

# Exploratory Data Analysis

---

Assignment 1

By

Kiran Yashawant Awari

# Data reading - 'application\_data.csv'

---

- Imported warnings
- Imported important libraries – Numpy, Pandas, Matplotlib, Seaborn etc.
- Read CSV file 'application\_data.csv' for first analysis as df
- Understood shape of file by using df.shape



# Data Cleansing

---

- Identified missing values in columns using function `isnull()`
- Removed columns with more than 40 % missing values
- For column 'OCCUPATION\_TYPE' replaced all missing values with mode
- Replaced missing values in column 'EXT\_SOURCE\_3' with median value
- Replaced missing values in right skewed 6 columns with respective median values
- Replaced missing values in columns 'AMT\_ANNUITY', 'AMT\_GOODS\_PRICE', 'EXT\_SOURCE\_2' with median value
- Replaced missing values in column 'NAME\_TYPE\_SUITE' with mode

# Univariate Analysis on Categorical Columns

---

1. application for cash loans is more than that of revolving loan
2. number of female applicants are more than that of male applicants
3. most of applicants do not own a car
4. maximum number of applicants own a real estate property
5. number of Unaccompanied applicants is more
6. maximum number of working professionals applied for loan
7. Most of applicants have secondary education
8. marital status of maximum number of applicants is married
9. Housing type of maximum applicants is House/ Apartment
10. Most of applicants are laborers
11. Applications are carried on all week days where on weekends applications are less as compared to other weekdays

# Correlation Analysis at Target = 0

---

There is a strong correlation between the following TOP 10 numeric variables for TARGET=0:

- DAYS\_EMPLOYED & FLAG\_EMP\_PHONE
- OBS\_60\_CNT\_SOCIAL\_CIRCLE & OBS\_30\_CNT\_SOCIAL\_CIRCLE
- FLOORSMAX\_AVG & FLOORSMAX\_MEDI
- YEARS\_BEGINEXPLUATATION\_AVG & YEARS\_BEGINEXPLUATATION\_MEDI
- FLOORSMAX\_MEDI & FLOORSMAX\_MODE
- AMT\_GOODS\_PRICE & AMT\_CREDIT
- FLOORSMAX\_MODE & FLOORSMAX\_AVG
- YEARS\_BEGINEXPLUATATION\_MODE & YEARS\_BEGINEXPLUATATION\_AVG
- YEARS\_BEGINEXPLUATATION\_MEDI & YEARS\_BEGINEXPLUATATION\_MODE
- REGION\_RATING\_CLIENT\_W\_CITY & REGION\_RATING\_CLIENT



# Correlation Analysis at Target = 1

---

There is a strong correlation between the following TOP 10 numeric variables for TARGET=1

1 FLAG\_EMP\_PHONE & EMP\_YRS

2 DAYS\_EMPLOYED & FLAG\_EMP\_PHONE

3 OBS\_30\_CNT\_SOCIAL\_CIRCLE & OBS\_60\_CNT\_SOCIAL\_CIRCLE

4 FLOORSMAX\_AVG & FLOORSMAX\_MEDI

5 YEARS\_BEGINEXPLUATATION\_AVG & YEARS\_BEGINEXPLUATATION\_MEDI

6 FLOORSMAX\_MEDI & FLOORSMAX\_MODE

7 FLOORSMAX\_AVG & FLOORSMAX\_MODE

8 AMT\_GOODS\_PRICE & AMT\_CREDIT

9 YEARS\_BEGINEXPLUATATION\_AVG & YEARS\_BEGINEXPLUATATION\_MODE

10 YEARS\_BEGINEXPLUATATION\_MEDI & YEARS\_BEGINEXPLUATATION\_MODE

# Reading CSV file 'previous\_application.csv'

---

- Read csv file 'previous\_application.csv' as inp0
- Understood data by using functions describe() and info()

# Data cleansing of 'previous\_application.csv'

---

- Identified missing values by using `isnull()` function
- Removed columns having missing values more than 40 %
- For columns `AMT_ANNUITY`, `AMT_GOODS_PRICE` & `CNT_PAYMENT` replaced missing values with median values
- Replaced missing values in column `AMT_CREDIT` with median value
- Replaced missing values of `PRODUCT_COMBINATION` with Cash as it is mode of the column



# Univariate analysis on previous\_application.csv

---

- Most number of loans were Cash Loans
- Maximum number of applications are on Monday, Tuesday and Wednesday
- Most of applications had one application per loan contract
- Most had "XAP" purpose of loan application (Analysis cannot be performed)
- Maximum number of applications were approved
- Most of loans were repaid by "Cash through the Bank"
- Maximum number of loans were rejected because of reason "XAP" (Analysis cannot be performed)
- Most previous loans were applied by people who were unaccompanied
- Most applicants were applied by Repeater clients
- Maximum number of XNA goods observed (Analysis cannot be performed)
- Most of them were for POS
- Most were product type "XNA" (not applicable) (Analysis cannot be performed )
- Most applications had clients acquired from Credit and cash offices
- Most were from "XNA" (not applicable) seller industry (Analysis cannot be performed)
- Most of applications fell under "XNA" (not applicable) interest group (Analysis cannot be performed )
- Most applications had the Product Combination as "Cash"

# Merging 2 dataframes

---

- Merged `application_data.csv` and `previous_application.csv` by `pd.merged` function