# Credit EDA Case Study

## Exploratory Data Analysis

Team
Nistha Kumar & Gaurav Rana

# Table of Content

**1** Understanding of the problem statement

**2** Overall EDA Approach

**3** Data Quality Check

**4** Data Analysis

**5** EDA Summary (Key Patterns)
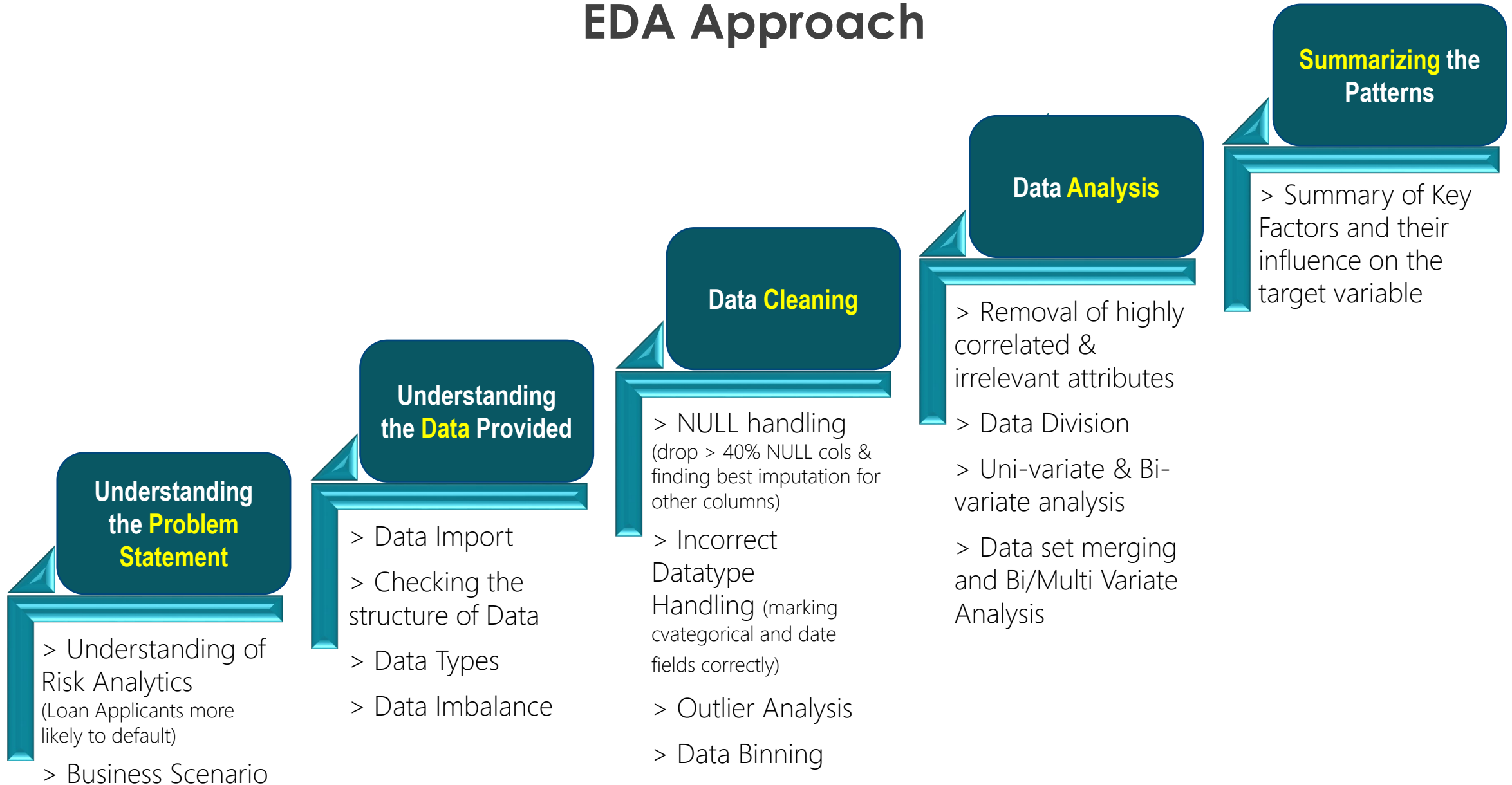
# Problem Statement

## Business Definition

→ A "consumer finance company" wants to minimise the risk of losing money while lending to customers

→ There are two risks:

- 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

## Data Project definition

→ "Target" variable is what tells whether a customer has defaulted or not

→ Find all the "applicant", "loan", "previous application" related variables that have a correlation/influence on the target variable (strong indicators of default)

→ Provide a summary of all the variables that have a trend against the target variable, so that the company can utilise this knowledge for its portfolio and risk assessment

| application_data | | | previous_application | |
|---|---|---|---|---|
| key | SK_ID_CURR | <-- join--> | key | SK_ID_CURR |
| target variable | TARGET | | target variable | NAME_CONTRACT_STATUS |
| attribute1 | | | attribute1 | |
| attribute2 | | | attribute2 | |
| attribute3 | | | attribute3 | |
| … | | | … | |
| … | | | … | |

# Overall EDA Approach

**Understanding the Problem Statement**

> Understanding of Risk Analytics (Loan Applicants more likely to default)

> Business Scenario

**Understanding the Data Provided**

> Data Import

> Checking the structure of Data

> Data Types

> Data Imbalance

**Data Cleaning**

> NULL handling (drop > 40% NULL cols & finding best imputation for other columns)

> Incorrect Datatype Handling (marking cvategorical and date fields correctly)

> Outlier Analysis

> Data Binning

**Data Analysis**

> Removal of highly correlated & irrelevant attributes

> Data Division

> Uni-variate & Bi-variate analysis

> Data set merging and Bi/Multi Variate Analysis

**Summarizing the Patterns**

> Summary of Key Factors and their influence on the target variable

# Data Quality check (1/2)

| | NULL analysis | NULL Imputation | Datatype Correction | Date conversion | Outlier Analysis | Binning |
|---|---|---|---|---|---|---|
| **Approach** | • Find % age of NULLs in every column<br>• Decided on a threshold that any column with more than 40% values as NULLs would be dropped | • For rest of the columns (NULL %age between 1% and 40%) – analysis needed on whether and how to impute<br>• For Categorical : Mode imputation<br>• For Continuous: Median or Mean | • Checking data types - all categorical columns (based on unique values) should be of object data type<br>• All continuous columns should be integers/float) | • There were many dates with negative values<br>• They were converted into datetime formats by adding the negative number to 01/01/2020 as a reference date. | • Found the topmost 5 most skewed columns for Outlier Analysis<br>• Boxplot and distplot analysis on whether outliers need to be removed, binned or retained | • Binning some variables to make the continuous variable to categorical. |
| **Application Data** | • 49 columns found with > 40% NULLs and hence dropped | 1 Categorical to be imputed by MODE (OCCUPATION_TYPE)<br>1 Continuous : imputed by MEAN (EXT_SOURCE_3)<br>6 continuous : imputed by MEADIAN (AMT_REQ*) | • Some columns converted to object type: TARGET, FLAG_MOBIL, FLAG_EMAIL, FLAG_DOCUMENT_* | • Columns converted : ['DAYS_BIRTH','DAYS_EMPLOYED','DAYS_REGISTRATION','DAYS_ID_PUBLISH','DAYS_LAST_PHONE_CHANGE'] | > OBS_30/60_CNT_SOCIAL_CIRCLE : values over 24.0 should be deleted (99.99 percentile value)<br>> AMT_REQ_CREDIT_BUREAU_H/D: have outliers but there is possibility of that data, no action<br>> AMT_REQ_CREDIT_BUREAU_QRT : values over 8 should be deleted. (99.999 percentile value) | 1.AMT_INCOME_TOTAL : binned it into :"Very Low", "Low", "Medium", "High", "Very High" categories.<br>2. HOUR_APPR_PROCESS_START_CAT binned into "Morning", "Afternoon", "Evening" |
| **Previous Applications** | • 11 columns found with > 40% NULLs and hence dropped | • Three columns to be imputed by median:<br>• CNT_PAYMENT AMT_ANNUITY AMT_GOODS_PRICE | • One column converted to Object type: NFLAG_LAST_APPL_IN_DAY | • Column converted : DAYS_DECISION | ➢ CNT_PAYMENT : no need to handle outliers, there is a possibility of this data<br>➢ SELLERPLACE_AREA : values over 120000 can be deleted. 99.999% values are within this range<br>➢ DAYS_DECISION : values over 2913 should be deleted. 99.90% values are within this range | 1. HOUR_APPR_PROCESS_START_CAT binned into "Morning", "Afternoon", "Evening" |

# Data Quality check (2/2)

Some Sample work done for Data Quality
(refer to notebooks attached for details)

Analyzing distributions & skewness to decide whether MEAN or MEDIAN statistic to be used

Binning data to handle outliers

NULL Imputation for Categorical Variables through MODE

Looking at the categorical variables and correcting datatypes

Conversion to Correct Datetime formats
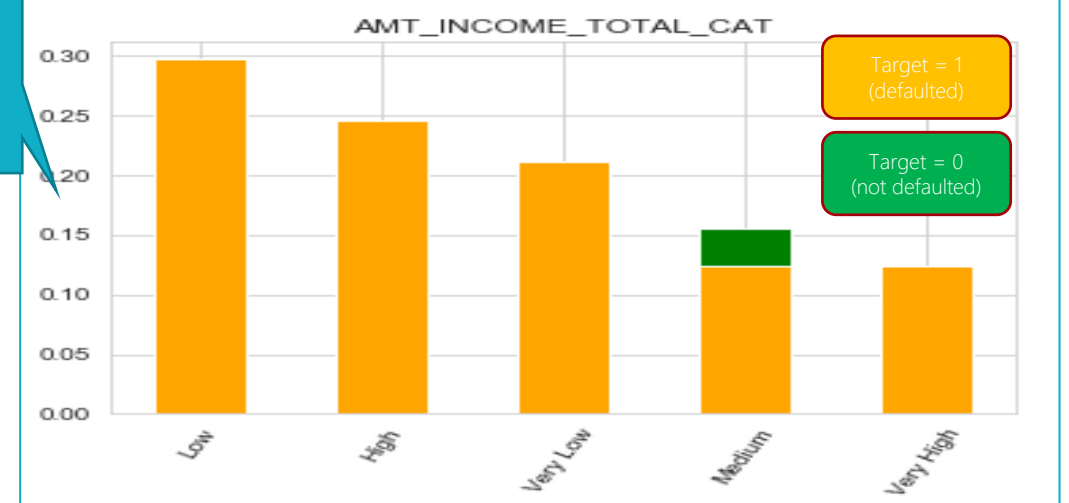
Boxplots for outlier analysis and approach to treat

## Categorical Variables influencing Target Variable



More Pensioners are less likely to default whereas Working applicants are expected to default more

Maximum people not defaulting are of Higher education type. So people with Higher Education can be considered for approval.

Applicants with Medium Income Category are most likely to not default.

Applicants are mostly staying in region with rating 2 for both target 0 and 1. Also people staying in region 3 are more likely to default compared to region 1
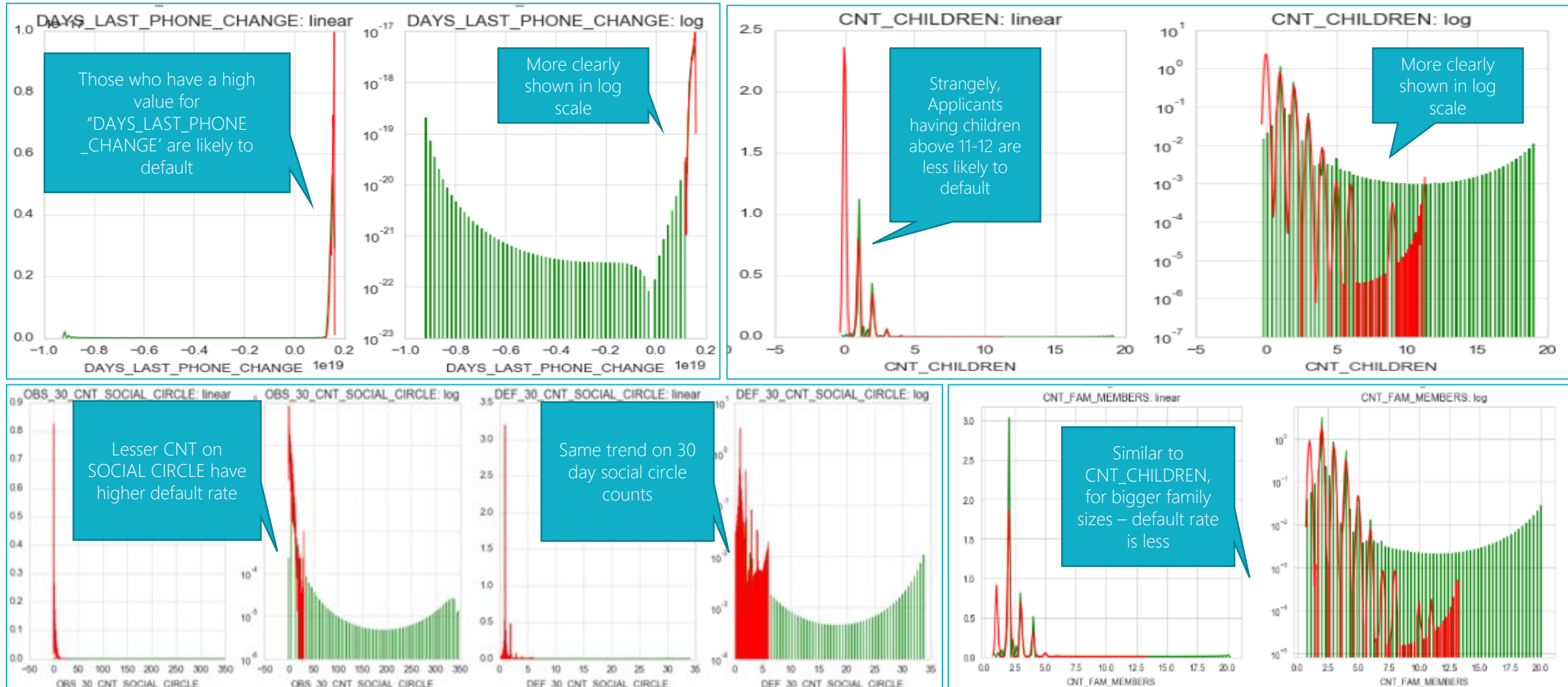
Target = 1 (defaulted)

Target = 0 (not defaulted)

7

Target = 1 (defaulted)  Target = 0 (not defaulted)

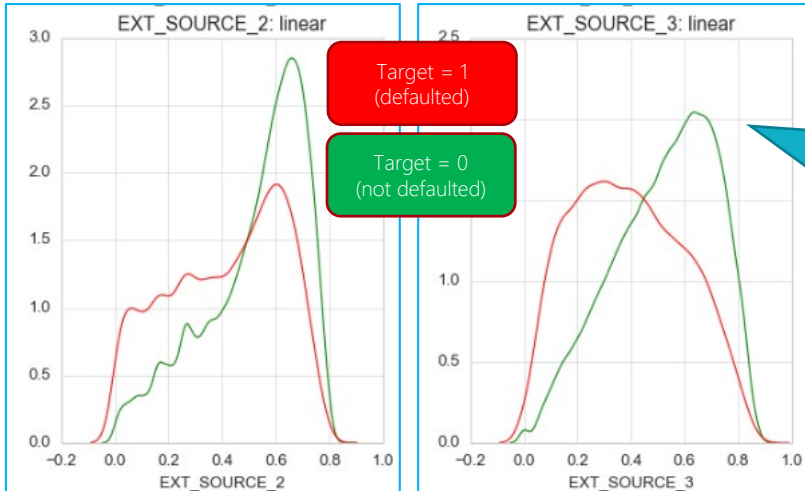## Continuous Variables influencing Target Variable (both liner and logarithmic trend shown)



Those who have a high value for "DAYS_LAST_PHONE _CHANGE' are likely to default

More clearly shown in log scale

Strangely, Applicants having children above 11-12 are less likely to default

More clearly shown in log scale

Lesser CNT on SOCIAL CIRCLE have higher default rate

Same trend on 30 day social circle counts

Similar to CNT_CHILDREN, for bigger family sizes – default rate is less

## Continuous Variables influencing Target Variable (contd.)



Target = 1 (defaulted)

Target = 0 (not defaulted)

Both SOURCE 2 & SOURCE 3 have influence – higher values have lesser default rate

Slide # 7 shows working people have more chances to default. Here we see that working people who donot default (left) have higher AMT_CREDIT values than one ones who do(right)

If you see the PAIRPLOT analysis, combination of AMT_CREDIT and CNT_CHILDREN has an impact i.e.
Lower children & higher credit
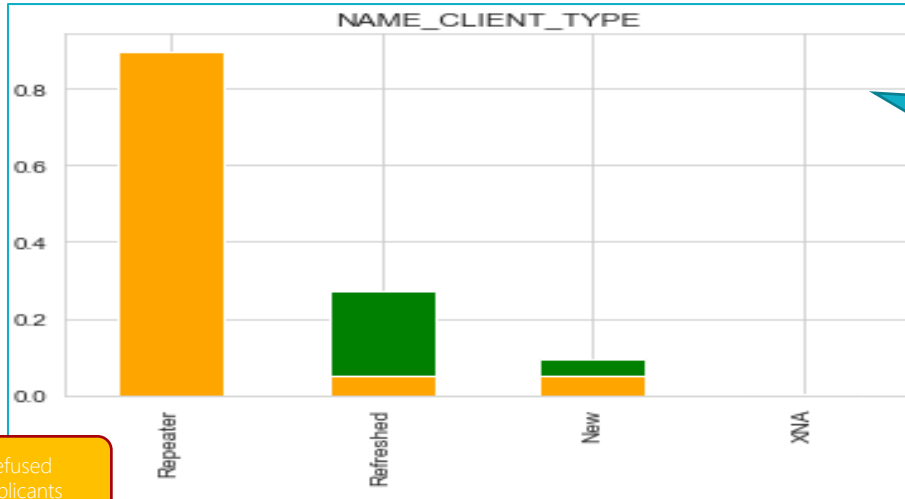Or higher children & lower credit – both have a tendency to default
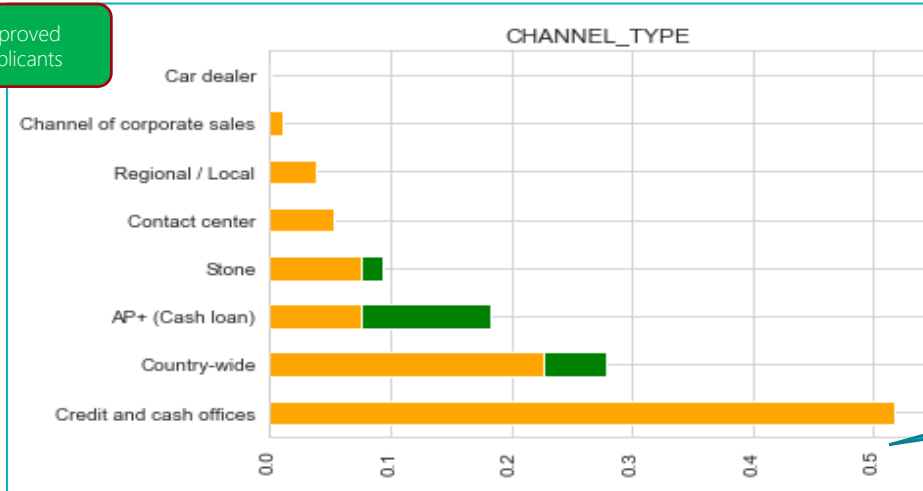
9

# Data Analysis
# Previous Applications (1/2)

## Categorical Variables influencing NAME_CONTRACT_STATUS Variable

> Applicants with portfolio status Cash are more likely to be refused for the loan whereas those of type POS are mostly approved.



NAME_CLIENT_TYPE
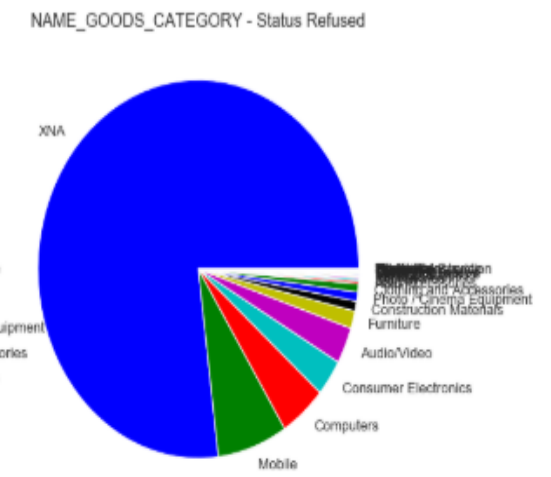
Clients of type "Refreshed" have most trend of being approved
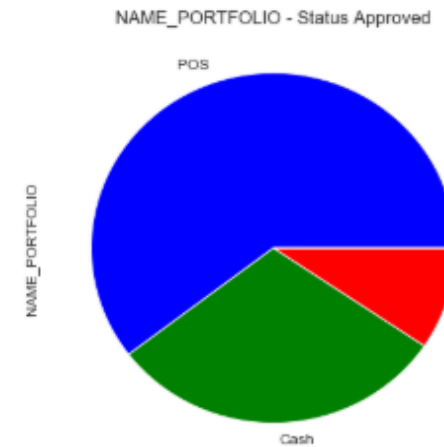
Refused Applicants

Approved Applicants

> Chances of Loan getting rejected are more if the applicant dosent specify the Goods_Category(XNA).
> Loans for Goods_Category Mobile are maximum for both approved and rejected loans

NAME_PORTFOLIO - Status Approved

NAME_PORTFOLIO - Status Refused

CHANNEL_TYPE

Channel type AP+(Cash Loan) have the highest value of loans getting approved. So this can be a good parameter for the bank to decide on the loan status
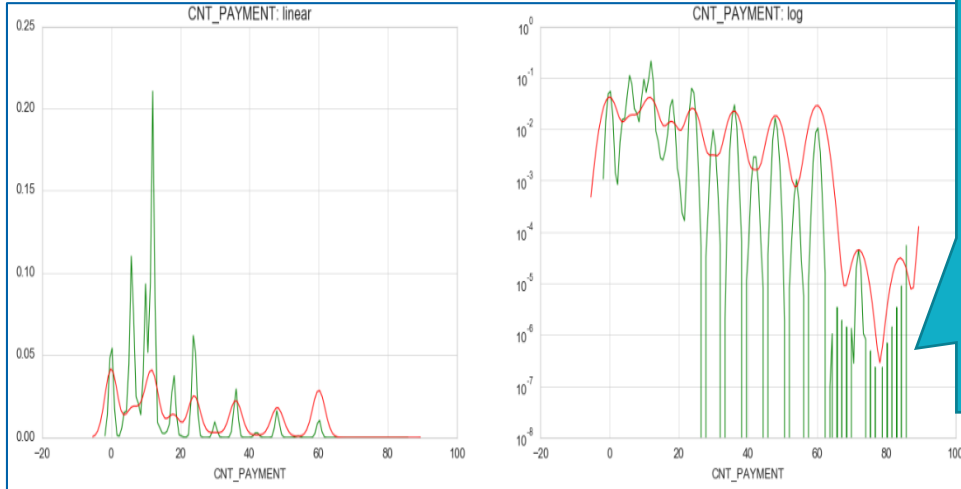
NAME_GOODS_CATEGORY - Status Approved

NAME_GOODS_CATEGORY - Status Refused

# Data Analysis
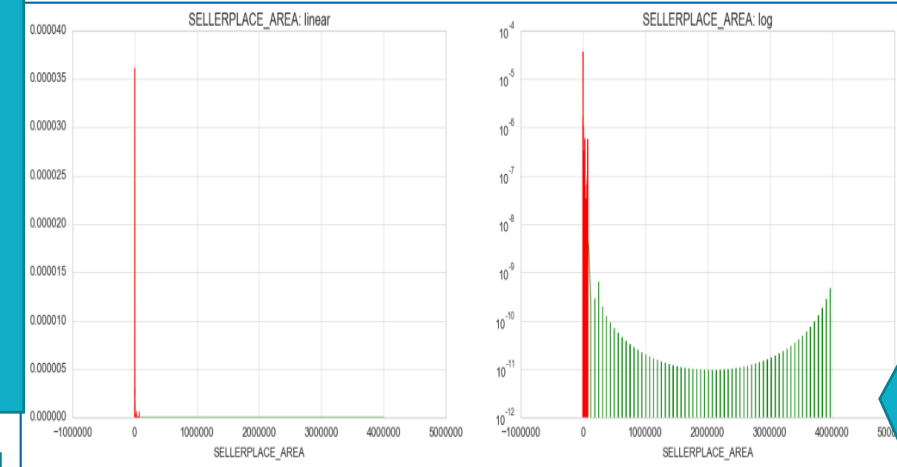# Previous Applications (2/2)

Refused Applications

Approved Applications

## Continuous Variables influencing NAME_CONTRACT_STATUS Variable



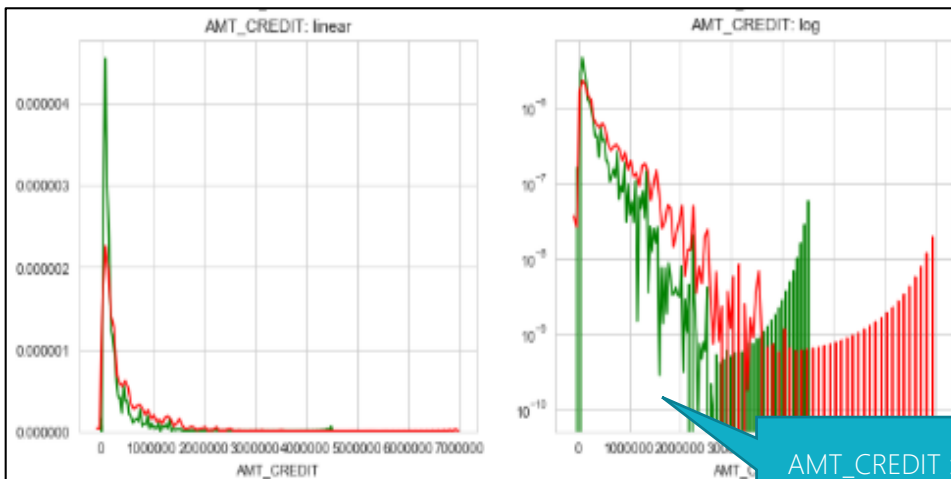CNT_PAYMENT : Term of previous credit at application of the previous application

➢ Maximum approved previous applications are there where the term is between 10-14

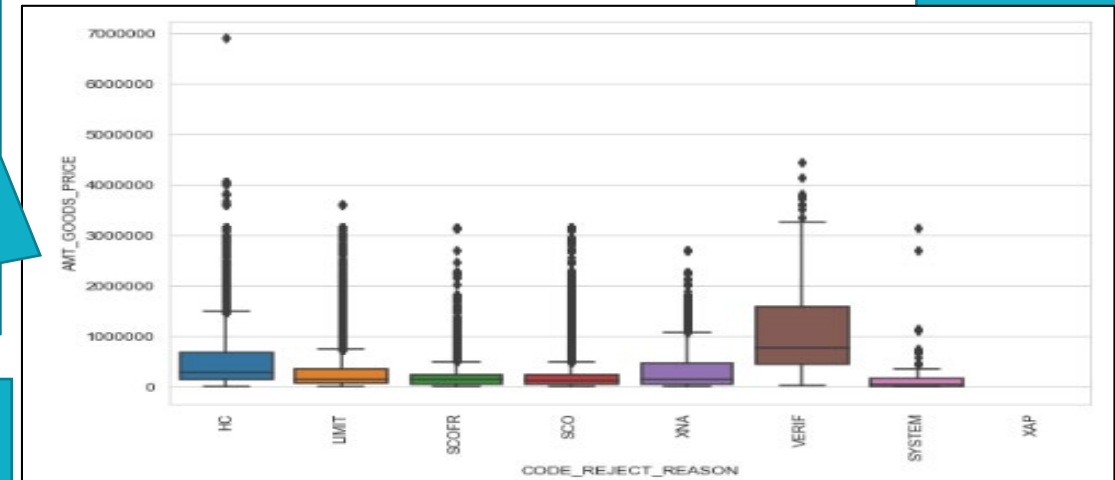SELLERPLACE_AREA : this is "Selling area of seller place of the previous application"

As can be seen (graph on the left is liner view and on the right is the logarithmic view)

➢ Records with Lower SELLERPLACE_AREA have a high REFUSAL RATE

CODE_REJECT_REASON : Maximum of the rejected loans are because of the reason: "HC" but "VERIF" have the highest AMT_GOODS_PRICE and high AMT_GOODS_PRICE results in more rejection of the loan.

AMT_CREDIT : Maximum rejected loans are of Higher AMT_CREDIT

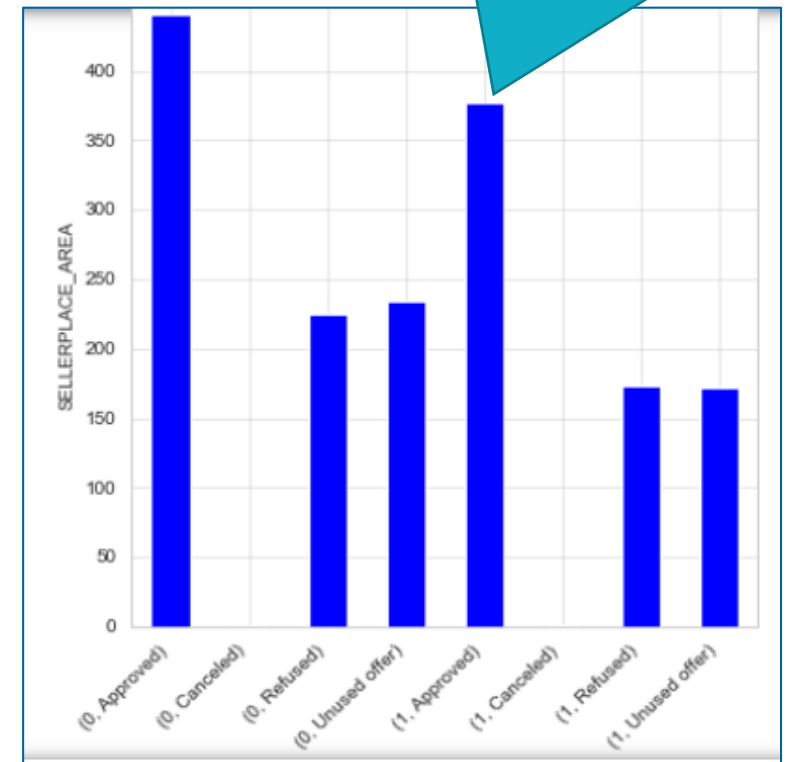## Impact of variables – on a combination of "TARGET" and "NAME_CONTRACT_STATUS"

Note: first 4 bars show 'non defaults' – different statuses, and last 4 bars show "default" – different statuses. Bar height is the MEAN value

➢ EXT source 2 and 3 have an impact on the default rate.
➢ As can be seen, default cases have consistently lower values for both these variables

➢ As seen from this analysis also, Lower SELLERPLACE_AREA has a higher tendency for default
(comparison between approval – default and non-default cases)
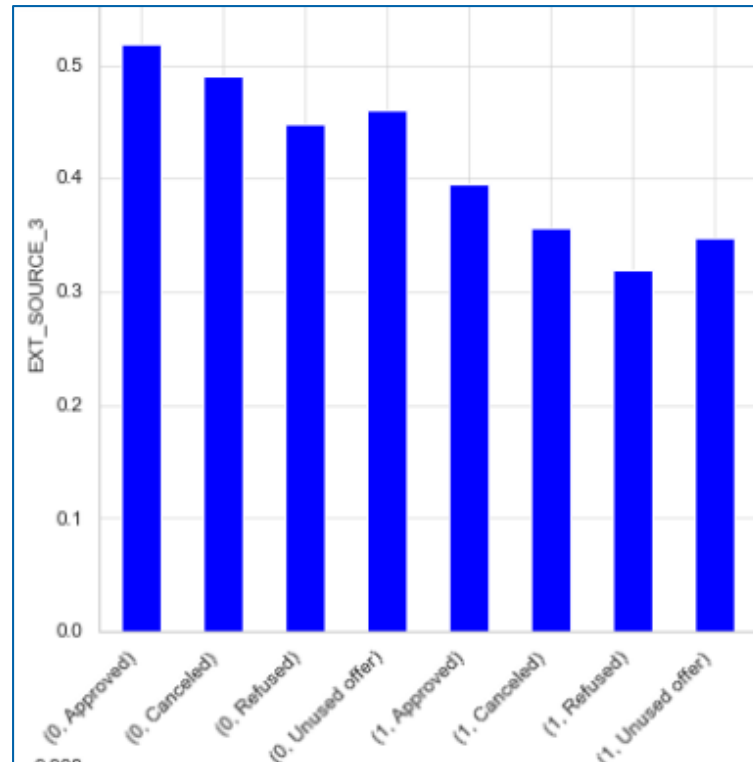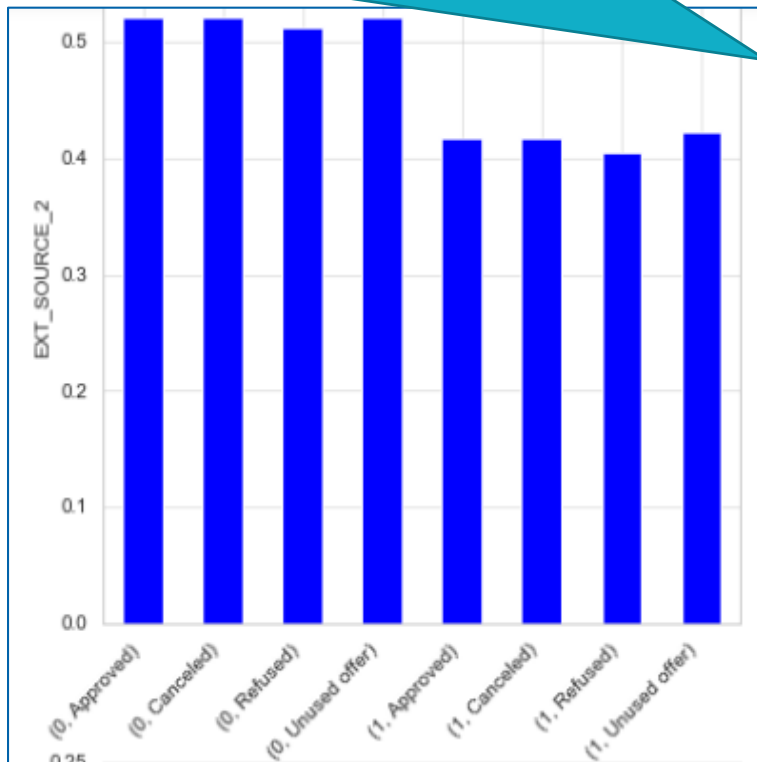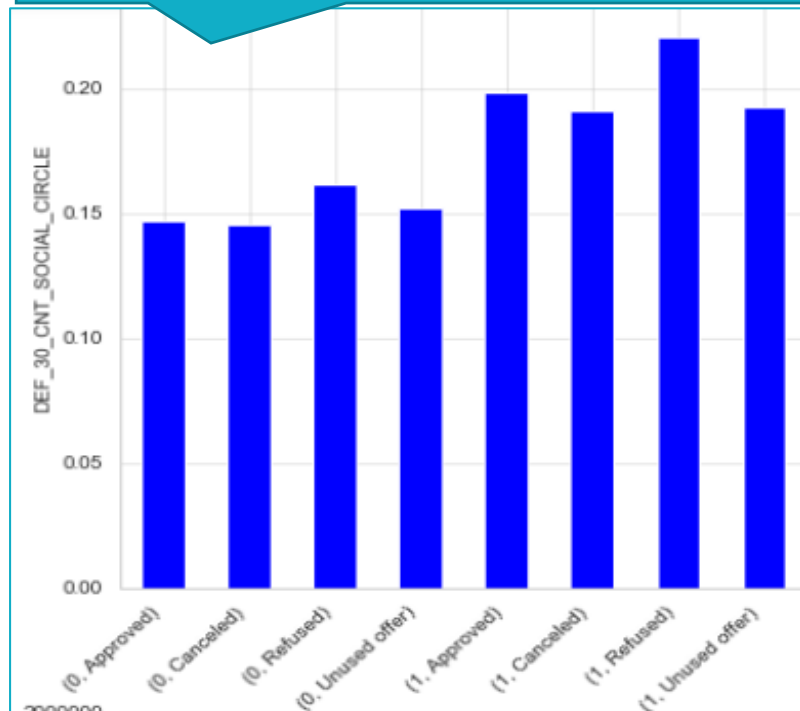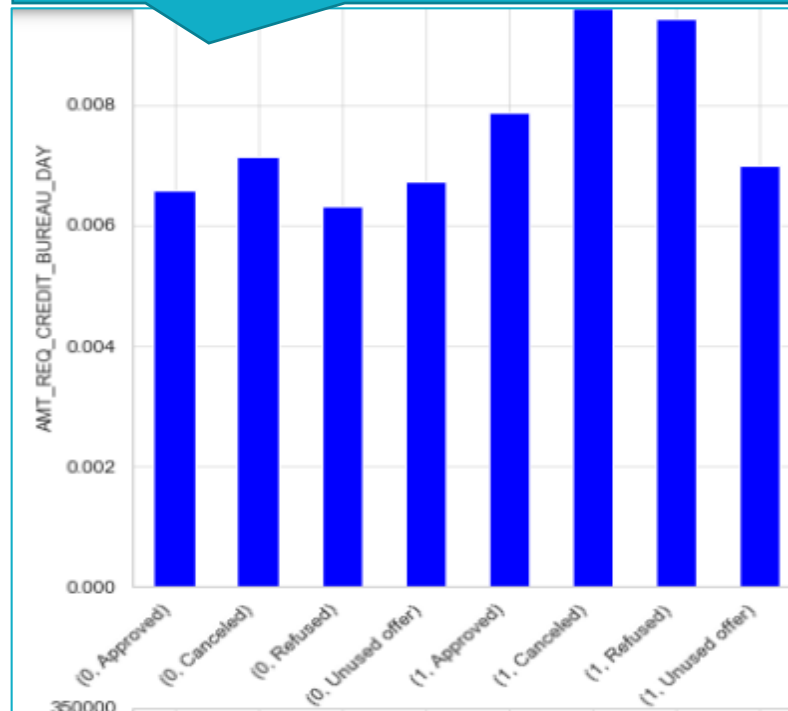
## Impact of variables – on a combination of "TARGET" and "NAME_CONTRACT_STATUS"

Note: first 4 bars show 'non defaults' – different statuses, and last 4 bars show "default" – different statuses. Bar height is the MEAN value

➤ Social status count (def 30) – mean value for non defaults is lower than mean value for defaults

➤ It seems when NUMBER of enquiries to credit bureau is higher, then default chances are also higher

➤ Status to status comparison : mean of number of children is higher for defaulters than non-defaulters



Columns impacting the bank's decision:
AMT_CREDIT, AMT_ANNUITY, AMT_APPLICATION, AMT_GOODS_PRICE, CNT_FAMILY_MEMBERS, CNT_CHILDREN, OBS_60_CNT_SOCIAL_CIRCLE, OBS_30_CNT_SOCIAL_CIRCLE, DEF_30_CNT_SOCIAL_CIRCLE, DEF_60_CNT_SOCIAL_CIRCLE, NAME_EDUCATION_TYPE, NAME_INCOME_TYPE, REGION_RATING_CLIENT, AMT_INCOME_TOTAL, DAYS_LAST_PHONE_CHANGE, EXT_SOURCE_2, EXT_SOURCE_3, SELLERPLACE_AREA, AMT_REQ_CREDIT_BUREAU_DAY, FLAG_OWN_CAR, FLAG_OWN_REALTY

13

# Final EDA Summary



**Tendency to Default** (central node)

Connected nodes:
- Education
- Income Type
- Credit bureau enquires - Day
- Region Rating Client
- Seller place area
- Total Income Category
- Ext source 2 and 3
- Days last phone change
- Social circle counts
- Children count

Callout notes:
- Maximum people not defaulting are of Higher education type
- Working people are expected to default more than Pensioners
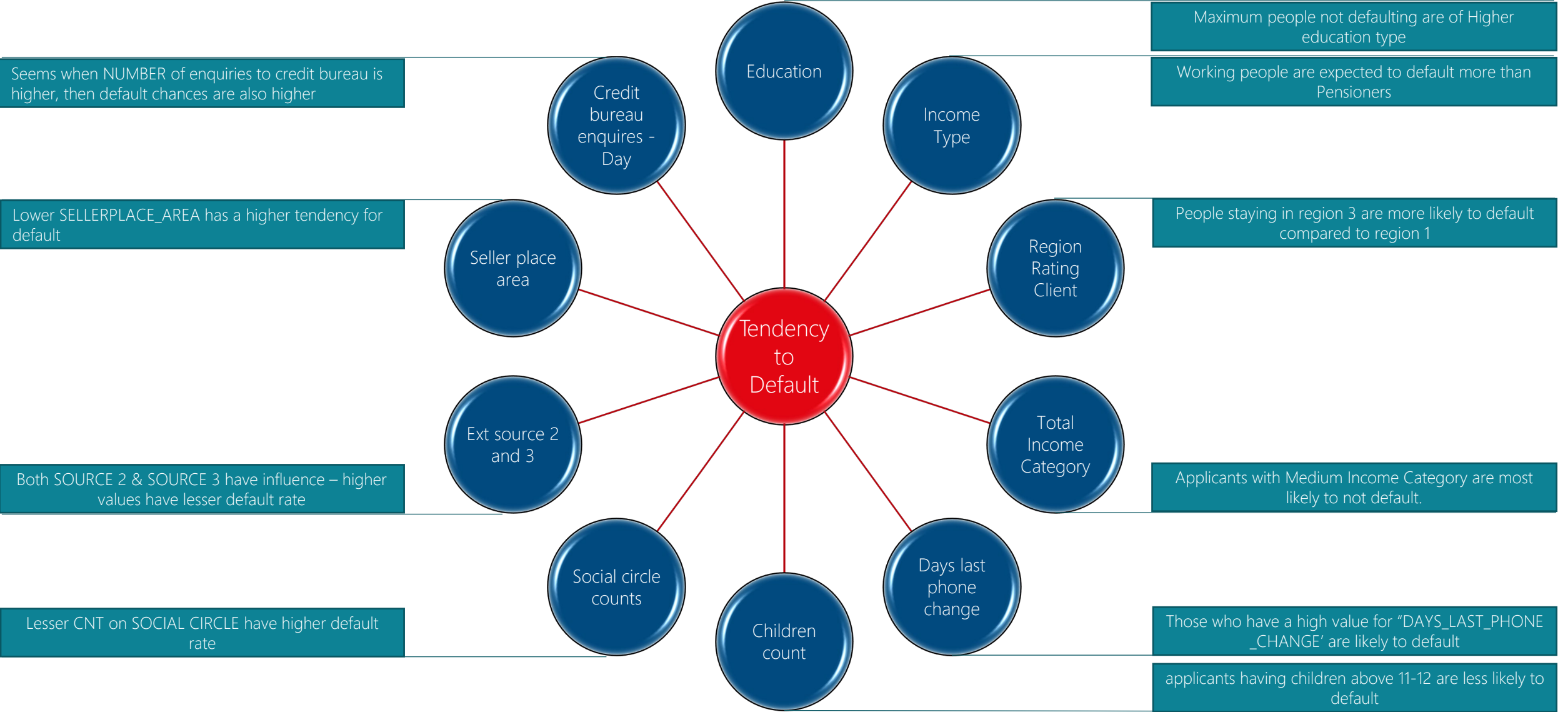- Seems when NUMBER of enquiries to credit bureau is higher, then default chances are also higher
- Lower SELLERPLACE_AREA has a higher tendency for default
- People staying in region 3 are more likely to default compared to region 1
- Both SOURCE 2 & SOURCE 3 have influence – higher values have lesser default rate
- Applicants with Medium Income Category are most likely to not default.
- Lesser CNT on SOCIAL CIRCLE have higher default rate
- Those who have a high value for "DAYS_LAST_PHONE _CHANGE' are likely to default
- applicants having children above 11-12 are less likely to default

# Thank You