# BANK LOAN CASE STUDY

-ANDRI

# Project Description

The main aim of this project is to identify patterns that indicate if a customer will have difficulty paying their installments. This information can be used to make decisions such as denying the loan, reducing the amount of loan, or lending at a higher interest rate to risky applicants. The company wants to understand the key factors behind loan default so it can make better decisions about loan approval.

▶ **When a customer applies for a loan, a company faces two risks:**

1. If the applicant can repay the loan but is not approved, the company loses business.

2. If the applicant cannot repay the loan and is approved, the company faces a financial loss.

# Approach and tech stack used

▶ Understanding the data.

▶ Cleaning / pre-processing the data and handling all NULL values and outliers.

▶ Merging the datasets/csvs to gain deeper insights.

▶ Data analysis and visualization.

▶ Gaining insights/hypothesis from the analysis.

I have used MS-Office 2019 for analysis and powerpoint presentation.
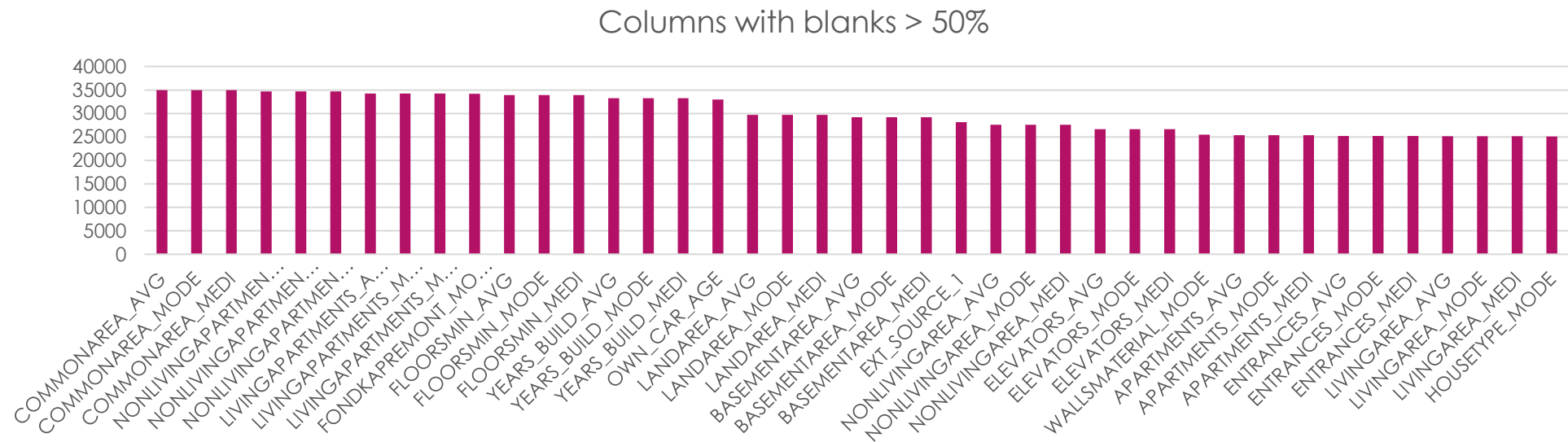
# Understanding the dataset

▶ Our dataset consists of three .csv files providing information of various aspects of the loan application mentioned below:

1. application_data.csv: gives information about current loan applications.

2. previous_application.csv: has information about client's previous loan data.

3. columns_description.csv: contains information about columns present in the above two datas.

# Application data file

▶ Rows: 50000

▶ Columns: 122

▶ Columns with blank cells: 67

▶ Columns with high number of blank cells (>50% of total rows): 41

▶ Columns with less number of blank cells (<50% of total rows): 26s

# Data cleaning

▶ Identifying the columns (41 in number) with a high number of blank cells will be **dropped** as imputation would not work on it.

Columns with blanks > 50%

# Data cleaning



Columns with blanks < 50%

# Data cleaning

▶ Out of the remaining 26 columns, certain columns are not required for analysis such as:

| FLOORSMAX_AVG |
|---|
| FLOORSMAX_MODE |
| FLOORSMAX_MEDI |
| YEARS_BEGINEXPLUATATION_AVG |
| YEARS_BEGINEXPLUATATION_MODE |
| YEARS_BEGINEXPLUATATION_MEDI |
| TOTALAREA_MODE |
| EMERGENCYSTATE_MODE |
| EXT_SOURCE_3 |
| EXT_SOURCE_2 |

▶ Dropping these 10 columns, we are left with 16 columns that need to be imputed with either **mean, median or mode for numerical** and **mode for categorical** columns.

# Data Cleaning

▶ Dropping the blank rows from the below columns as there is just one blank cell.

| |
|---|
| AMT_ANNUITY |
| CNT_FAM_MEMBERS |
| DAYS_LAST_PHONE_CHANGE |

▶ Features to impute finally 13 columns:

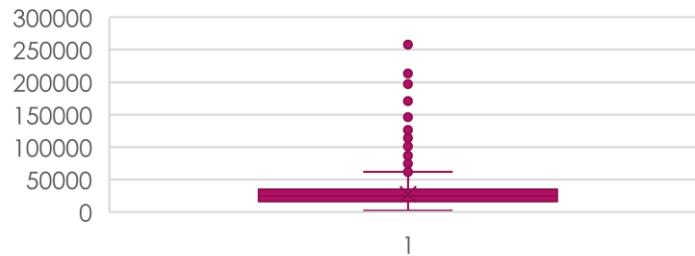| |
|---|
| AMT_REQ_CREDIT_BUREAU_HOUR |
| AMT_REQ_CREDIT_BUREAU_DAY |
| AMT_REQ_CREDIT_BUREAU_WEEK |
| AMT_REQ_CREDIT_BUREAU_MON |
| AMT_REQ_CREDIT_BUREAU_QRT |
| AMT_REQ_CREDIT_BUREAU_YEAR |
| NAME_TYPE_SUITE |
| OBS_30_CNT_SOCIAL_CIRCLE |
| DEF_30_CNT_SOCIAL_CIRCLE |
| OBS_60_CNT_SOCIAL_CIRCLE |
| DEF_60_CNT_SOCIAL_CIRCLE |
| OCCUPATION_TYPE |
| AMT_GOODS_PRICE |

# Mode Imputation:



NAME_TYPE_SUITE blank cells imputed with "Unaccompanied"



OCCUPATION_TYPE column has high number of blank cells (~15k) which is more than the highest occupation category hence imputed with "Unknown"

# Median imputation

- 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 : 1
- OBS_30_CNT_SOCIAL_CIRCLE : 0
- DEF_30_CNT_SOCIAL_CIRCLE : 0
- OBS_60_CNT_SOCIAL_CIRCLE : 0
- DEF_60_CNT_SOCIAL_CIRCLE : 0
- AMT_GOODS_PRICE : 450000

# Outliers removal



- Removing the outliers from the numerical columns
- CNT_CHILDREN in today's age can't be more than 3-4, in the dataset it goes up to 11, I have taken max 5 children.
- Removing vague values from DAYS_EMPLOYED **(~8k)** and AMT_INCOME_TOTAL that are too high, 1000 years for DAYS_EMPLOYED and 117000000 for AMT_INCOME_TOTAL.
- Removing the outliers from other numerical columns such as AMT_ANNUITY, AMT_GOODS_PRICE, AMT_CREDIT, DAYS_LAST_PHONE_CHANGE.

# Previous_application data file

- No. of columns: 37
- No. of rows: 50000
- No. of columns with blank: 15

**Number of blank cells in columns**

| Column | Blank cells |
|---|---|
| PRODUCT_COMBINATION | 8 |
| CNT_PAYMENT | 10592 |
| AMT_ANNUITY | 10592 |
| AMT_GOODS_PRICE | 10744 |
| NFLAG_INSURED_ON_APPROVAL | 19160 |
| DAYS_TERMINATION | 19160 |
| DAYS_LAST_DUE | 19160 |
| DAYS_LAST_DUE_1ST_VERSION | 19160 |
| DAYS_FIRST_DUE | 19160 |
| DAYS_FIRST_DRAWING | 19160 |
| NAME_TYPE_SUITE | 24243 |
| RATE_DOWN_PAYMENT | 25198 |
| AMT_DOWN_PAYMENT | 25198 |
| RATE_INTEREST_PRIVILEGED | 49834 |
| RATE_INTEREST_PRIMARY | 49834 |

# Data cleaning

- Dropping columns with more than 50% of blank cells:

| |
|---|
| RATE_INTEREST_PRIMARY |
| RATE_INTEREST_PRIVILEGED |
| AMT_DOWN_PAYMENT |
| RATE_DOWN_PAYMENT |

- Dropping these 4 columns, we are left with 33 columns, out of which we can impute the remaining columns:

| |
|---|
| NAME_TYPE_SUITE |
| DAYS_FIRST_DRAWING |
| DAYS_FIRST_DUE |
| DAYS_LAST_DUE_1ST_VERSION |
| DAYS_LAST_DUE |
| DAYS_TERMINATION |
| |
| NFLAG_INSURED_ON_APPROVAL |
| AMT_GOODS_PRICE |
| AMT_ANNUITY |
| CNT_PAYMENT |
| PRODUCT_COMBINATION |

# Data cleaning

▶ Dropping unnecessary columns which don't provide any information:

1. NAME_TYPE_SUITE
2. WEEKDAY_APPR_PROCESS_START
3. HOUR_APPR_PROCESS_START
4. FLAG_LAST_APPL_PER_CONTRACT
5. NFLAG_LAST_APPL_IN_DAY

# Median imputation

1. AMT_ANNUITY
2. AMT_GOODS_PRICE
3. DAYS_FIRST_DRAWING
4. DAYS_FIRST_DUE
5. DAYS_LAST_DUE_1ST_VERSION
6. DAYS_LAST_DUE DAYS_TERMINATION

# MODE IMPUTATION

▶ NFLAG_INSURED_ON_APPROVAL

# Data cleaning: Dropping

▶ Dropping the blank cells from "PRODUCT_COMBINATION" feature as the blanks were just 8 in number and dropping would be better than imputation.

# Data cleaning: Custom imputation

- On analyzing the contract status for the CNT_PAYMENT, most of the contracts were approved, hence imputing CNT_PAYMENT with median would be better.

Contract Status

# Outliers

### AMT_ANNUITY

### AMT_APPLICATION

### DAYS_FIRST_DUE

### CNT_PAYMENT

### DAYS_DECISION

- There are many outliers in the columns like AMT_ANNUITY, AMT_APPLICATION etc.
- Few outliers in CNT_PAYMENT.
- Few outliers in DAYS_DECISION are such that it indicates that the decision time taken is high, which is not a good practice.

# Removing outliers

- I have cleaned the outliers from the columns such as DAYS_FIRST_DRAWING, DAYS_FIRST_DUE, DAYS_LAST_DUE_1ST_VERSION, DAYS_LAST_DUE, DAYS_TERMINATION as the outliers were unrealistic values.



DAYS_FIRST_DRAWING



DAYS_FIRST_DUE

# Data imbalance

As we can see there is data imbalance and the number of defaulters is way less than the number of re-payers in this dataset provided.



Count of target

4026, 8%

45973, 92%

0
1

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Age

Most of the loan applicants are from age group 30-40 years. As the age increases the rate of defaulting decreases.

## Univariate analysis

### AGE



## Univariate segmented analysis

### Age vs Defaulters

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Years Employed

Most of the loan applicants are from 0-5 years of experience. As the experience increases the rate of defaulting decreases

## Univariate analysis



Years Employed

## Univariate segmented analysis



Years Employed vs Defaulters

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Contract type

## Univariate analysis

▶ Highest are cash loans.

**Count of Contract Type**



## Univariate segmented analysis

▶ High defaulters are from Cash loans

▶ Cash loans: 9.2%

▶ Revolving loans: 5.1%

**Contract Type vs Defaulters**

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Gender

## Univariate analysis

▶ Most of the females have taken loan.

**Count of Gender**



## Univariate segmented analysis

▶ Female defaulters: 7.5%

▶ Male defaulters: ~11%

**Gender vs Defaulters**

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: NAME_TYPE_SUITE

## Univariate analysis

▶ Most of the people were unaccompanied while taking loan.



NAME_TYPE_SUITE

## Univariate segmented analysis

▶ Here I have shown the categories with significant numbers.



Company vs Defaulters

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Education

## Univariate analysis

▶ Most of the customers have secondary and higher education.

### Education



## Univariate segmented analysis

▶ Highest and lowest defaulters have lower secondary (16%) and academic degree(0%) respectively.

### Total

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Family Status

## Univariate analysis

▶ Highest number of people who took loans are married.

### Family Status



## Univariate segmented analysis

▶ Highest default percent is in civil marriage and single people ~10%

### Total

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Housing Type

## Univariate analysis

▶ People who took loans had houses/apartments.

### Housing type



## Univariate segmented analysis

▶ Highest defaulters are from rented apartments and with parents ~11%

### Total

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Occupation

## Univariate analysis

▶ Laborers are the highest loan taking occupation.

## Univariate segmented analysis

▶ Highest defaulter percent is for low skilled labourers ~17%.



Occupation



Total

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Children

## Univariate analysis

► Most of the people taking loans are either having no child or one/two children.

**Children**



## Univariate segmented analysis

► The highest defaulters are people having more number of children 25% of defaulters.

**Total**

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: Region_Rating_Client

## Univariate analysis

▶ Most of the people taking loans are from region 2.

### Region rating client



## Univariate segmented analysis

▶ Region rating 3 have the highest defaulters.

### Total

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS:No. of family members

## Univariate analysis

► Most of the people taking loans have 2 family members.

No. of family members



## Univariate segmented analysis

► Highest percent of defaulters are from people having more family members.

Total

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS

▶ The most of the loans are taken between 245000 to 645000.

▶ Highest defaulters are from 445000 – 645000.



AMT_CREDIT

# UNIVARIATE ANALYSIS:

Occupation type



Business entity Type 3 takes the highest amount of loans.

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS

▶ Most of the defaulters are from the income range of 25650 – 1025650.



INCOME

# UNIVARIATE/SEGMENTED UNIVARIATE ANALYSIS: previous_application

## Univariate analysis

▶ Most of the previous loan takers are for XNA category.

### NAME_CASH_LOAN_PURPOSE



## Univariate segmented analysis

▶ Most of the loan defaulters are from XNA category.

### Total

# UNIVARIATE ANALYSIS: previous_application

Previous application Name Contract Status



- Approved
- Canceled
- Refused
- Unused offer

Approximately equal number of previous loans were either approved or refused.

# Bivariate/Segmented Bivariate Analysis

▶ We can see the credit amount is highest for people in age bracket 51-61 years so the insight is that the 51-61 year people default less than the others.

### Age vs Avg Credit Amount

# Bivariate/Segmented Bivariate Analysis

▶ We see pensioners having the least amount credited but also that they don't default.

Income type vs amount credit based on defaulters

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 582178.8996 | 553675.7117 | 202500 | 589739.8718 | 595881.9866 | 539246.7 | 537158.948 | 508529.2396 |
| 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 |
| Commercial associate | | Pensioner | State servant | | Student | Working | |

# Correlation among the features: Defaulters



▶ Top 5 correlations are from:

1. OBS_30_CNT_SOCIAL_CIRCLE – OBS_60_CNT_SOCIAL_CIRCLE

2. AMT_CREDIT – AMT_GOODS_PRICE

3. REGION_RATING_CLIENT – REGION_RATING_CLIENT_W_CITY

4. CNT_CHILDREN – CNT_FAM_MEMBERS

5. REG_REGION_NOT_WORK_REGION – LIVE_REGION_NOT_WORK_REGION

# Correlation among the features: Repayers



▶ Top 5 correlations are from:

1. OBS_30_CNT_SOCIAL_CIRCLE - OBS_60_CNT_SOCIAL_CIRCLE

2. AMT_CREDIT - AMT_GOODS_PRICE

3. REGION_RATING_CLIENT - REGION_RATING_CLIENT_W_CITY

4. CNT_CHILDREN - CNT_FAM_MEMBERS

5. REG_REGION_NOT_WORK_REGION - LIVE_REGION_NOT_WORK_REGION

# Insights and hypothesis:

- Elderly people with good experience years are less likely to default on a loan, hence can be considered good customers for loans.

- Most people take cash loans and hence there is a high defaulter rate for cash loans.

- Females take more loans than males but the default rate for males is higher, hence the banks can consider giving out loans to females more than men.

- People with less educational qualifications tend to default more than the ones having academic degrees.

- Most of the single or civil marriage people default more than any other family status.

- Taking a high amount as the loan has a high default rate, this could be considered while giving out loans.

- People with region rating 3 default more than any other region rating.

- People with less income default more.

- Pensioners have less amount of credit but they don't default much, hence can be considered for providing loans in the future.

# Results:

- This project was of a great help to work with huge dataset with data cleaning and handling outliers using the domain knowledge.

- This project helped me understand how to merge two excel sheets and use add-ins like data analyse for correlation.

# Drive Links:

- Application_data: https://docs.google.com/spreadsheets/d/1f4LFl0Z9NhbGHCyuxTpe6rzyQKI_BKbr/edit?usp=sharing&ouid=115109770037321084146&rtpof=true&sd=true

- Previous_application: https://docs.google.com/spreadsheets/d/1PdOBodjo_yOuVJycu9P0xddODKUKUSB8/edit?usp=sharing&ouid=115109770037321084146&rtpof=true&sd=true

- Merged sheets: https://docs.google.com/spreadsheets/d/1aAJemTTQC8RxI_jPQmlz1PSiB3x9Jom-/edit?usp=sharing&ouid=115109770037321084146&rtpof=true&sd=true

- Correlation: https://docs.google.com/spreadsheets/d/1SboxQXWvU-F6tPMJZDoCIlje-Hpk6-Ck/edit?usp=sharing&ouid=115109770037321084146&rtpof=true&sd=true

- Analysis: https://docs.google.com/spreadsheets/d/1al0q0--yOniq_62NPtfyDc9WAu2xjJsR/edit?usp=sharing&ouid=115109770037321084146&rtpof=true&sd=true

# ThankYou