Cairo University
Faculty of Engineering
Computer Engineering Department

# Bank Loan Default Risk Analysis
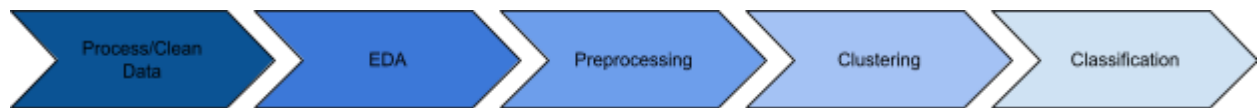
## Project Report

# Problem Description

The project aims to analyze risk factors in loan applications to help banks minimize losses that happen due to defaults. As there are some applications with no credit history or missing guarantees this leads to high risks to banks and lending companies. This results in increased default rates.

By analysing previous and current loan applications, this project will aim to provide insights to how banks and lending companies can assess risk and improve lending strategies.

# Pipeline



## Process/Clean Data

In this phase, the raw dataset was cleaned and preprocessed to ensure data quality and consistency.

1. **Initial Exploration:**
   Exploring the shape of the dataset.
   Basic information about the dataset was displayed, including column names, data types. The number of unique values for each column was calculated to understand the diversity of categorical and numerical features.
2. **Data Cleaning:**
   Null Value Handling: Columns with more than 40% null values were identified and removed to ensure data quality.
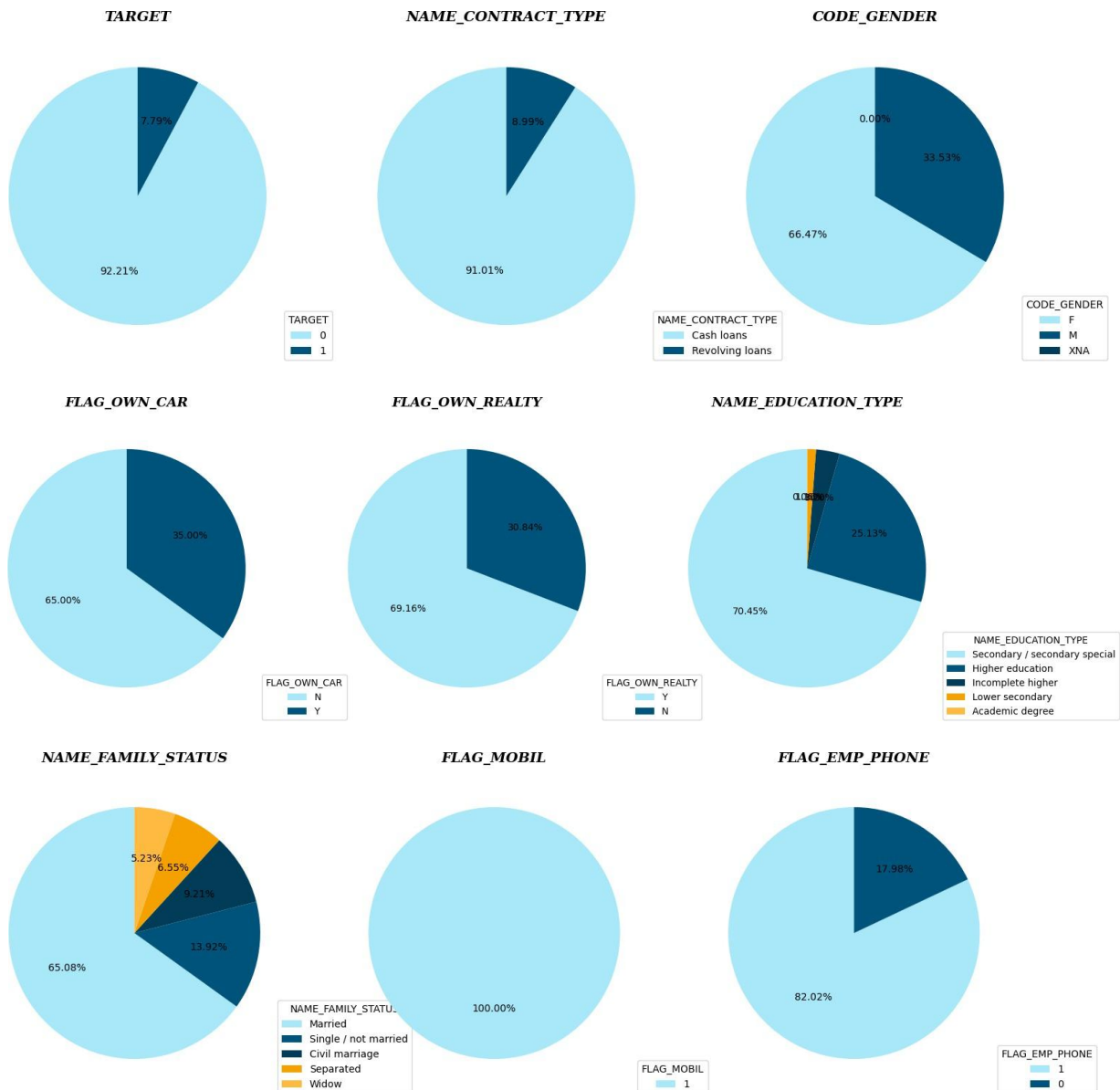   Duplicates Removal: check for duplicate rows and drop it.
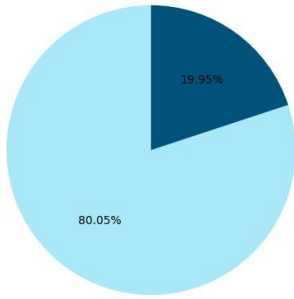   Row Removal: Rows containing any null values were dropped to maintain consistency in the dataset.

# EDA
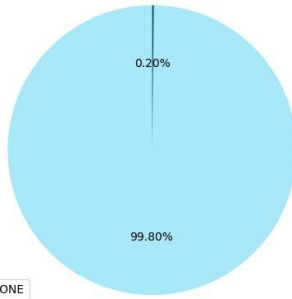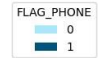
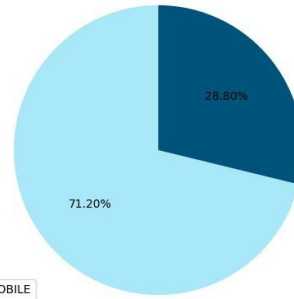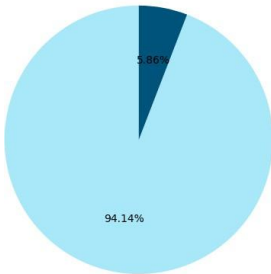## Univariate Analysis (Applications)

### Categorical Variables



TARGET

NAME_CONTRACT_TYPE

CODE_GENDER

FLAG_OWN_CAR

FLAG_OWN_REALTY

NAME_EDUCATION_TYPE

NAME_FAMILY_STATUS

FLAG_MOBIL

FLAG_EMP_PHONE

## FLAG_WORK_PHONE

19.95%

80.05%

FLAG_WORK_PHONE
- 0
- 1

## FLAG_CONT_MOBILE

0.20%

99.80%

FLAG_CONT_MOBILE
- 1
- 0

## FLAG_PHONE

28.80%

71.20%

FLAG_PHONE
- 0
- 1

## FLAG_EMAIL

5.86%

94.14%

FLAG_EMAIL
- 0
- 1

## REGION_RATING_CLIENT

10.03%

15.50%

74.48%

REGION_RATING_CLIENT
- 2
- 3
- 1

## REGION_RATING_CLIENT_W_CITY

10.66%

14.14%

75.20%

REGION_RATING_CLIENT_W_CITY
- 2
- 3
- 1

## REG_REGION_NOT_LIVE_REGION

1.42%

98.58%

REG_REGION_NOT_LIVE_REGION
- 0
- 1

## REG_REGION_NOT_WORK_REGION

4.88%

95.12%

REG_REGION_NOT_WORK_REGION
- 0
- 1

## LIVE_REGION_NOT_WORK_REGION

3.94%

96.06%

LIVE_REGION_NOT_WORK_REGION
- 0
- 1

## REG_CITY_NOT_LIVE_CITY

7.47%

92.53%

REG_CITY_NOT_LIVE_CITY
- 0
- 1

## REG_CITY_NOT_WORK_CITY

22.53%

77.47%

REG_CITY_NOT_WORK_CITY
- 0
- 1

## LIVE_CITY_NOT_WORK_CITY

17.71%

82.29%

LIVE_CITY_NOT_WORK_CITY
- 0
- 1

## FLAG_DOCUMENT_2

100.00%

FLAG_DOCUMENT_2
0

## FLAG_DOCUMENT_3

28.35%

71.65%

FLAG_DOCUMENT_3
1
0

## FLAG_DOCUMENT_4

0.01%

99.99%

FLAG_DOCUMENT_4
0
1

## FLAG_DOCUMENT_5

1.47%

98.53%

FLAG_DOCUMENT_5
0
1

## FLAG_DOCUMENT_6

8.79%

91.21%

FLAG_DOCUMENT_6
0
1

## FLAG_DOCUMENT_7

0.01%

99.99%

FLAG_DOCUMENT_7
0
1

## FLAG_DOCUMENT_8

8.16%

91.84%

FLAG_DOCUMENT_8
0
1

## FLAG_DOCUMENT_9

0.36%

99.64%

FLAG_DOCUMENT_9
0
1

## FLAG_DOCUMENT_10

0.00%

100.00%

FLAG_DOCUMENT_10
0
1

## FLAG_DOCUMENT_11

0.33%

99.67%

FLAG_DOCUMENT_11
0
1

## FLAG_DOCUMENT_12

0.00%

100.00%

FLAG_DOCUMENT_12
0
1

## FLAG_DOCUMENT_13

0.38%

99.62%

FLAG_DOCUMENT_13
0
1

### FLAG_DOCUMENT_14

0.30%

99.70%

FLAG_DOCUMENT_14
- 0
- 1

### FLAG_DOCUMENT_15

0.12%

99.88%

FLAG_DOCUMENT_15
- 0
- 1

### FLAG_DOCUMENT_16

1.00%

99.00%

FLAG_DOCUMENT_16
- 0
- 1

### FLAG_DOCUMENT_17

0.03%

99.97%

FLAG_DOCUMENT_17
- 0
- 1

### FLAG_DOCUMENT_18

0.82%

99.18%

FLAG_DOCUMENT_18
- 0
- 1

### FLAG_DOCUMENT_19

0.06%

99.94%

FLAG_DOCUMENT_19
- 0
- 1

### FLAG_DOCUMENT_20

0.06%

99.94%

FLAG_DOCUMENT_20
- 0
- 1

### FLAG_DOCUMENT_21

0.03%

99.97%

FLAG_DOCUMENT_21
- 0
- 1

### AMT_REQ_CREDIT_BUREAU_HOUR

0.08%

99.39%

AMT_REQ_CREDIT_BUREAU_HOUR
- 0.0
- 1.0
- 2.0
- 3.0
- 4.0

### Insights:

- 7.79% of the clients face defaults.
- There are more females in the dataset. (~66.5%)
- Most of the clients (are married/don't own cars/own realty). (~65% - 70%)
- Most of the clients live in regions of middle rating.
- Most of the clients are with education ("Secondary / secondary special", "Higher education") level.
- Most clients don't deliver less than 2 documents.
- Most of the loans are cash (~91%).
- Most of the clients live in the city they work in.

**Documents_count**

| Documents_count |
|---|
| 1 |
| 0 |
| 2 |
| 3 |
| 4 |

88.28% | 9.11% | 0.05%

**OWN_CAR_REALTY**

| OWN_CAR_REALTY |
|---|
| NY |
| YY |
| NN |
| YN |

45.05% | 24.11% | 19.95% | 10.89%

**Documents_count**

88.3% | 9.1% | 2.6% | 0.1%

## Insights:

- From all the documents 88% of clients deliver only 1.
- Nearly 45% of clients have a realty but don't have a car.

**AMT_REQ_CREDIT_BUREAU_WEEK** and **AMT_REQ_CREDIT_BUREAU_QRT**

**Insights:**

- Most clients don't have children. (~70%)
- Most clients apply for the loan alone then with family in 2nd proportion.
- Nearly half of the clients are working in a standard job with fixed income.
- Most of clients have privately owned apartment or flat (~90%)
- ~90% of clients don't have any defaults in their social network in the last 60 days
-

## Numerical Variables

**Insights**:
- A large portion of clients modified their phone records fewer than 125 days before applying.

As we saw that some histograms have outliers that made the histogram unreadable, we plotted again with removing outliers

# Multivariate Analysis (Applications)



Correlation Matrix

# Univariate Analysis (Previous Applications)

## Categorical Variables

### NAME_CONTRACT_TYPE



- 55.48% Consumer loans
- 37.03% Cash loans
- 7.49% Revolving loans

NAME_CONTRACT_TYPE
- Consumer loans
- Cash loans
- Revolving loans

### FLAG_LAST_APPL_PER_CONTRACT



- 99.72% Y
- 0.28% N

FLAG_LAST_APPL_PER_CONTRACT
- Y
- N

### NFLAG_LAST_APPL_IN_DAY



- 99.79% 1
- 0.21% 0

NFLAG_LAST_APPL_IN_DAY
- 1
- 0

### NAME_CONTRACT_STATUS



- 79.75% Approved
- 19.33% Refused

NAME_CONTRACT_STATUS
- Approved
- Refused
- Canceled
- Unused offer

### NAME_PAYMENT_TYPE



- 79.99% Cash through the bank
- 19.27% XNA

NAME_PAYMENT_TYPE
- Cash through the bank
- XNA
- Non-cash from your account
- Cashless from the account of the employer

### NAME_CLIENT_TYPE



- 67.71% Repeater
- 23.63% New
- 8.59% Refreshed
- 0.07% XNA

NAME_CLIENT_TYPE
- Repeater
- New
- Refreshed
- XNA

### NAME_PORTFOLIO



- 55.44% POS
- 37.03% Cash
- 7.49% Cards

NAME_PORTFOLIO
- POS
- Cash
- Cards
- Cars

### NAME_PRODUCT_TYPE



- 55.48% XNA
- 32.47% x-sell
- 12.06% walk-in

NAME_PRODUCT_TYPE
- XNA
- x-sell
- walk-in

### NAME_YIELD_GROUP



- 30.93% middle
- 28.35% high
- 25.84% low_normal
- 7.49% XNA
- 7.39% low_action

NAME_YIELD_GROUP
- middle
- high
- low_normal
- XNA
- low_action

## Numerical Variables



**Insights:**
- A large portion of clients schedule loan repayment under 20 terms (ex. months).

As we saw that some histograms have outliers that made the histogram unreadable, we plotted again with removing outliers

# Bivariate Analysis (Applications)

## NAME_INCOME_TYPE



Working — TARGET 0: 90.78%, TARGET 1: 9.22%

State servant — TARGET 0: 94.42%, TARGET 1: 5.58%

Commercial associate — TARGET 0: 92.64%, TARGET 1: 7.36%

Pensioner — TARGET 0: 94.83%, TARGET 1: 5.17%

Unemployed — TARGET 0: 100.00%

Student — TARGET 0: 100.00%

Businessman — TARGET 0: 100.00%



Target Distribution by AMT_INCOME_TOTAL Bins

NAME_INCOME_TYPE



Target Distribution by CREDIT_TO_INCOME Bins

**Insights:**
- From the working and commercial associate clients there is (7% to 10%) who face defaults.
- As total income increase, default rates decreases
- Clients that ask for credit 2-5 x their income have a higher default rate.

Target Distribution by AMT_CREDIT Bins

**Insights:**
- The risk of default seems to peak in the middle ranges of loan amounts. Specifically, the highest default ratio appears to be around the 314,100-450,000 range, where the red line reaches its maximum at approximately 0.10 (10% default rate).



Target Distribution by DAYS_EMPLOYED Bins

**Insights:**
- Default risk increases as employment duration decreases. The highest default rates (~11%) appear in the 676-482, 482-307, and 307-152 bins, representing people with relatively short employment durations. (for the last bin, it likely contains anomalous values)

# ORGANIZATION_TYPE



| | | | | | |
|---|---|---|---|---|---|
| **Business Entity Type 3** | **Government** | **Other** | **Medicine** | **Business Entity Type 2** | **Self-employed** |
| 9.10% / 90.90% | 6.66% / 93.34% | 7.44% / 92.56% | 6.49% / 93.51% | 8.33% / 91.67% | 9.85% / 90.15% |
| **Housing** | **Kindergarten** | **Trade: type 7** | **Industry: type 11** | **Military** | **Transport: type 4** |
| 7.60% / 92.40% | 6.96% / 93.04% | 9.06% / 90.94% | 8.38% / 91.62% | 4.97% / 95.03% | 8.95% / 91.05% |
| **School** | **Services** | **Emergency** | **Security** | **Trade: type 2** | **University** |
| 5.49% / 94.51% | 6.71% / 93.29% | 7.11% / 92.89% | 9.69% / 90.31% | 6.93% / 93.07% | 4.92% / 95.08% |
| **Police** | **Construction** | **Business Entity Type 1** | **Industry: type 4** | **Agriculture** | **Restaurant** |
| 4.61% / 95.39% | 11.37% / 88.63% | 8.20% / 91.80% | 9.87% / 90.13% | 9.46% / 90.54% | 11.90% / 88.10% |
| **Transport: type 2** | **Hotel** | **Industry: type 7** | **Trade: type 3** | **Industry: type 3** | **Bank** |
| 7.48% / 92.52% | 6.38% / 93.64% | 8.13% / 91.87% | 9.79% / 90.21% | 10.41% / 89.59% | 4.85% / 95.15% |

Pie charts showing TARGET distribution (0/1) by sector:

| Industry: type 9 — 93.53% | Postal — 91.51% | Trade: type 6 — 96.09% | Industry: type 2 — 93.25% | Transport: type 1 — 95.29% | Transport: type 3 — 84.30% |
|---|---|---|---|---|---|
| Electricity — 93.87% | Industry: type 12 — 96.30% | Insurance — 94.17% | Industry: type 1 — 88.57% | Security Ministries — 95.27% | Mobile — 90.27% |
| Trade: type 1 — 92.67% | Industry: type 5 — 92.98% | Industry: type 10 — 92.39% | Legal Services — 92.95% | Advertising — 91.09% | Trade: type 5 — 94.74% |
| Cleaning — 91.04% | Industry: type 13 — 90.70% | Industry: type 8 — 93.75% | Realtor — 89.97% | Culture — 95.35% | Telecom — 92.65% |
| Religion — 92.65% | Industry: type 6 — 93.02% | Trade: type 4 — 96.15% | | | |

**Insights:**

- Borrowers employed in the Mobile (16.10%) and Hotel (13.65%) sectors show the highest default rates, while Government, Military, and Banking sectors maintain lower, more stable default rates (~9%).

# CODE_GENDER

### *M*



9.73%

90.27%

TARGET
- 0
- 1

### *F*



6.81%

93.19%

TARGET
- 0
- 1

## NAME_FAMILY_STATUS

### *Single / not married*



9.27%

90.73%

TARGET
- 0
- 1

### *Married*



7.40%

92.60%

TARGET
- 0
- 1

### *Widow*



5.53%

94.47%

TARGET
- 0
- 1

### *Civil marriage*



9.58%

90.42%

TARGET
- 0
- 1

### *Separated*



7.81%

92.19%

TARGET
- 0
- 1

**NAME_EDUCATION_TYPE**

| Secondary / secondary special | Higher education | Incomplete higher | Lower secondary | Academic degree |
|---|---|---|---|---|
| 8.60% / 91.40% | 5.34% / 94.66% | 8.26% / 91.74% | 10.48% / 89.52% | 2.11% / 97.79% |

Target Distribution by AGE Bins



**Insights:**

- Borrowers who are men (9.73%) show higher default rates, compared to women (~6.81%).
- Borrowers who are widowed show the lowest default rate (5.53%). Meanwhile, those in a civil marriage (9.58%) and single/not married (9.27%) categories have the highest default rates.
- There is an inverse relationship between education level and default risk. The lower the education level, the higher the chance of default.
- Younger clients have higher default rates.

Target Distribution by OBS_30_CNT_SOCIAL_CIRCLE Bins



Target Distribution by OBS_60_CNT_SOCIAL_CIRCLE Bins

**<u>Insights:</u>**

- Most clients (~85%) have 0-1 default in their social circle.
- As the number of acquaintances with defaults increases, the client's default rate increases slightly. The effect is small, but positive.

**NAME_HOUSING_TYPE**

| House / apartment | Rented apartment | Municipal apartment | With parents | Office apartment | Co-op apartment |
|---|---|---|---|---|---|
| 92.47% | 87.91% | 91.85% | 88.51% | 93.84% | 93.00% |

TARGET: 0, 1 (for each chart)

## Insights:

- Ownership (house, co-op) is a positive indicator for financial stability.
- Renting or living with parents indicates a higher risk, possibly due to less financial independence or lower income levels.

# FLAG_OWN_CAR

**N**



8.19%

91.81%

TARGET
0
1

**Y**



7.06%

92.94%

TARGET
0
1

# FLAG_OWN_REALTY

**Y**



7.67%

92.33%

TARGET
0
1

**N**



8.06%

91.94%

TARGET
0
1

**OWN_CAR_REALTY**

**NY**



7.95%

92.05%

TARGET
0
1

**YY**



7.15%

92.85%

TARGET
0
1

**YN**



6.86%

93.14%

TARGET
0
1

**NN**



8.72%

91.28%

TARGET
0
1

## Insights:
- Clients who doesn't have a car, or realty have higher defaults.



Target Distribution by CNT_FAM_MEMBERS Bins

## Insights:
- As the number of family members increases, the client's default rate increases slightly. The effect is small, but positive.



Target Distribution by Documents_count Bins

- Most clients submit only one or no documents. Submitting more documents is rare but slightly reduces the risk of default.

## REG_REGION_NOT_LIVE_REGION

*0*

7.77%

92.23%

TARGET
- 0
- 1

*1*

9.12%

90.88%

TARGET
- 0
- 1

## REG_REGION_NOT_WORK_REGION

*0*

7.75%

92.25%

TARGET
- 0
- 1

*1*

8.69%

91.31%

TARGET
- 0
- 1

# LIVE_REGION_NOT_WORK_REGION

### 0



7.77%

92.23%

TARGET
0
1

### 1



8.21%

91.79%

TARGET
0
1

# REG_CITY_NOT_LIVE_CITY

### 0



7.46%

92.54%

TARGET
0
1

### 1



11.93%

88.07%

TARGET
0
1

# REG_CITY_NOT_WORK_CITY

**0**



7.07%

92.93%

TARGET
- 0
- 1

**1**

10.27%

89.73%

TARGET
- 0
- 1

# LIVE_CITY_NOT_WORK_CITY

**0**

7.40%

92.60%

TARGET
- 0
- 1

**1**

9.62%

90.38%

TARGET
- 0
- 1

**Insights:**
- Clients' whose work region is not the same as live region have more risk of default.

## REGION_RATING_CLIENT



## REGION_RATING_CLIENT_W_CITY





Target Distribution by REGION_POPULATION_RELATIVE Bins

**Insights:**

- Clients from regions with higher ratings (Rating 3) have more the default rate compared to those with lower rating (Rating 1).

- Clients who live in regions with higher populations have fewer default risk.

## CURR_AMT_GT_PREV_MAX

### 0

7.77%

92.23%

TARGET
- 0
- 1

### 1

7.80%

92.20%

TARGET
- 0
- 1

Target Distribution by PREV_COUNT Bins



Target Distribution by PREV_ACCEPTED_RATIO Bins

Target Distribution by CURR_AMT_MEAN_RATIO Bins

**Insights:**
- Asking for a new loan larger than any one of the previous doesn't mean higher default risk.
- As the client has more previous loans, the default risk increases.
- As the count of previously accepted loans over the total ratio increases, the default risk decreases.
- Borrowers whose current loan is smaller than their historical mean (ratio < 0.63) have a lower default rate (~7%).
- Above this threshold, default risk seems to stabilize (around 7.5-8.5%).

# Preprocessing

1- Process Curr Data

**Application_data.csv** —> **processed_current_application.csv**
**col with Threshold for null & row contains nulls**

2- Process Prev Data

**previous_application.csv** → **processed_previous_application.csv**
**1- col with Threshold for null & row contains nulls**
**2- Generate New Features using (prev_features.sh) file**

```
NAME_CONTRACT_TYPE  -> 3 features (all features)
AMT_ANNUITY -> avg(AMT_ANNUITY)
AMT_APPLICATION -> avg(AMT_APPLICATION)
AMT_CREDIT -> avg(AMT_CREDIT)
AMT_GOODS_PRICE -> avg(AMT_GOODS_PRICE/AMT_APPLICATION)
NAME_CASH_LOAN_PURPOSE -> XAP, other
NAME_CONTRACT_STATUS -> 2 features (approved , refused)
DAYS_DECISION -> avg(DAYS_DECISION)
FLAG_LAST_APPL_PER_CONTRACT -> sum(0)
NFLAG_LAST_APPL_IN_DAY-> sum(0)
NAME_PAYMENT_TYPE -> 1 features (cash payment) -> XX
CODE_REJECT_REASON -> 3 features (XAP, HC , Limit)
NAME_CLIENT_TYPE -> 2 features (repeater , refreshed)
NAME_PORTFOLIO -> 3 features (POS , Cash , Cards)
CNT_PAYMENT -> mean(CNT_PAYMENT)
NAME_YIELD_GROUP -> avg(encoded (NAME_YIELD_GROUP))
SK_ID_CURR -> count(SK_ID_CURR)
HOUR_APPR_PROCESS_START -> mean(HOUR_APPR_PROCESS_START)
CHANNEL_TYPE -> top 3 (Credit and cash offices, Country-wide, Stone)
PRODUCT_COMBINATION -> top 3 (Cash, POS household with interest, POS mobile with interest)
```

3- Merge Data

**Processed_current_application.csv & processed_current_application.csv –> merged_application.csv**

**Merge on SK_ID_CURR**

4- Encoding_Outliers_FeatureSelection

**merged_application.csv  –> encoded_merged_application.csv**

**Encode the Features using (curr_application_features_encoding_methods.txt) file**

**Encoded_merged_application.csv → df_no_outliers.csv**

**Remove outliers using Zscore**

**Df_no_outliers.csv → high corr features (two files before and after )**

**→ featureSelected_encoded_merged_application.csv**

**→ X_train.csv , y_train.csv & X_test.csv , y_test.csv**
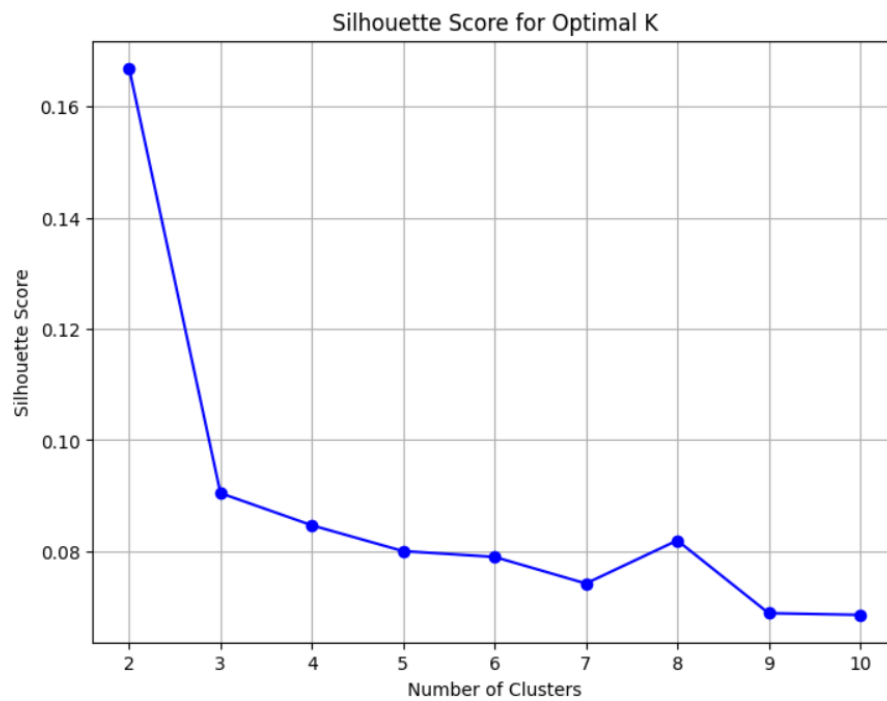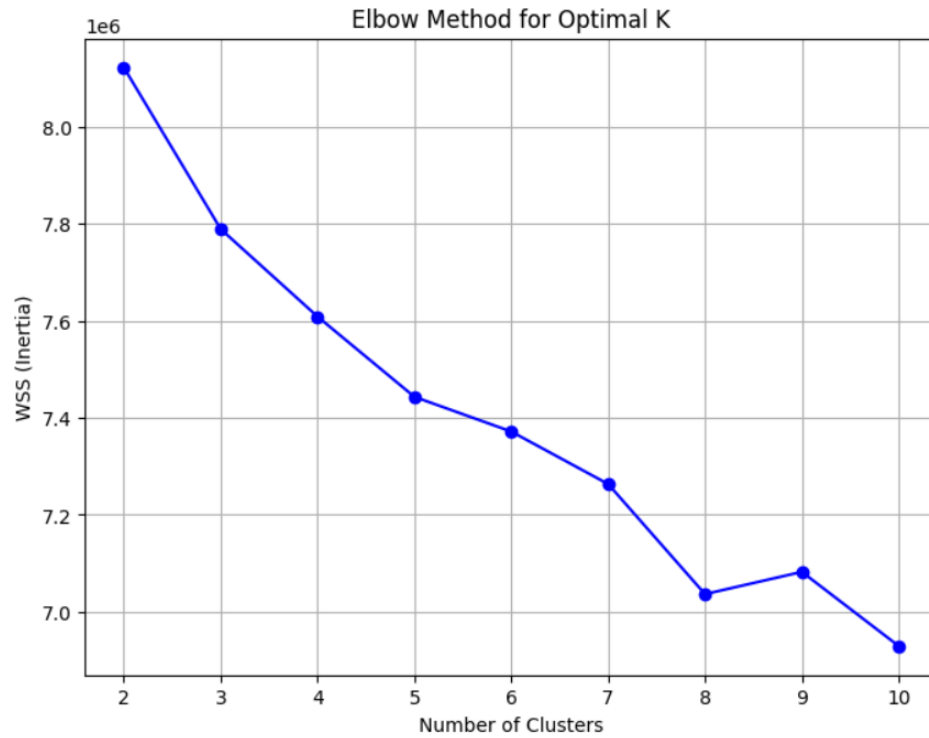
**1- Split data (train_test)**

**2- Study High Correlated Features (MulCorr)**

**3 -Low Correlated Features with target**

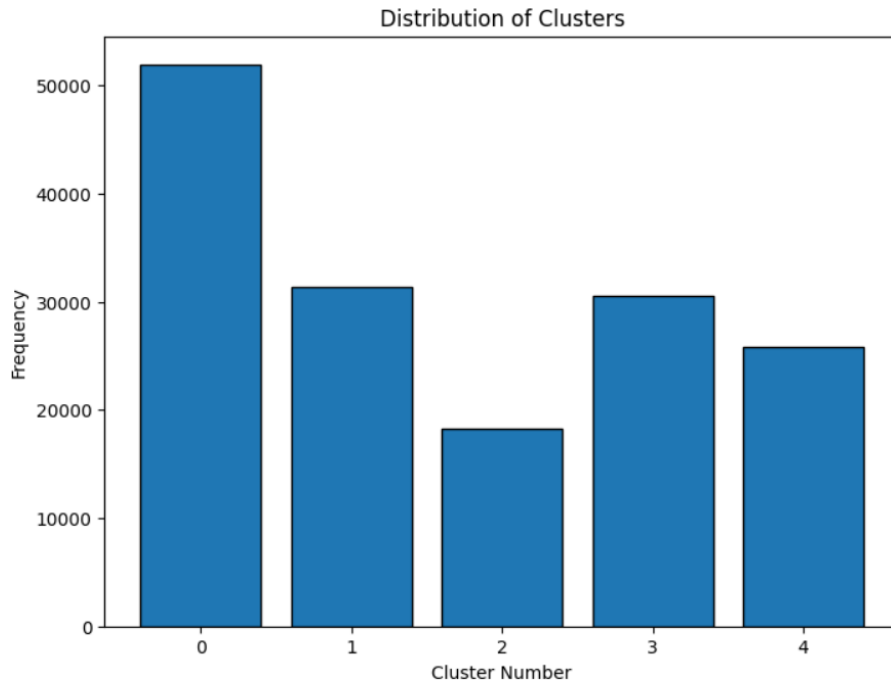# Clustering

**1- StandardScaler**

**2- choose best #clusters (WSS , Silhouette)**

Elbow Method for Optimal K



Silhouette Score for Optimal K

**the best number is not clear & we do not have info from the field --> we will choose 5**

**3 -Kmeans**

Distribution of Clusters

# Classification

**1- Split with stratification (biasedData)**

**2- TreeBased Models**

**RandomForest Results**

```
RandomForest AUC - Train: 0.6565381114570914, Validation:
0.647612960178587, Test: 0.6468449418542753
```

```
RandomForest Accuracy - Train: 0.9205224339274739, Validation:
0.921956926810244, Test: 0.9232210948103148
```

```
RandomForest F1 Score - Train: 0.882428173078826, Validation:
0.8845199005620704, Test: 0.8863642273919851
```

**Gradient Boosting Results**

```
Gradient Boosting AUC - Train: 0.7073215529220184, Validation:
0.677147390499939, Test: 0.6761514600234316
```

```
Gradient Boosting Accuracy - Train: 0.9205777504609711, Validation:
0.921913593623546, Test: 0.923199551907624
```

```
Gradient Boosting F1 Score - Train: 0.8825761677401183, Validation:
0.8844982692049035, Test: 0.8863534729873809
```

**Decision Tree Results**

```
Decision Tree AUC - Train: 0.4504909005748263, Validation:
0.44214254129043434, Test: 0.45441470539669887

Decision Tree Accuracy - Train: 0.9205593116164721, Validation:
0.921956926810244, Test: 0.923199551907624

Decision Tree F1 Score - Train: 0.8825549170153326, Validation:
0.884519905620704, Test: 0.8863534729873809
```

**XGBoost Results**

```
XGBoost AUC - Train: 0.5033, Validation: 0.5010, Test: 0.5010

XGBoost Accuracy - Train: 0.9210, Validation: 0.9220, Test: 0.9232

XGBoost F1 Score - Train: 0.8837, Validation: 0.8849, Test: 0.8867
```

## 3- LogReg

**Logistic Regression Results**

```
Logistic Regression AUC - Train: 0.6942, Validation: 0.6905, Test: 0.6859

Logistic Regression Accuracy - Train: 0.9204, Validation: 0.9219, Test:
0.9231

Logistic Regression F1 Score - Train: 0.8826, Validation: 0.8848, Test:
0.8865
```

## 4- K-fold CV (XGBoost)

**paramGrid**

```python
paramGrid = (ParamGridBuilder()

    .addGrid(xgb.max_depth, [6, 10])

    .addGrid(xgb.learning_rate, [0.1, 0.3])

    .addGrid(xgb.n_estimators, [100,200])

    .build())
```

**Best Model Results**

```
Cross-Validated XGBoost F1 - Train+Val: 0.9985, Test: 0.8875
```

# Enhancements and future work

Applying DL model and  Deployment