# Current Application Data

## Importing data
- Shape – 307511 x 122

## Preparing for Analysis - I
- Dropping Columns with missing percentage greater than 50 % - 41 columns
- 6 columns had around 13 % missing values and were not imputed for analysis as they were critical values and imputing them with 0(the mean, mode and median) would give a very wrong input
- Converting to right datatypes – Days column, Family Members and Region Rating W City

## Preparing for Analysis - II
- Binning of continuous numerical values like Income total, emi, credit, days birth
- Dividing the dataset into two based on TARGET variable
- Taking random sample of 5000 for two dataset for the analysis and now data is balanced for analysis

## Analysis
- Tried to analysis using Univariate and Bivariate approach using numerical and categorical columns to establish any pattern or trend with a specific feature

# Previous Application Data

## Preparing Data for Merge

- Cleaning the Previous Data – Dropping columns with missing greater than 50 %, converted the days column to absolute
- Grouped Columns with respect to **SK_ID_CURR**
- Selected the columns those are required by performing the median/mean operation on respect columns
- In columns having categorical values, used pd.crossTab and extracted their value counts as new numeric columns
- Merged the current and processed previous dataframe using inner join on SK_ID_CURR
- Shape - 291057 rows × 94 columns

## Preparing for Analysis - I

- Binning process to create categorical columns for better analysis
- Split the dataset into two dataset based on TARGET value
- Took random sample of 5000 for analysis

## Analysis

- Tried to analysis using Univariate and Bivariate approach using numerical and categorical columns to establish any pattern or trend with a specific feature of applicants previous loan application history to current loan difficulties in paying EMI