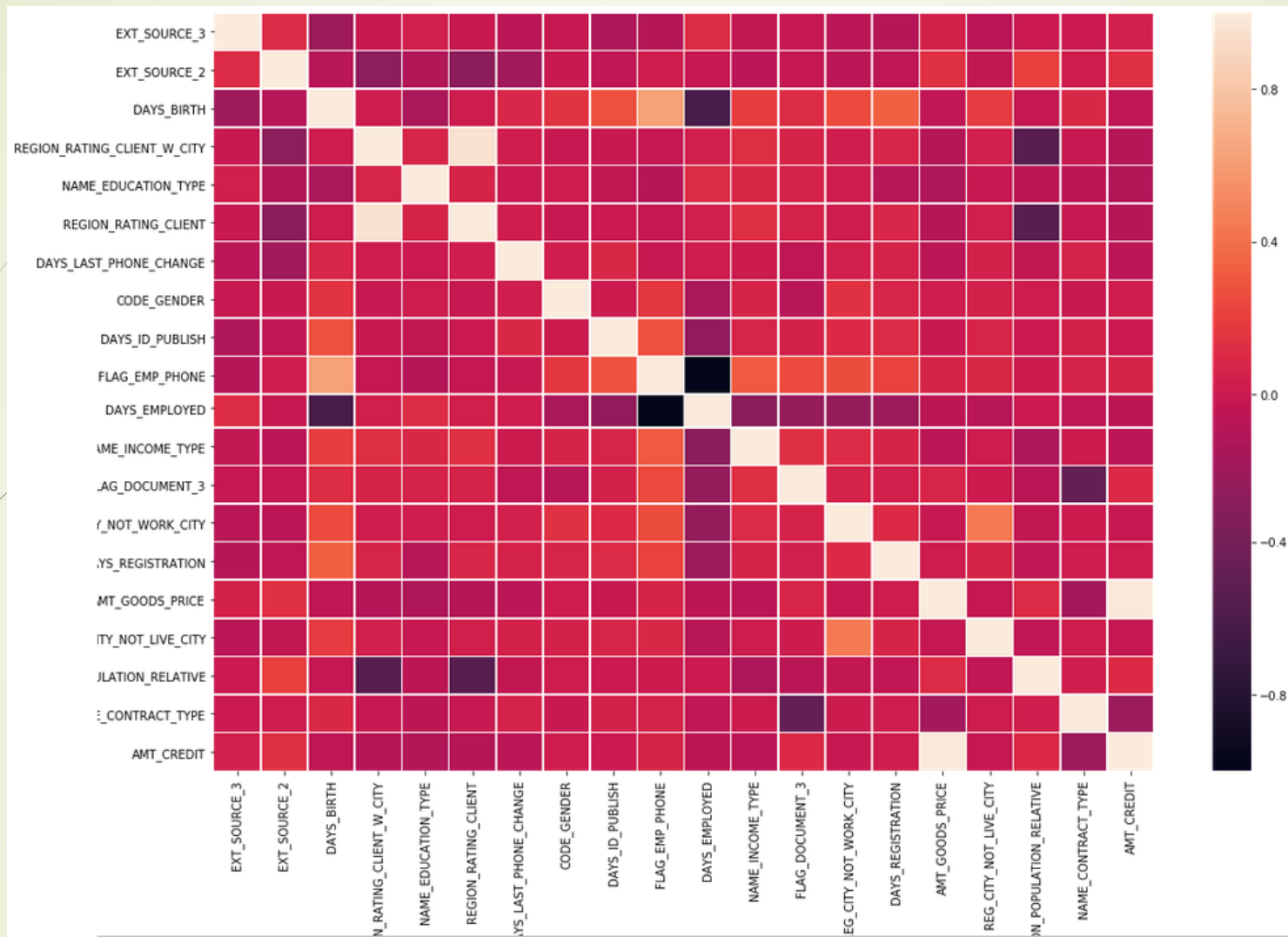Laborers and Drivers followed by sales staff

Bivariate Analysis of numerical columns of the clients who are not able to pay

Bivariate Analysis of numerical columns of the clients who are able to pay

- Here we could see that AMT_ANNUITY and AMT_CREDIT are highly correlated, if one goes up so s the other and vise versa. For both the sets but much more correlated for the Clients who are able to pay.

- - We could also see that AMT_CREDIT and AMT_GOODS_PRICE are also highly correlated, if one goes up so s the other and vise versa. For both the sets but are more correlated to the Clients who are are able to pay.

- - Similarly AMT_ANNUITY and AMT_GOODS_PRICE are also highly correlated, if one goes up so s the other and vise versa. For both the sets but are more correlated to the Clients who are able to pay.

- - For AMT_ANNUITY and AMT_INCOME_TOTAL the correlation is much higher for the Clients who are able to pay. Similar is the case with AMT_CREDIT and AMT_INCOME_TOAL.

- We can have look on the values and there correction. Lighter the color more corrected values.

- - 'FLAG_EMP_PHONE' is highly correlated with 'DAYS_BIRTH' i.e if one goes up so is the other.

- Similarly 'REG_CITY_NOT_LIVE_CITY' is highly correlated with 'REG_CITY_NOT_WORK_CITY'.

- There are some which are highly negatively correlated that means that if one increases other decreases. e.g. 'DAYS_EMPLOYED' and 'FLAG_EMP_PHONE' etc.
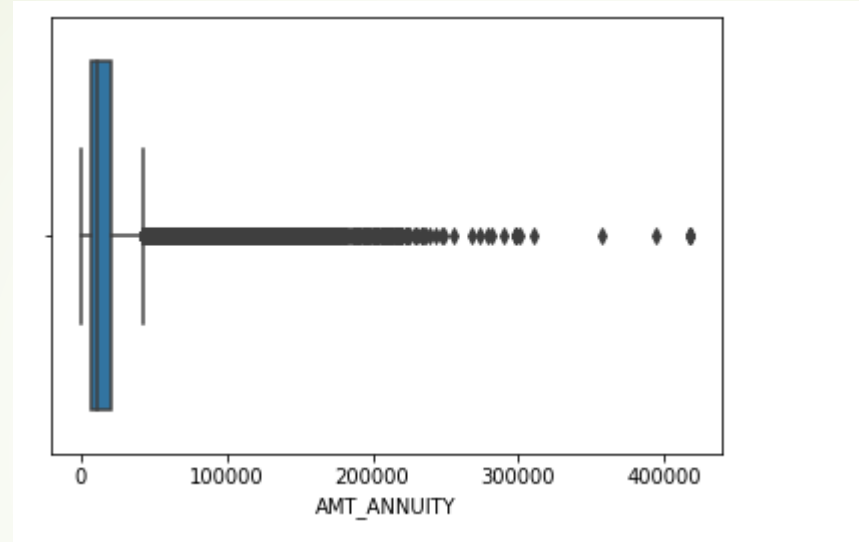
# The top 10 correlation for the Client with payment difficulties and all other cases (Target variable)

'EXT_SOURCE_3', 'EXT_SOURCE_2', 'NAME_EDUCATION_TYPE', 'DAYS_LAST_PHONE_CHANGE', 'CODE_GENDER', 'DAYS_ID_PUBLISH', 'NAME_INCOME_TYPE', 'REGION_POPULATION_RELATIVE', 'NAME_CONTRACT_TYPE', 'AMT_CREDIT'
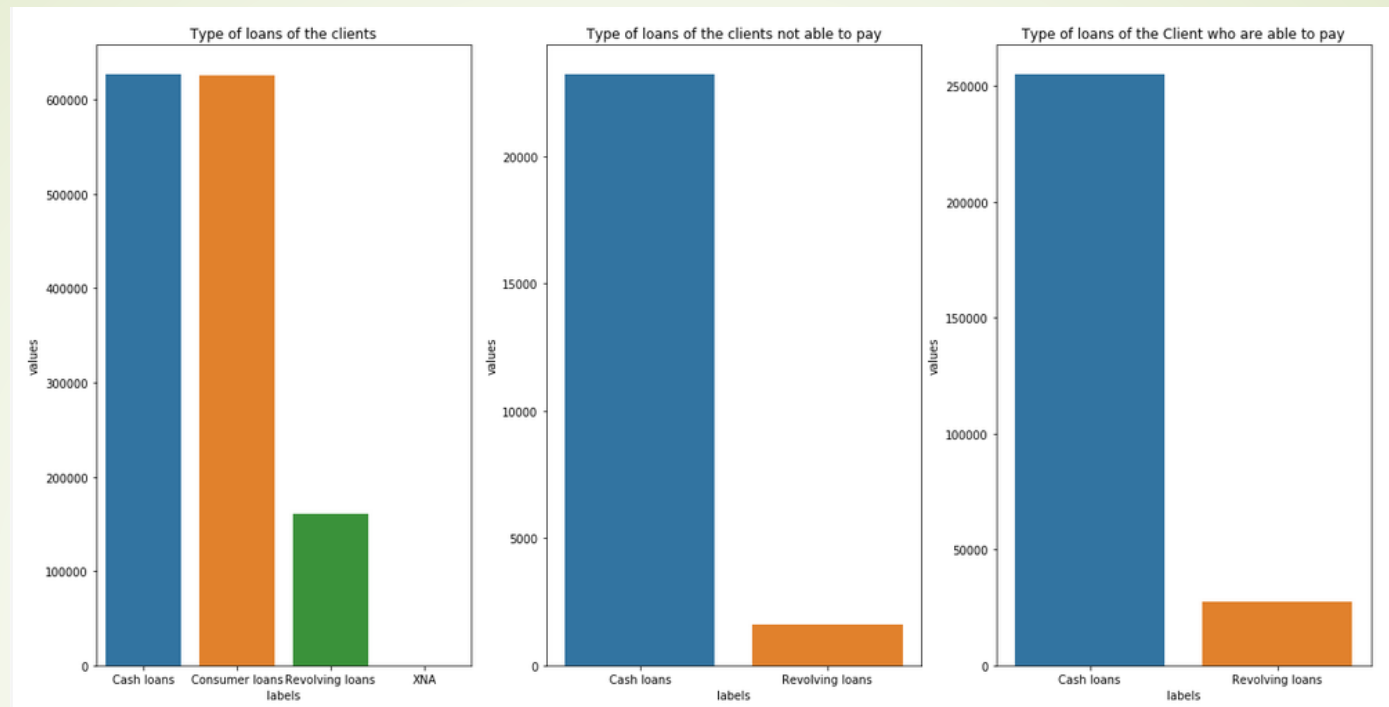
Loan defaulting is directly(positively) correlated to NAME_EDUCATION_TYPE,DAYS_LAST_PHONE_CHANGE,CODE_GENDER,DAYS_ID_PUBLISH,NAME_INCOME_TYPE

Loan defaulting is inverserly (negativly) correlated to EXT_SOURCE_3,EXT_SOURCE_2,REGION_POPULATION_RELATIVE,NAME_CONTRACT_TYPE,AMT_CREDIT

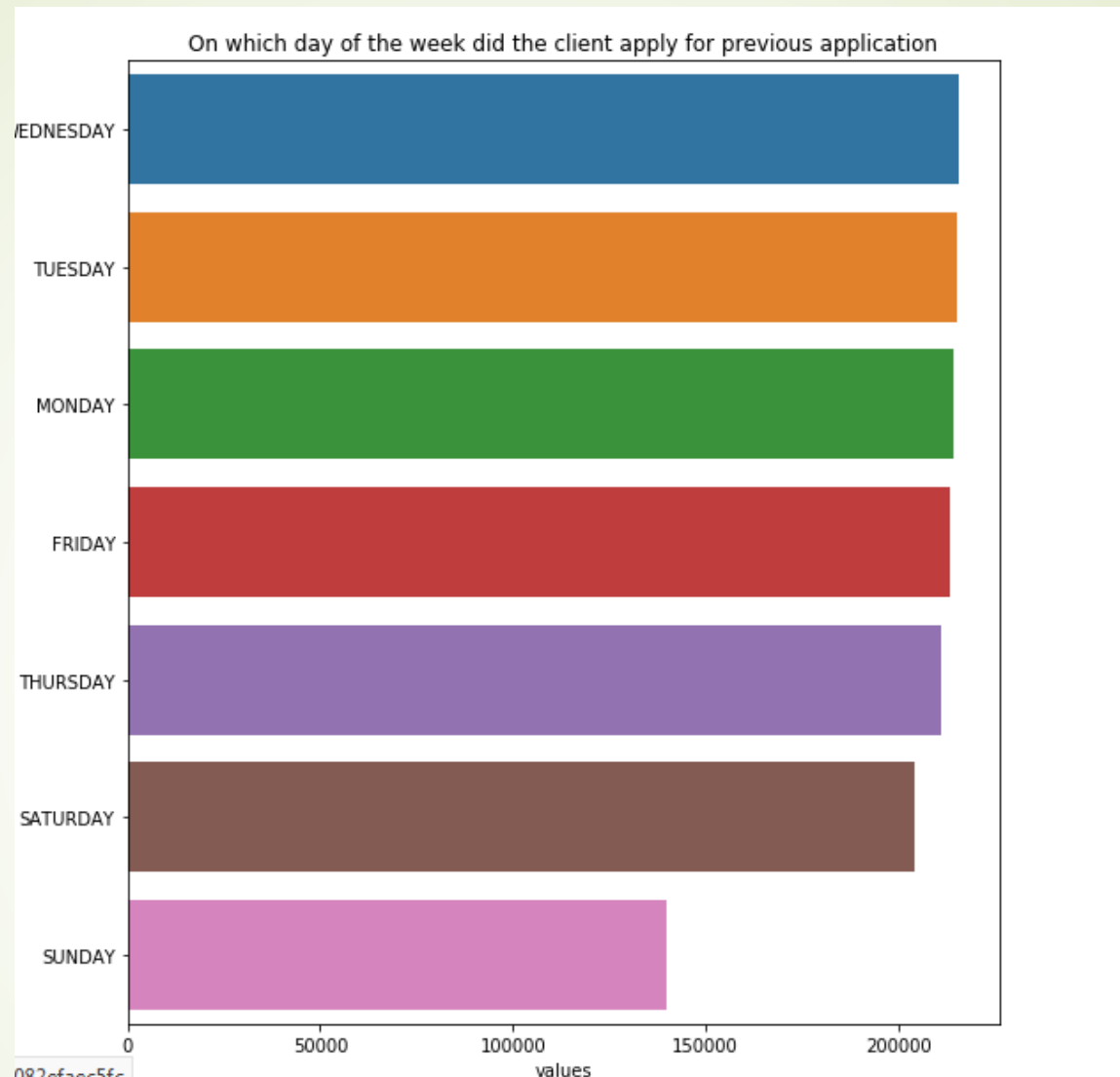# Analysis Previous Application
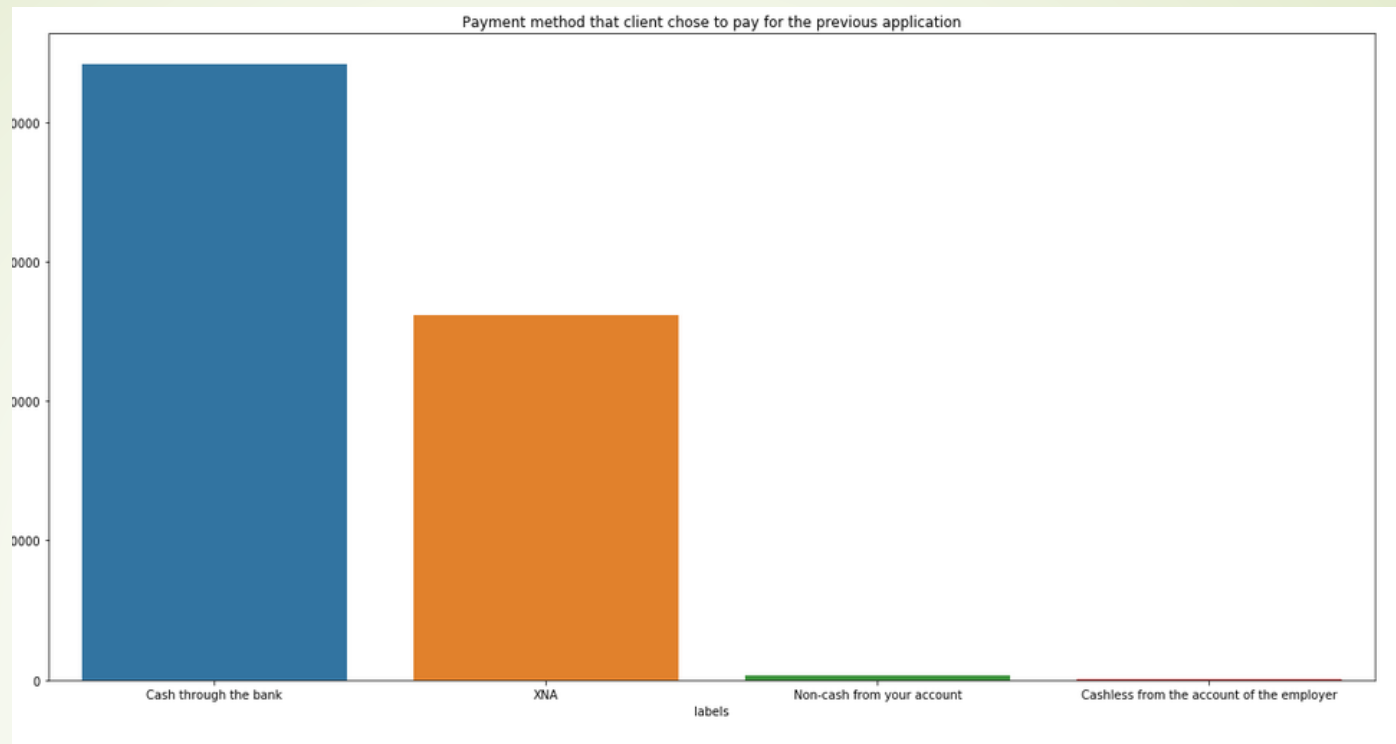


- ere we could easily see in the above plot that this column has outliers. Because those are the points which are way far away from the most of the data range as could be seen in the graph.

- - This column represent the annuity of the previous application.

- - So the outliers could be a mistake in which we have to rectify it(if we could) or remove it. In other case outlier could be a variance in the data which we could confirm by asking the company or by checking the other similar data for that employee.
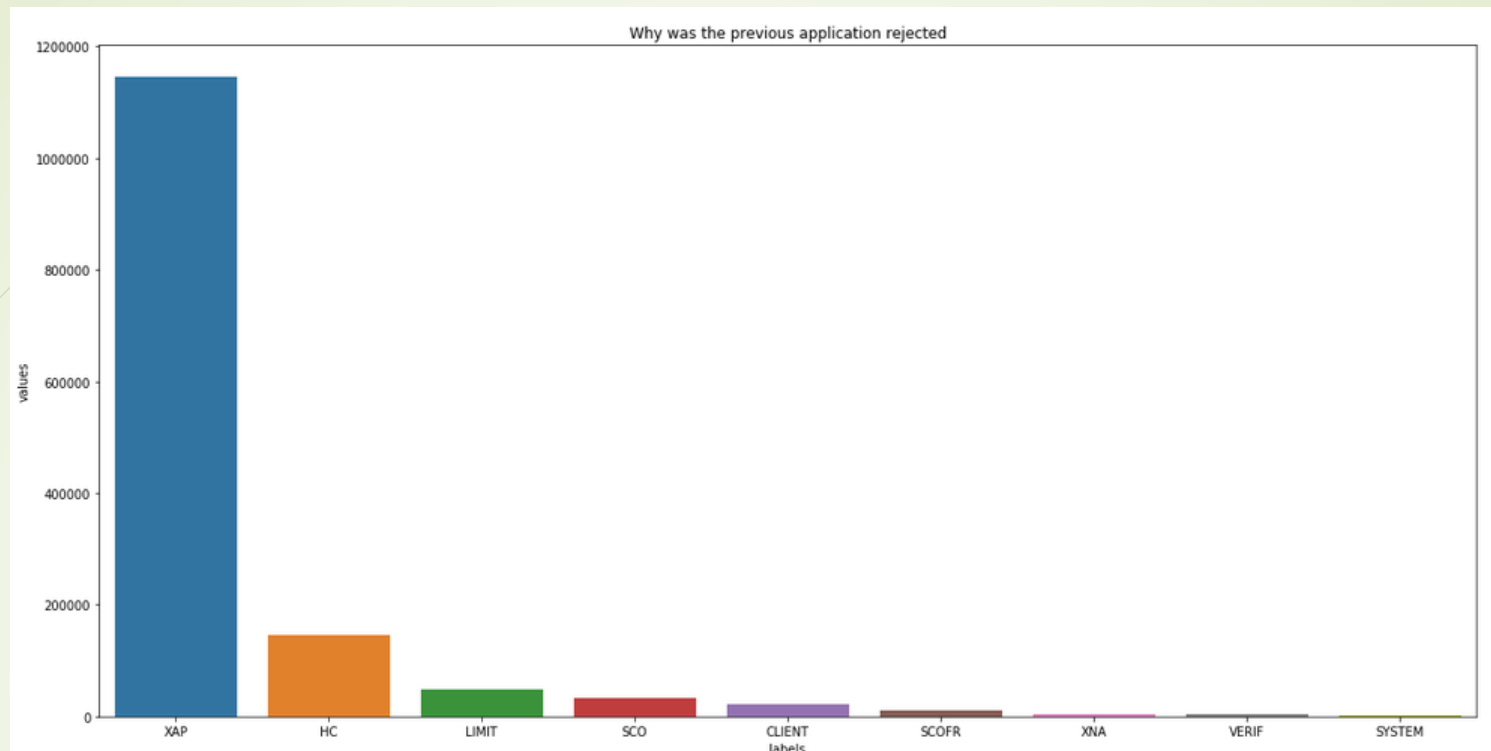
- We could see that the Cash loans are the large amount of which loans are unable to being paid

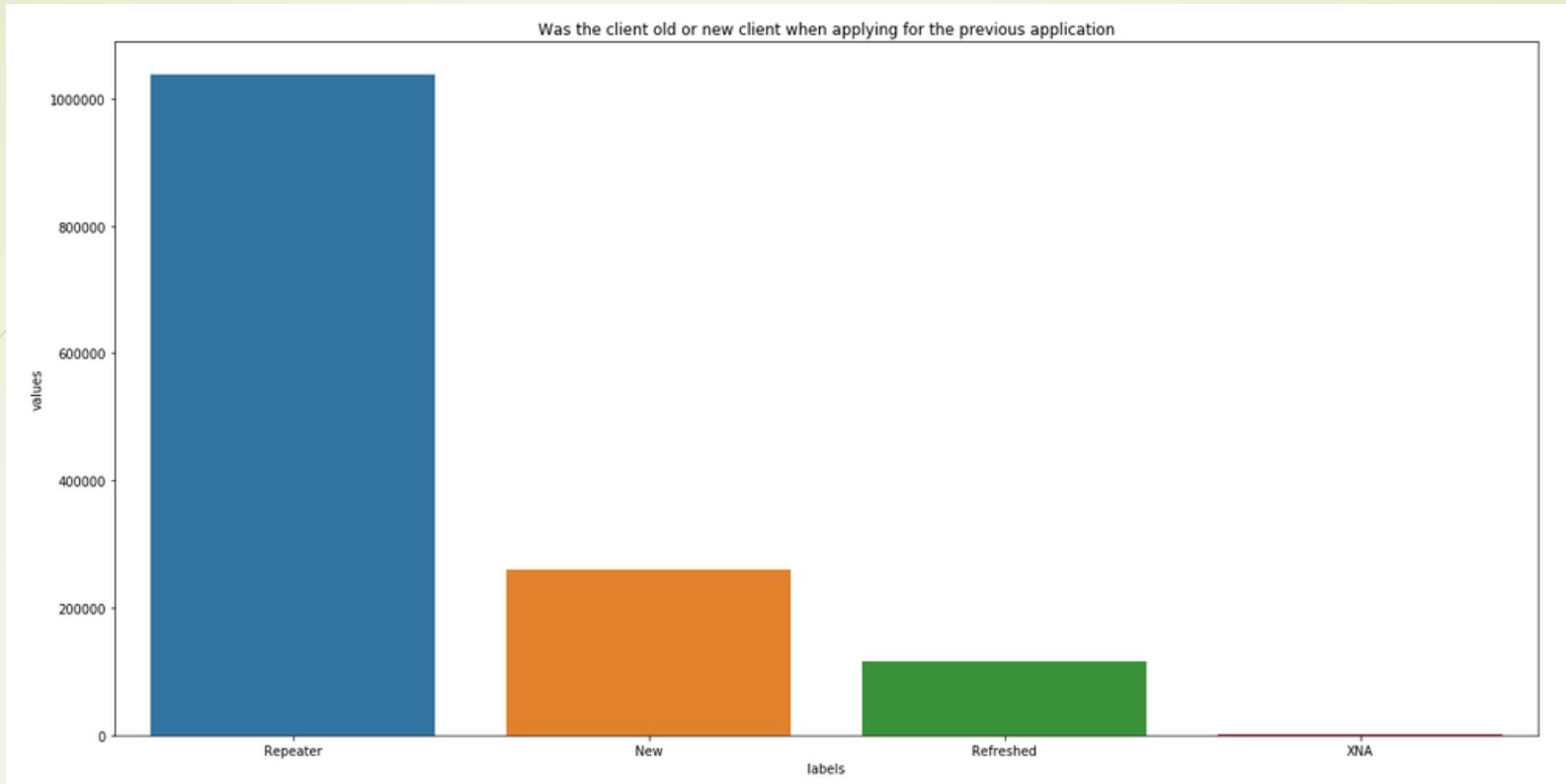On which day of the week did the client apply for previous application

- Here we could see that except for Sunday loans are taken almost equally on every day of the week, but on Sunday there is quite decrease in the amount of loan taken. Which may be due the Weekend Holiday.
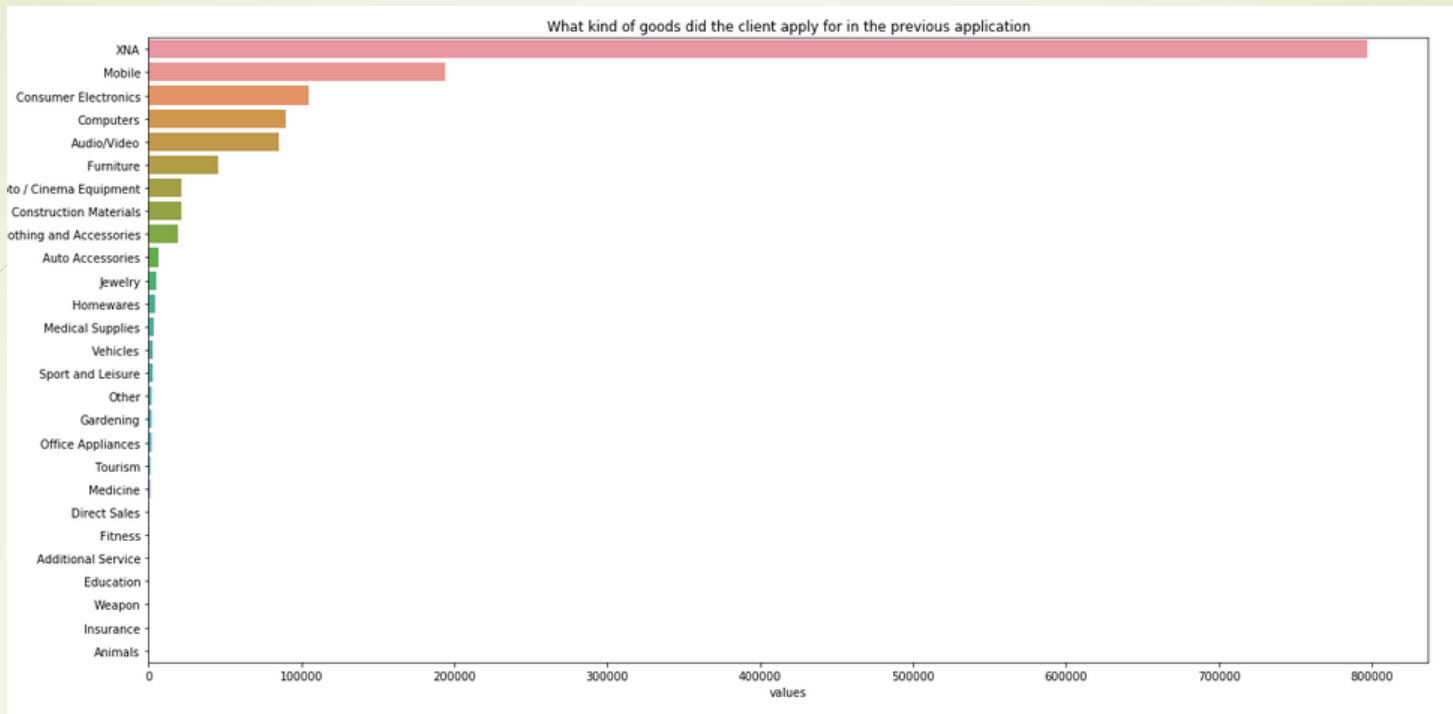
Payment method that client chose to pay for the previous application

■ We could see that most of the payment in previous application is done by cash through bank. Which tells us that people prefer the traditional method of paying by cash even the online method is available.
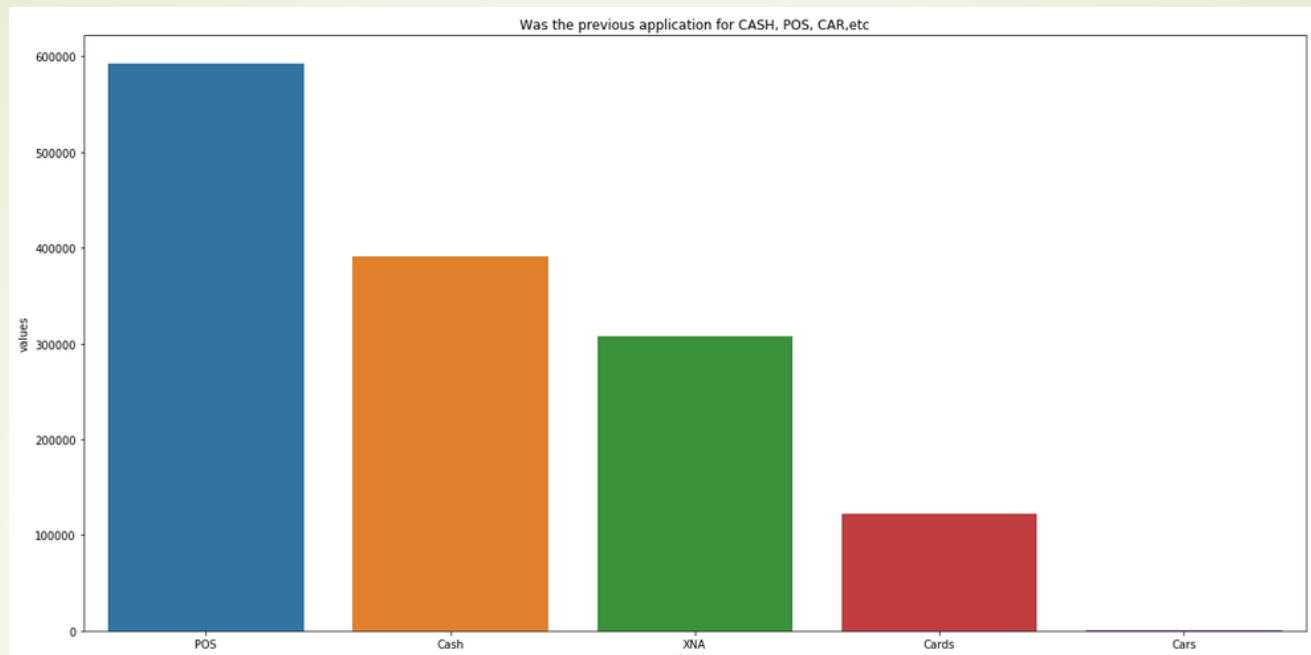
- The most of the previous application got rejected due in 'XAP' process, followed by 'HC' and the least rejection in 'SYSTEM'.
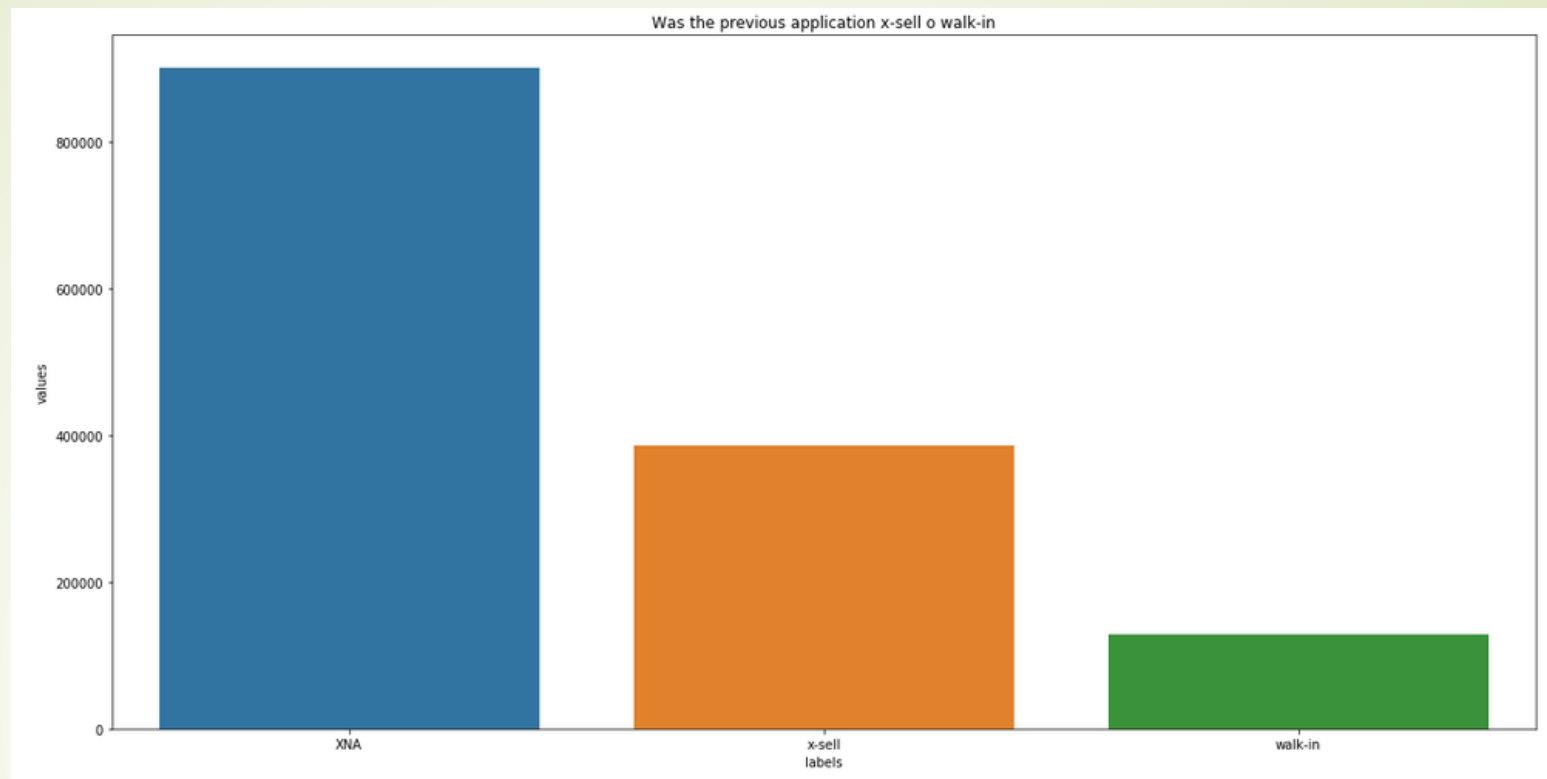
Most of the client in previous application was Repeater, which told us that their trust on the company is good. Very few number of clients are nw compare to the repeater.

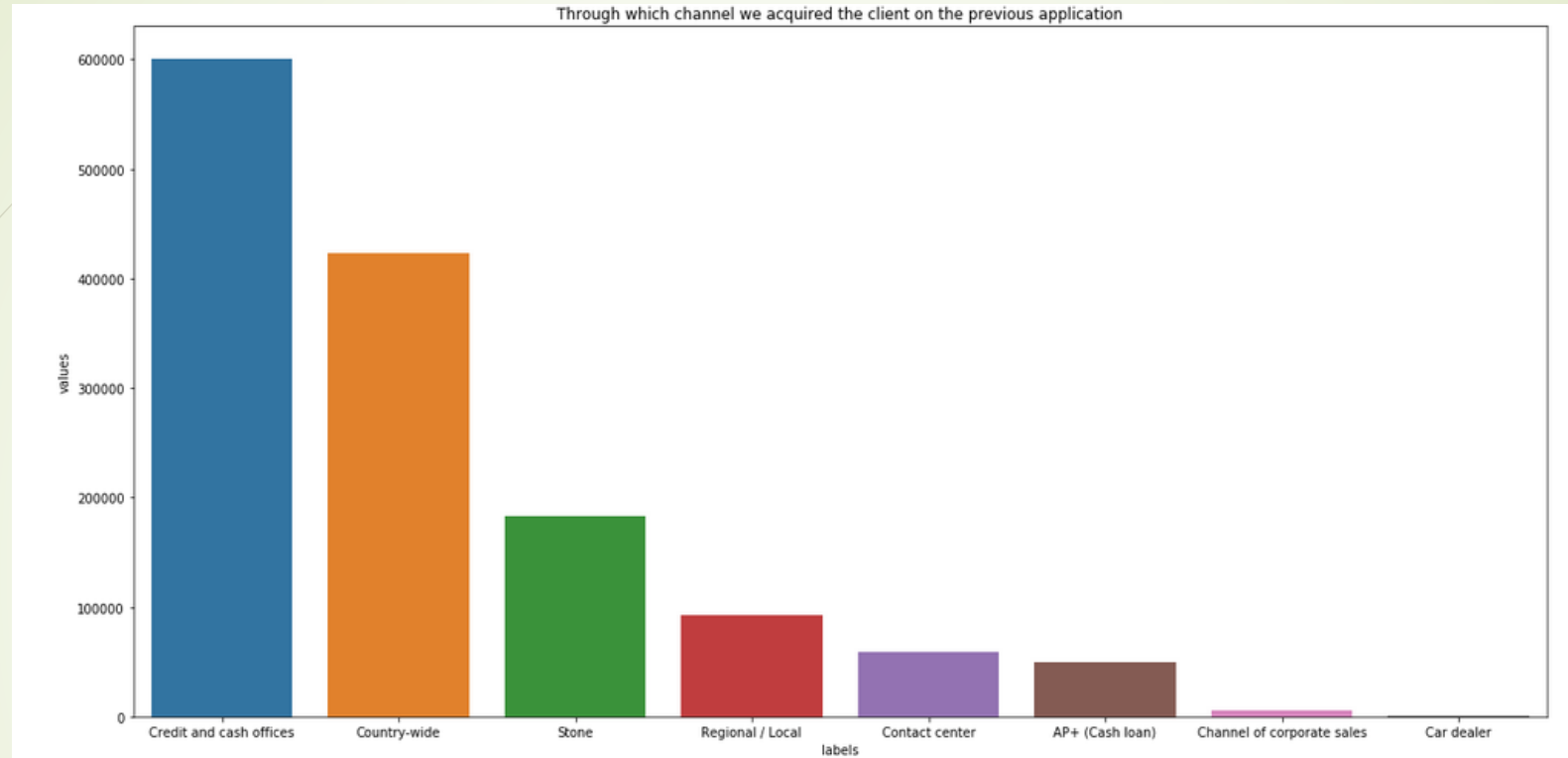What kind of goods did the client apply for in the previous application

- Most of the goods for which client took loan are for Mobile followed by Consumer electronics and Computers.

- This shows trend that most no of loans are taken for the electronics products and which tells us where bank market could flourish.
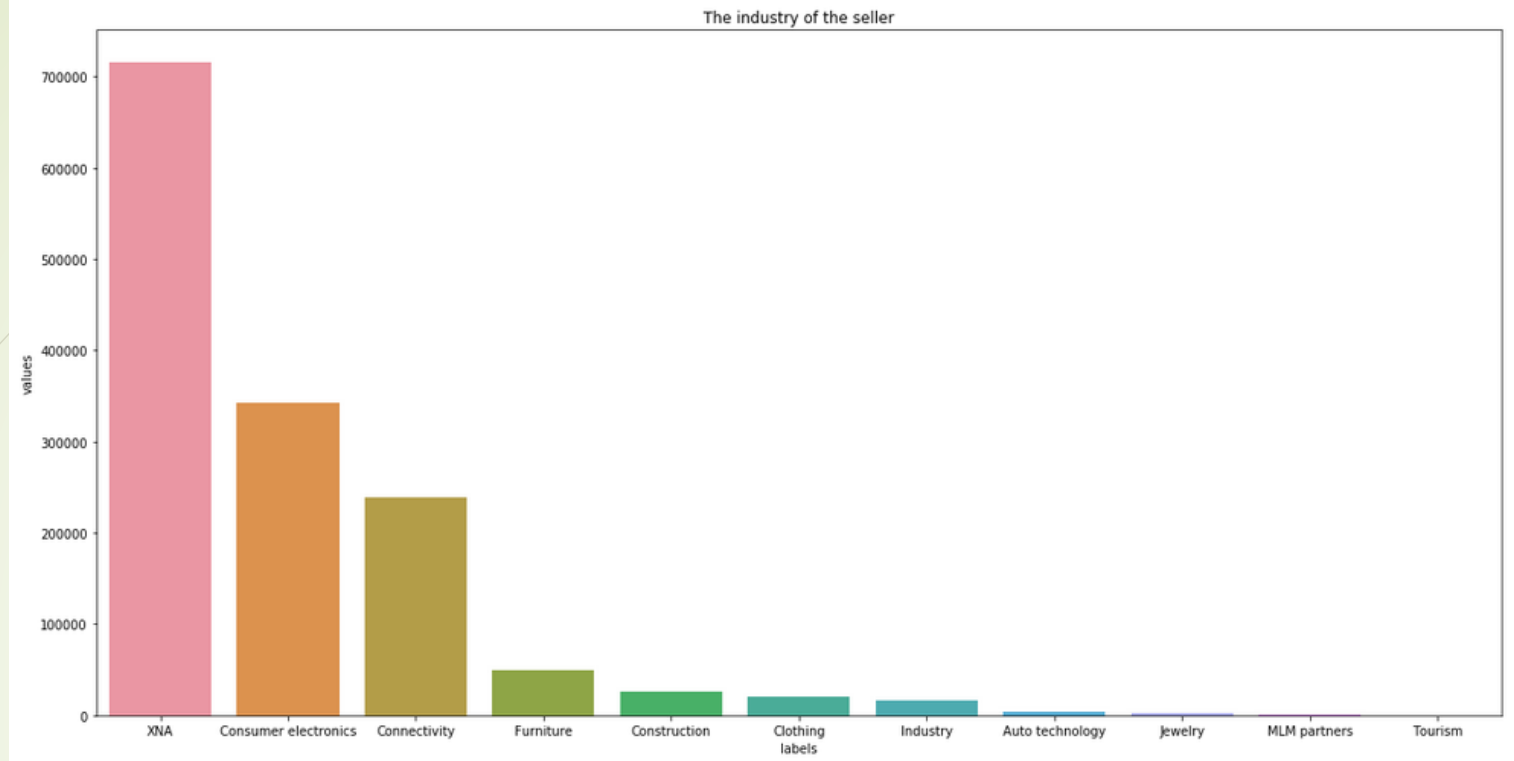
Was the previous application for CASH, POS, CAR,etc

- This tells us that most of the previous application are for 'POS' followed by 'Cash'.

- This tells us nearly twice the application of previous application are for 'x-sell' than of 'walk-in'

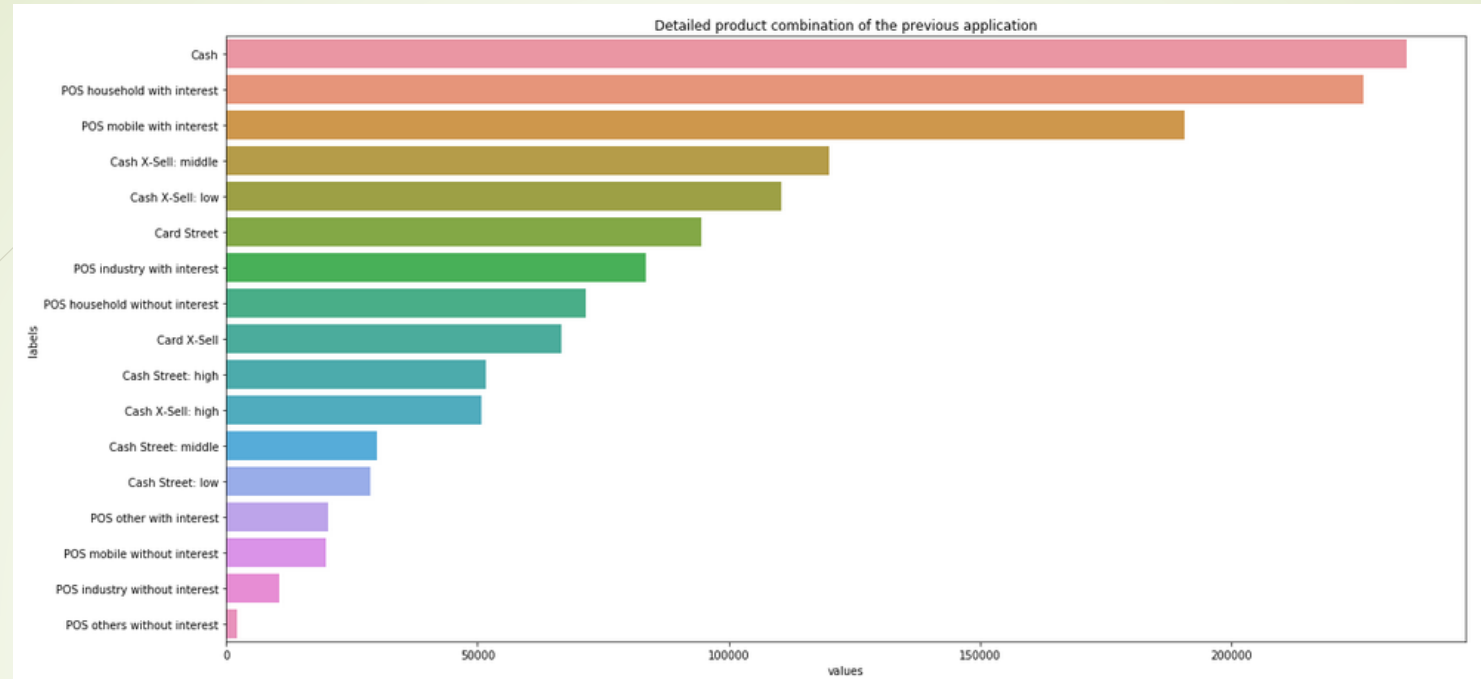Through which channel we acquired the client on the previous application

- Most of the client for previous application are acquired through 'Credit and Cash offices' followed by 'Country-wide' and least from 'Car dealer'.
- This tells us where to focus on in future for the market.
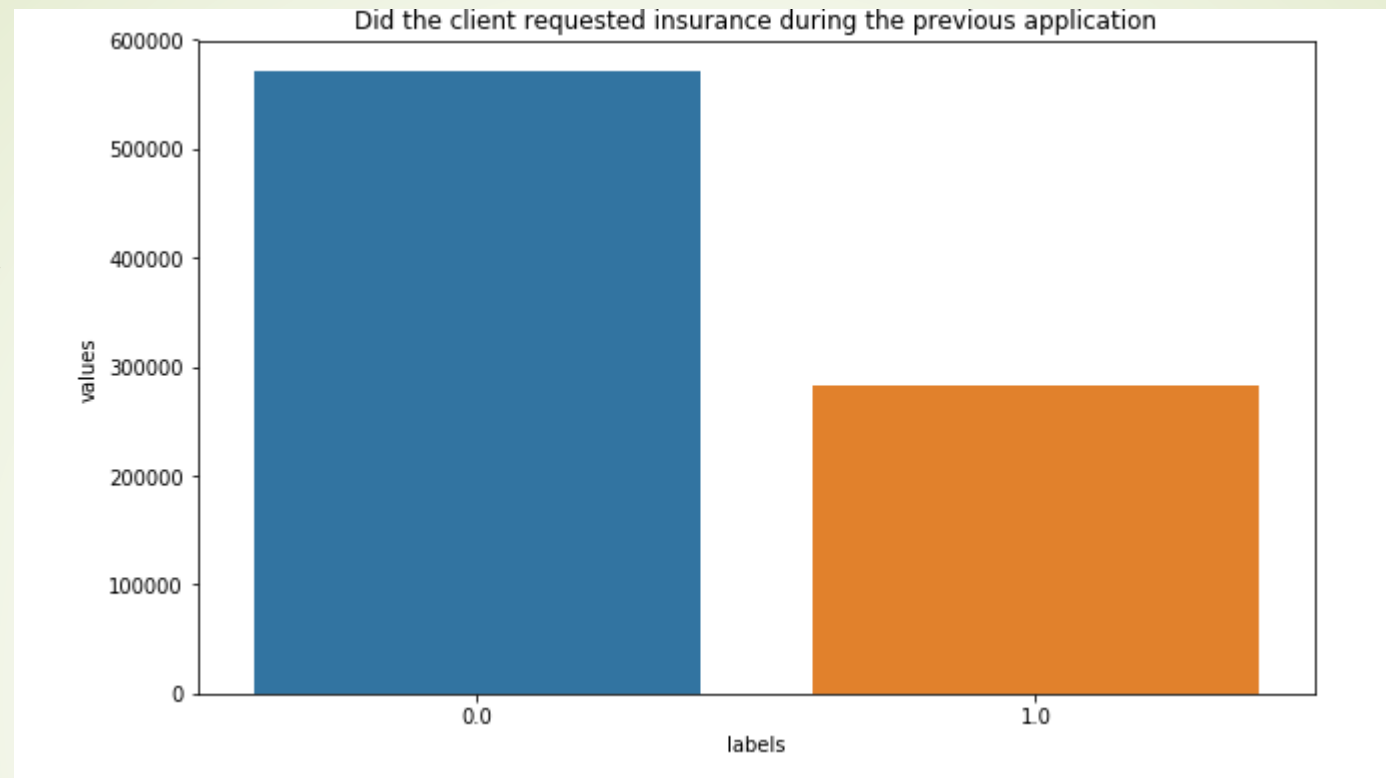
The industry of the seller

- In previous application most of the clients industry was 'Consumer electronics' followed by 'Connectivity' and the least one was 'Tourism'.
- This tells us in which industry we should focus more to get the more Clients.
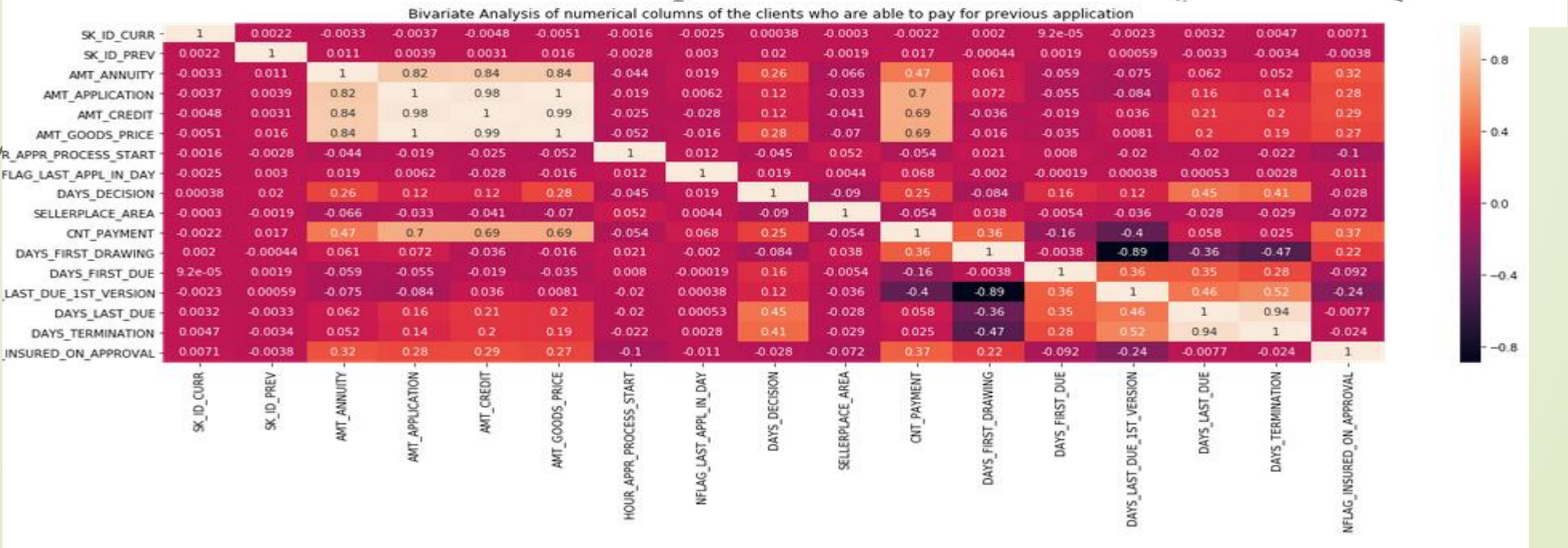
Grouped interest rate into small medium and high of the previous application
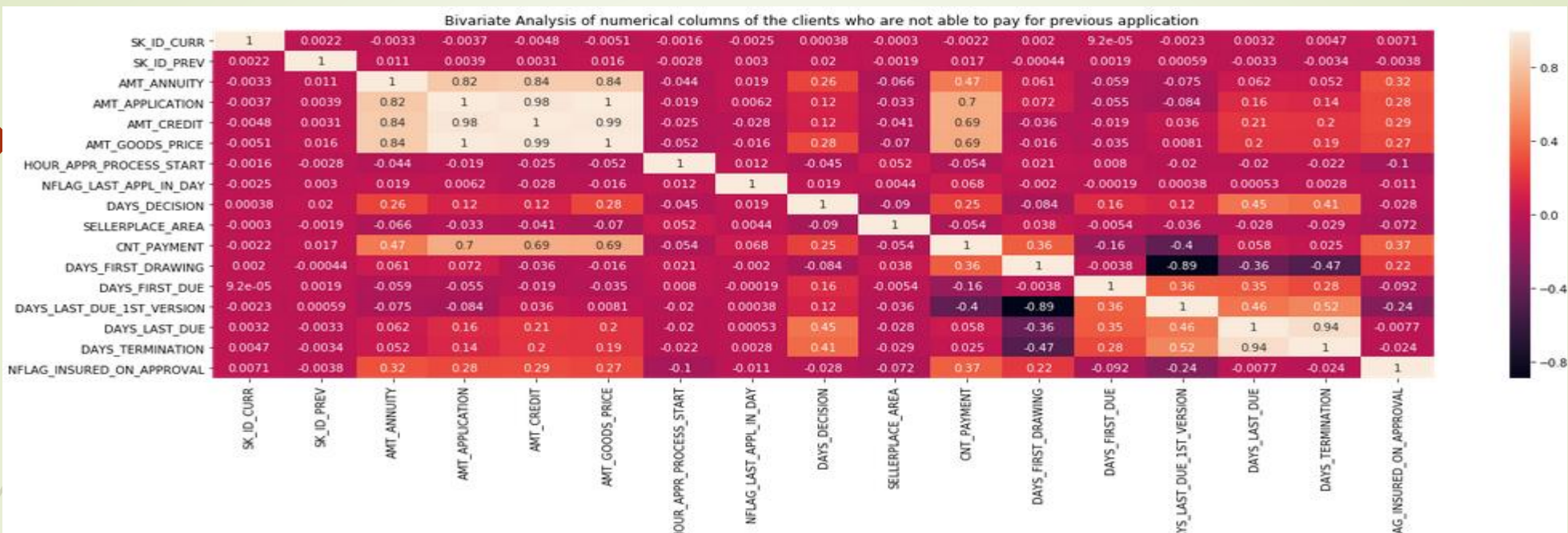
- Most the prefious applicaiton has medium interest rate followed by high interst rate, and the least one was low_action interest rate.

Detailed product combination of the previous application
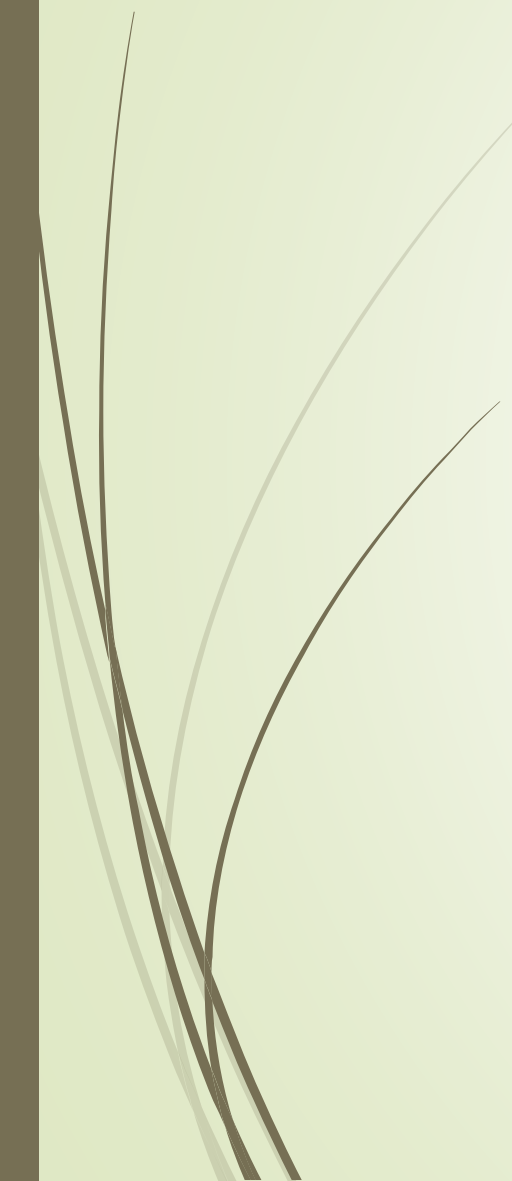
For the Detailed product combination of the previous application the most of it is in the form of 'cash' followed by 'POS household with interest'.

Did the client requested insurance during the previous application

- Nearly twice the number of Client do not required insurance during the previous application.

- This tells us approximately how many client would required the insurance now.

Bivariate Analysis of numerical columns of the clients who are not able to pay for previous application


Bivariate Analysis of numerical columns of the clients who are able to pay for previous application

- Here we could see that AMT_ANNUITY and AMT_CREDIT are highly correlated, if one goes up so s the other and vise versa.

- For both the sets but much more correlated for the Clients who are able to pay.

- We could also see that AMT_CREDIT and AMT_GOODS_PRICE are also highly correlated, if one goes up so s the other and vise versa.

- For both the sets but are more correlated to the Clients who are able to pay.

- Similarly AMT_ANNUITY and AMT_GOODS_PRICE are also highly correlated, if one goes up so s the other and vise versa.

- For both the sets but are more correlated to the Clients who are able to pay.

- For AMT_ANNUITY and AMT_INCOME_TOTAL the correlation is much higher for the Clients who are able to pay.

- Similar is the case with AMT_CREDIT and AMT_INCOME_TOAL.

- We can see that there are group of value have a similar correlation AMT_ANNUITY,AMT_APPLICATION,AMT_CREDIT,AMT_GOODS_PRICE for a segment of values with highly dependent factors.

- In the similar manner there are a lot for columns which are corrected for the case the Clients who are able to pay and the Clients who are not able to pay.

- DAYS_LAST_DUE_1ST_VERSION and DAYS_FIRST_DRAWING are highly correlated negatively. If one goes up other goes down and vise versa.

# THANK YOU