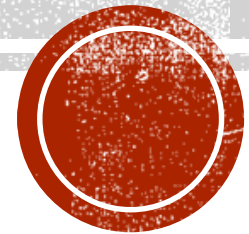


# HOME CREDIT DEFAULT RISK

Predict how capable a loan applicant is of repaying the loan.



# AGENDA



DATA UNDERSTANDING



DATA ANALYSIS AND VISUALIZATION



DATA PREPROCESSING



DATA MODELLING



DATA EVALUATION



FUTURE EXTENSION



DEMO



1



2



**TARGET AUDIENCE**



# MOTIVATION

- Whenever a customer has no credit history, lenders consider it risky to lend loans as there is no track record.
- Even a trustworthy candidate might be rejected.
- This forces the customers to lend loan from small vendors with high interest rate.
- The Bank loses a good customer and the customer deals with high rates





# DATA UNDERSTANDING





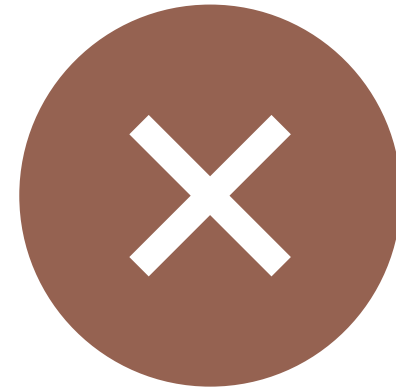
# DATASET



122 FEATURES WITH 307511  
ENTRIES



**SK\_ID\_CURR** IS USED AS A  
UNIQUE IDENTIFIER OF AN  
APPLICATION



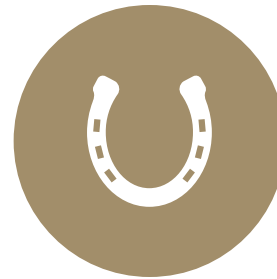
**TARGET** COLUMN INDICATES  
IF LOAN IS REPAYED OR NOT



# WHY DATA SCIENCE



Extracting information manually from the raw data could be tedious and error prone



High chance of mispredictions

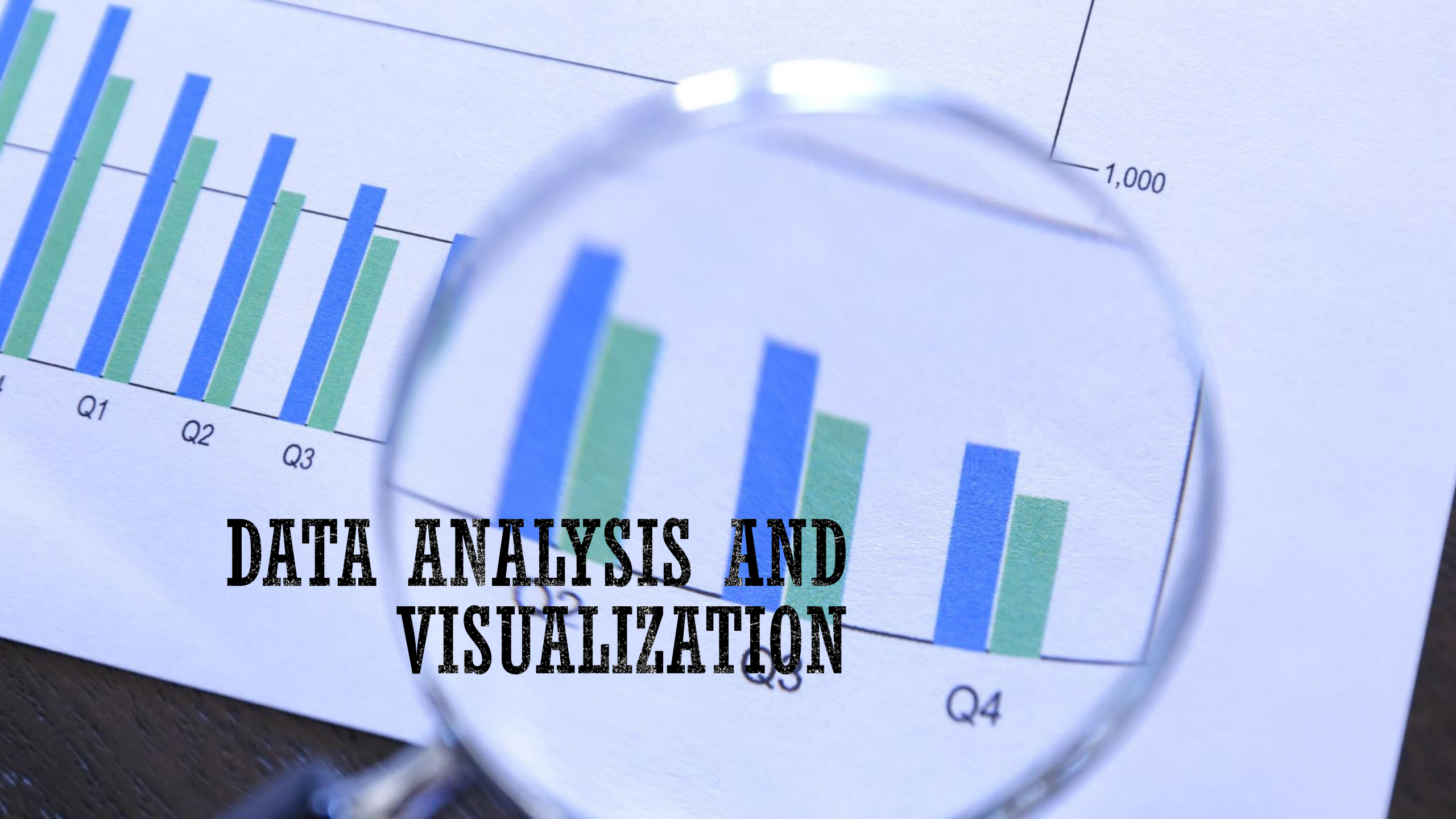


Decision should be based on previous knowledge



Good candidate of supervised classification problem





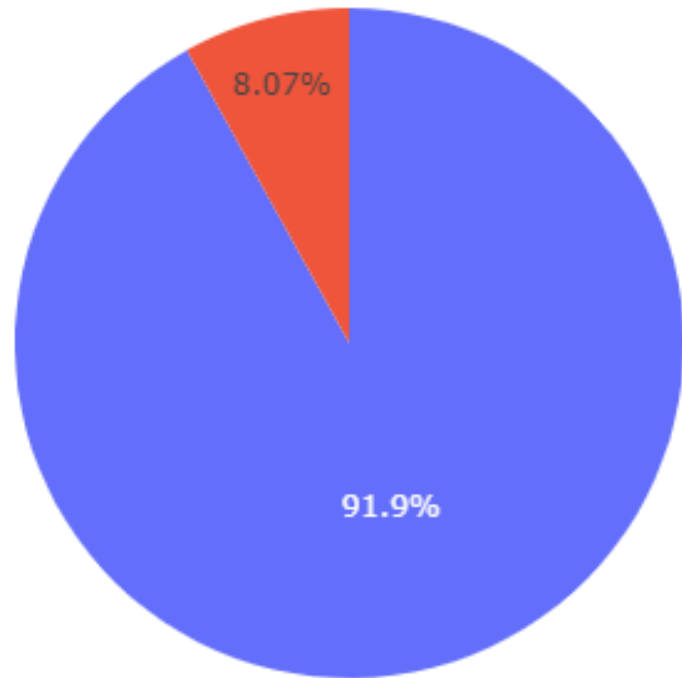
# DATA ANALYSIS AND VISUALIZATION



# DUPLICATES

No duplicates are  
present in the entire  
dataset.





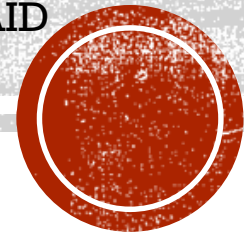
IMBALANCED DATA



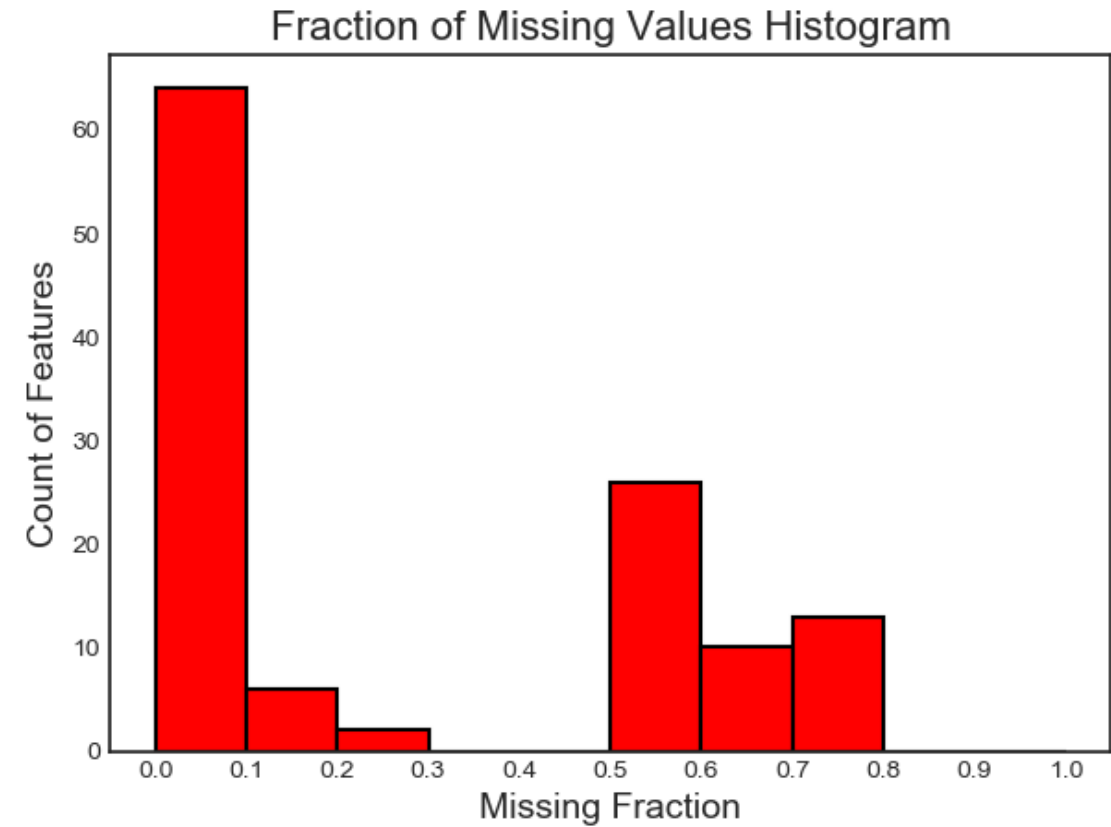
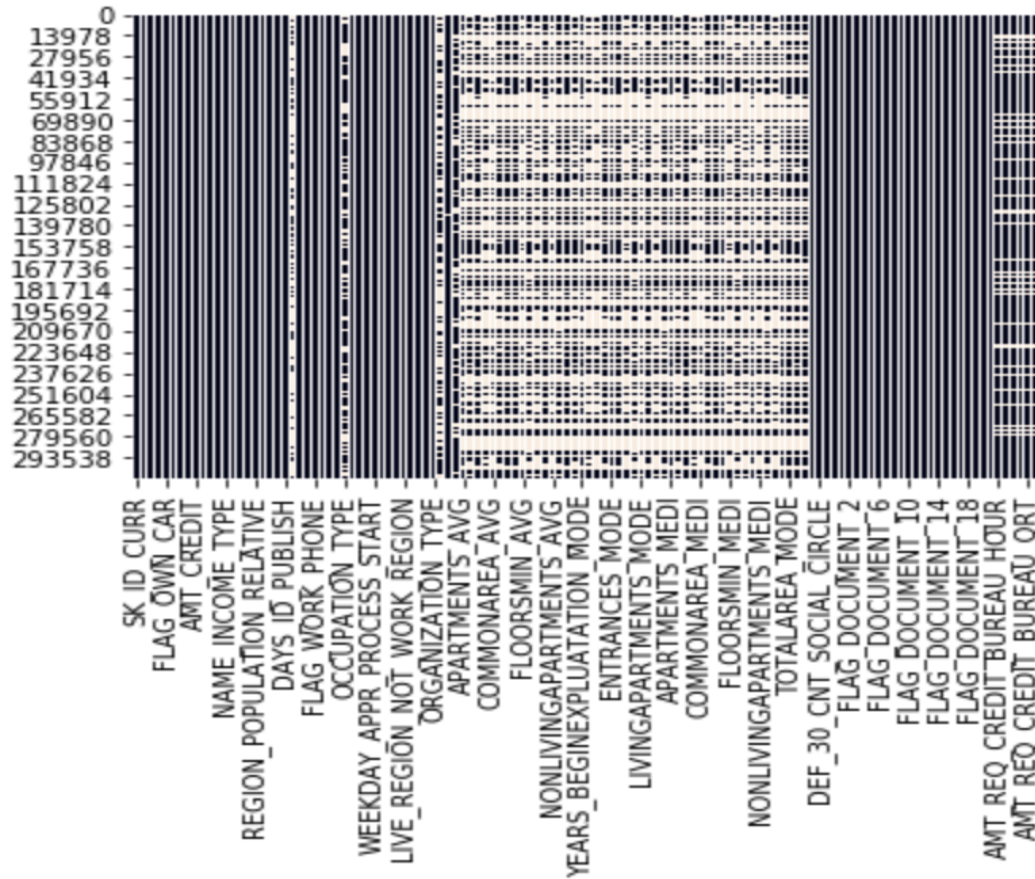
## CLASS DISTRIBUTION

0 - REPAID

1 - NOT REPAID



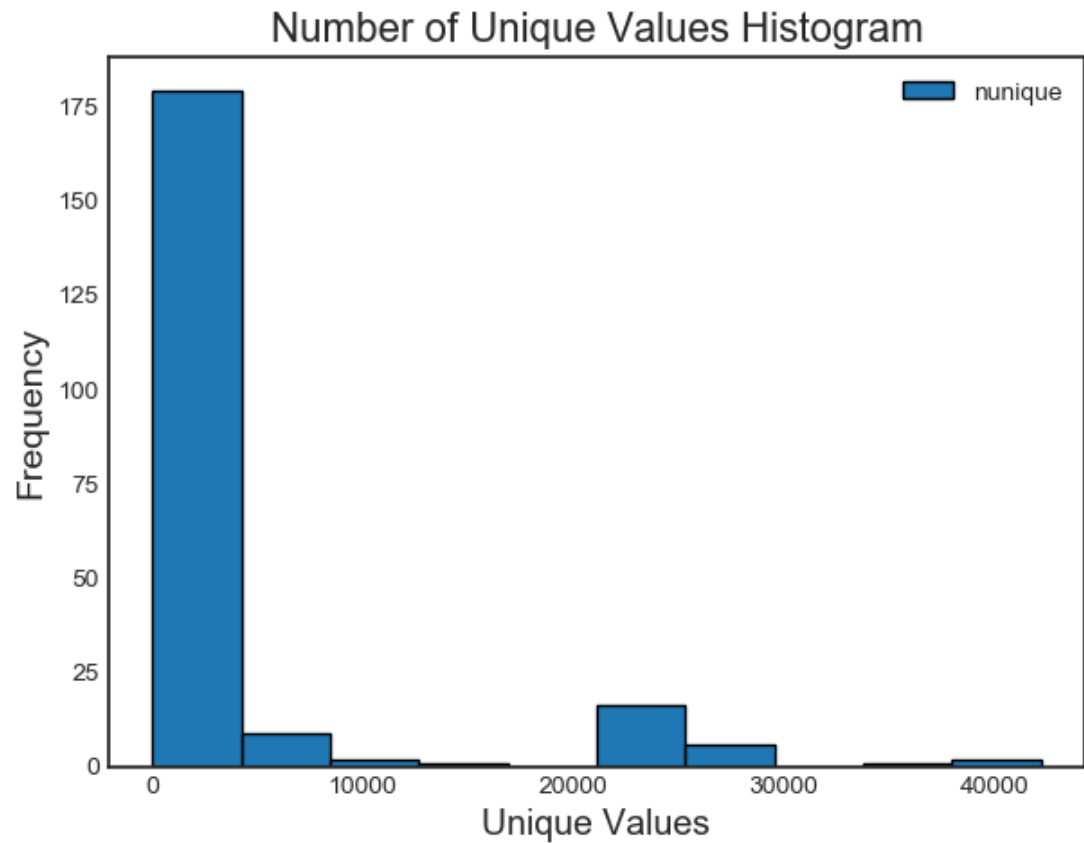
# MISSING VALUES



23 feature contains more than 60 percent of missing value



# SINGLE UNIQUE VALUE

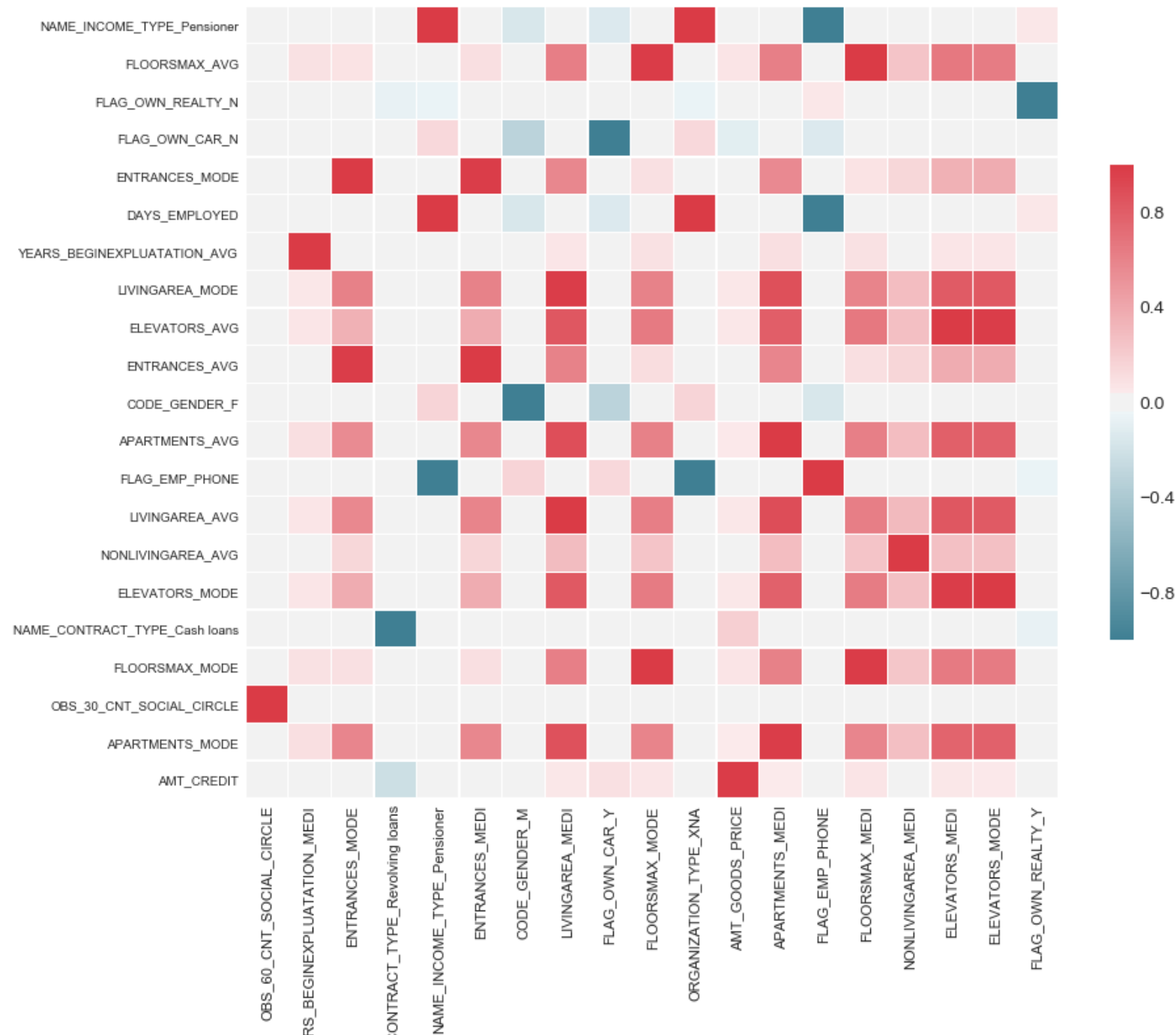


3 features with a single unique value





Correlations Above Threshold



# FEATURES CORRELATION

- Quantifies the degree to which a relationship between two variables
- Shows correlations above threshold 0.975
- There are 19 features with a correlation magnitude greater than the threshold
- Lighter green indicates the higher correlation.

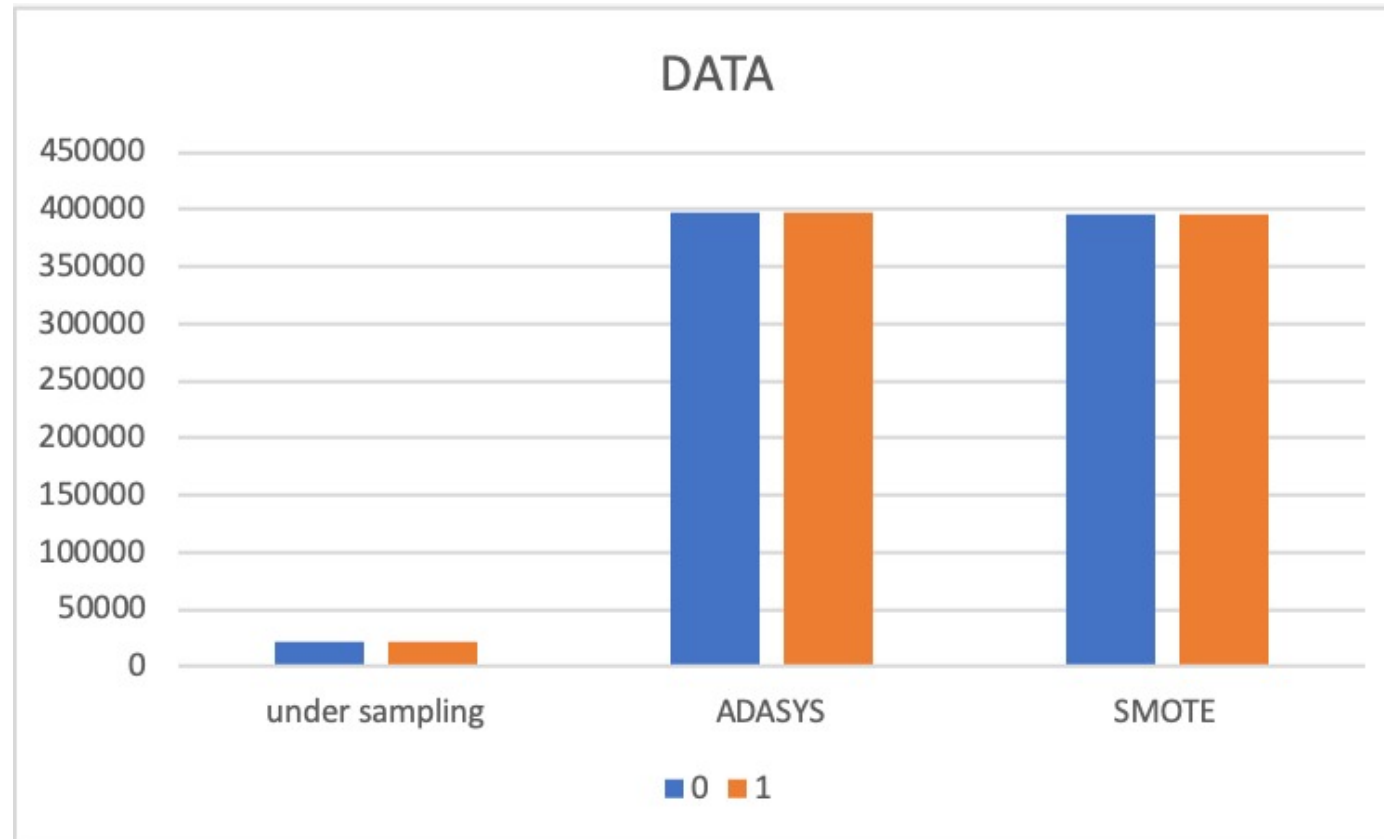






# DATA PREPROCESSING





**DEALING  
WITH  
IMBALANCED  
DATA**



**BEST RESULT : Under Sampling**

# FILLING MISSING VALUES

- Mean
- Median
- Mode
- Iterative Imputation

**BEST RESULT :** Iterative Imputation



# FEATURE SELECTION

---

Removed 23  
features which  
contain more than  
60 percent missing  
value

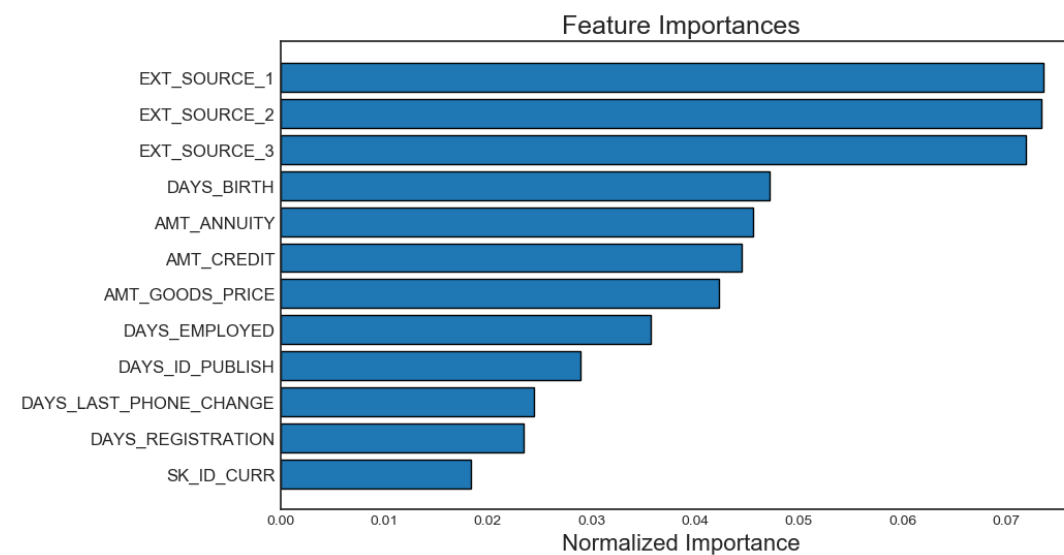
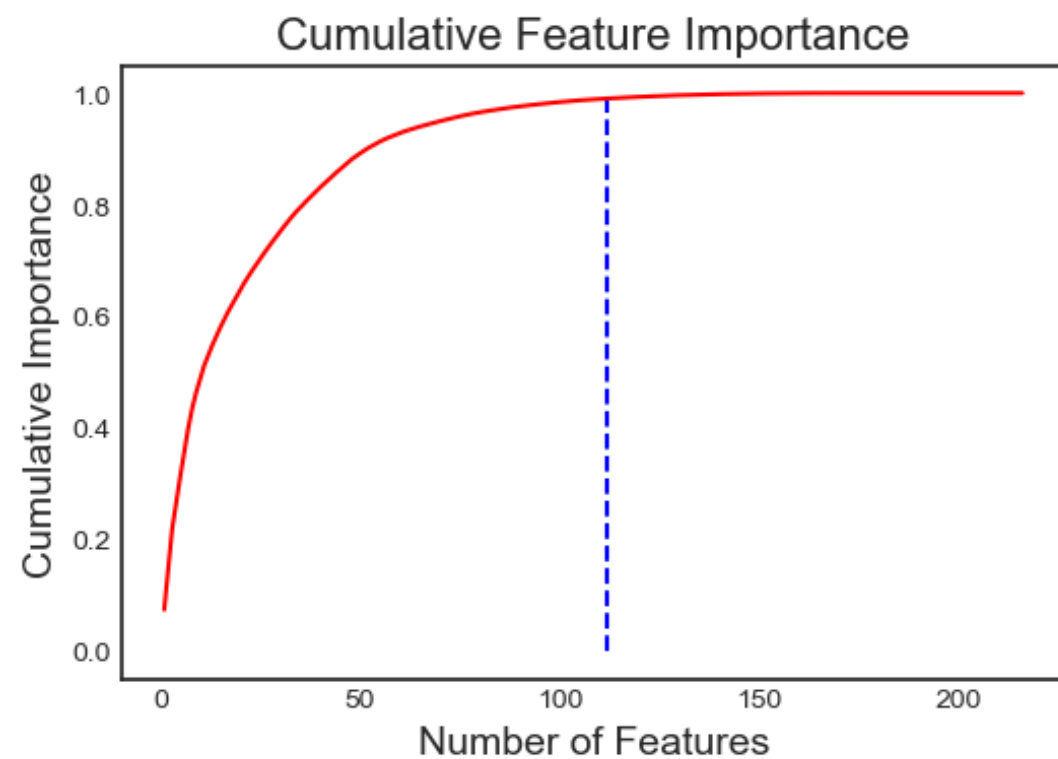
Removed 3 features  
with single unique  
value

Dropped highly  
correlated features

Removed 105  
features that do not  
contribute to  
cumulative  
importance of 0.99.

Removed 47  
features of zero  
importance after  
one-hot encoding.





Most important features



# ENCODING THE CATEGORICAL VARIABLE:

---



ONE HOT ENCODING



LABEL ENCODING





# DATA MODELING

# MODELS

**Logistic  
Regression**

**Gaussian Naive  
Bayes**

**Decision Tree  
Classifier**

**Random Forest**

**Bagging  
Classifier**

**Gradient  
Boosting**

**AdaBoost**

**Grid Search**

**LightGBM**

**Hyper-  
parameter  
Tuning using  
Grid Search**

## **Stacking**

- **Logistic Regression**
- **Bagging Classifier**
- **AdaBoost Classifier**
- **RandomForestClassifier**



## EXCEPTIONATION

Oversampling will give high accuracy

Tree-based models will perform better

Grid Search will improve performance

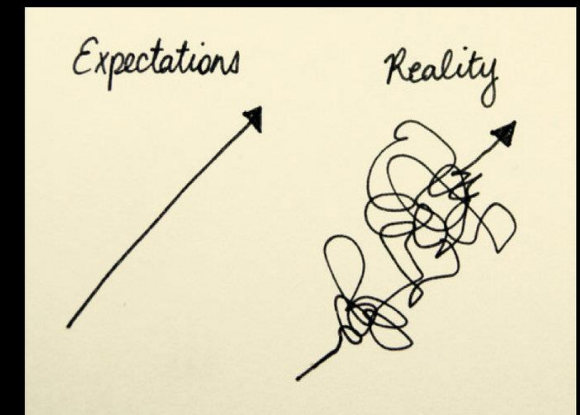
## REALITY

Under sampling proves out to be a better approach

Best result was given by Logistic Regression

Default parameters for Logistic Regression resulted in the same performance

# OBSERVATIONS





# DATA EVALUATION





# EVALUATION METRICS



**Accuracy**



**Precision**



**Recall**



**F1 score**



**ROC curve**



# WHICH METRIC TO USE ?

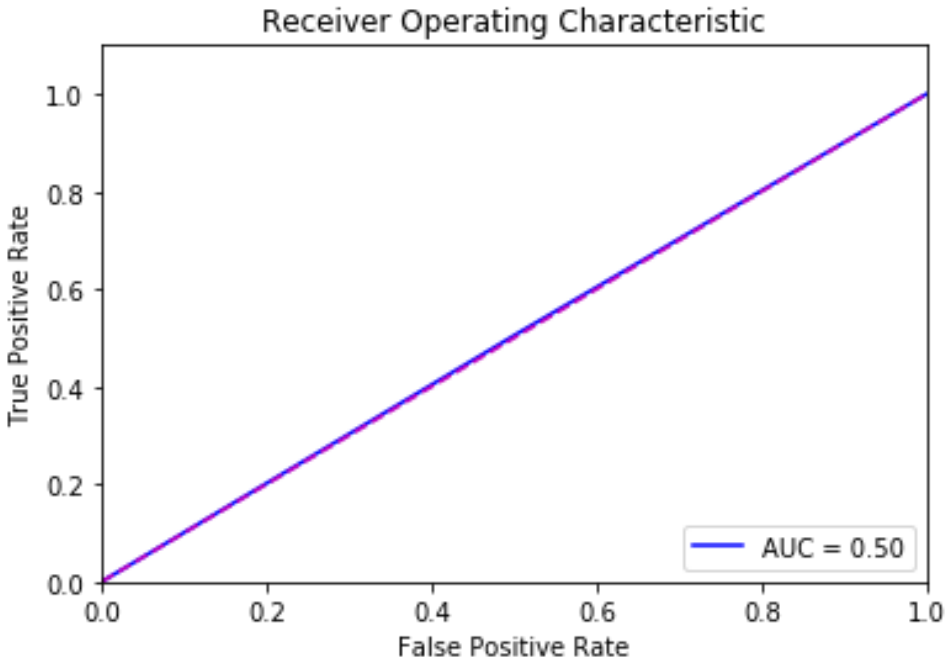
- A high precision value implies that the model returned substantially more relevant results than the irrelevant ones whereas a high recall value implies that the model returned most of the relevant results.
- Would like to have a model with higher recall value rather than the precision.
- Precision cannot be below a threshold to make sure the Home credit will not take every loaner as a person who can't have a consistent repay capability since a precision can indicate many false positives



# BEST EVALUATION METRIC FOR THE MODEL



Accuracy	ROC	Precision	Recall
0.498	0.502	0.09	0.54



**RANDOM MODEL**



# BUILDING THE CLASSIFIER

 Split the data (Stratified Shuffle Split cross-validator )

 Under sample

 One-Hot encoding of categorical features

 Impute missing values

 Scale (Standard Scaler)

 Feature Selection and filtering

 Classify



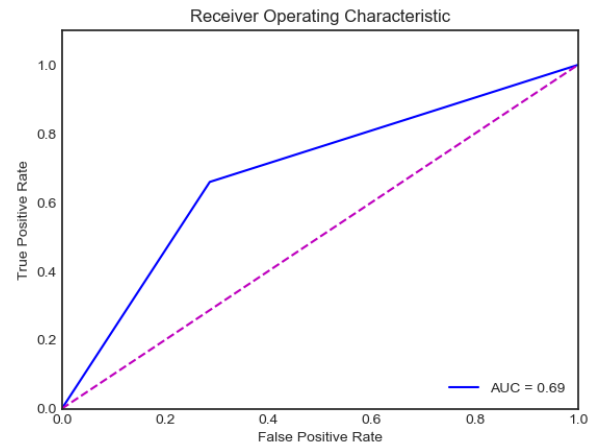


Model	Accuracy	ROC	Precision	Recall
Logistic Regression	0.688	0.682	0.16	0.68
Stacking	0.709	0.686	0.17	0.66

**BEST MODELS**

# AUC

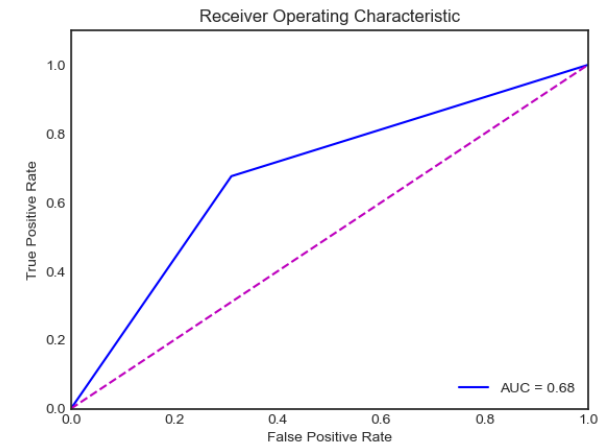
## Logistic Regression



Precision score : 0.16  
Recall score : 0.68

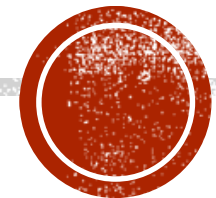
	precision	recall	f1-score	support
0	0.9604	0.6898	0.8029	42403
1	0.1607	0.6762	0.2597	3724

## Stacking



Precision score : 0.17  
Recall score : 0.66

	precision	recall	f1-score	support
0	0.9598	0.7135	0.8185	42403
1	0.1682	0.6595	0.2680	3724



**BEST MODEL**



**LOGISTIC REGRESSION**



# FUTURE EXTENSION

- Re-sampling skewed data by the scale of client's profile or other useful information available from the data.
- Collect additional data about our features in order to train the model better
- Explore the remaining datasets available to train the model better
- Improve UI
- Move the training to distributed environment and using cloud services



# TIME TO EVALUATE

Applicant 1: 'SK\_ID\_CURR': 100005,  
'NAME\_CONTRACT\_TYPE': 'Cash loans', 'CODE\_GENDER': 'M',  
'FLAG\_OWN\_CAR': 'N',  
'FLAG\_OWN\_REALTY': 'Y',  
'CNT\_CHILDREN': 0,  
'AMT\_INCOME\_TOTAL': 99000.0,  
'AMT\_CREDIT': 222768.0,  
'AMT\_ANNUITY': 17370.0,  
'AMT\_GOODS\_PRICE': 180000.0,  
'NAME\_TYPE\_SUITE': 'Unaccompanied',  
'NAME\_INCOME\_TYPE': 'Working',  
'NAME\_EDUCATION\_TYPE': 'Secondary / secondary special',  
'NAME\_FAMILY\_STATUS': 'Married',  
'NAME\_HOUSING\_TYPE': 'House / apartment',  
'REGION\_POPULATION\_RELATIVE': 0.035792000000000004,  
'DAYS\_BIRTH': -18064,  
'DAYS\_EMPLOYED': -4469,  
'DAYS\_REGISTRATION': -9118.0,  
'DAYS\_ID\_PUBLISH': -1623,



'OWN\_CAR\_AGE': nan, '  
FLAG\_MOBIL': 1,  
'FLAG\_EMP\_PHONE': 1,  
'FLAG\_WORK\_PHONE': 0,  
'FLAG\_CONT\_MOBILE': 1,  
'FLAG\_PHONE': 0,  
'FLAG\_EMAIL': 0,  
'OCCUPATION\_TYPE': 'Low-skill Laborers',  
'CNT\_FAM\_MEMBERS': 2.0,  
'REGION\_RATING\_CLIENT': 2,  
'REGION\_RATING\_CLIENT\_W\_CITY': 2,  
'WEEKDAY\_APPR\_PROCESS\_START': 'FRIDAY',  
'HOUR\_APPR\_PROCESS\_START': 9,  
'REG\_REGION\_NOT\_LIVE\_REGION': 0,  
'REG\_REGION\_NOT\_WORK\_REGION': 0,  
'LIVE\_REGION\_NOT\_WORK\_REGION': 0,  
'REG\_CITY\_NOT\_LIVE\_CITY': 0,  
'REG\_CITY\_NOT\_WORK\_CITY': 0,  
'LIVE\_CITY\_NOT\_WORK\_CITY': 0, '  
ORGANIZATION\_TYPE': 'Self-employed',  
'EXT\_SOURCE\_1': 0.5649902017969249,  
'EXT\_SOURCE\_2': 0.2916555320093651,  
'EXT\_SOURCE\_3': 0.4329616670974407,  
'APARTMENTS\_AVG': nan,

'BASEMENTAREA\_AVG': nan,  
'YEARS\_BEGINEXPLUATATION\_AVG': nan,  
'YEARS\_BUILD\_AVG': nan,  
'COMMONAREA\_AVG': nan,  
'ELEVATORS\_AVG': nan,  
'ENTRANCES\_AVG': nan,  
'FLOORSMAX\_AVG': nan,  
'FLOORSMIN\_AVG': nan,  
'LANDAREA\_AVG': nan,  
'LIVINGAPARTMENTS\_AVG': nan,  
'LIVINGAREA\_AVG': nan,  
'NONLIVINGAPARTMENTS\_AVG': nan,  
'NONLIVINGAREA\_AVG': nan,  
'APARTMENTS\_MODE': nan, '  
BASEMENTAREA\_MODE': nan,  
'YEARS\_BEGINEXPLUATATION\_MODE': nan,  
'YEARS\_BUILD\_MODE': nan,  
'COMMONAREA\_MODE': nan,  
'ELEVATORS\_MODE': nan,  
'ENTRANCES\_MODE': nan,  
'FLOORSMAX\_MODE': nan,  
'FLOORSMIN\_MODE': nan,  
'LANDAREA\_MODE': nan,  
'LIVINGAPARTMENTS\_MODE': nan,

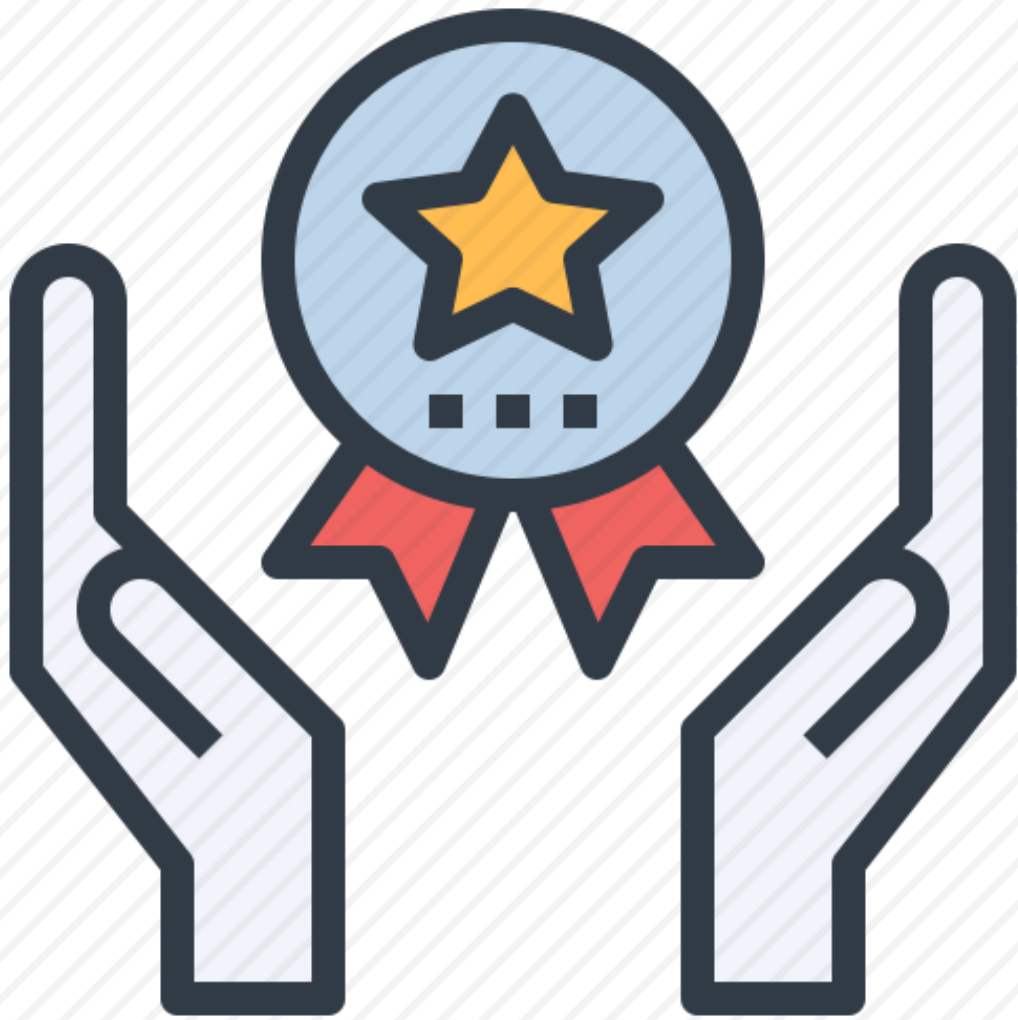




'LIVINGAREA\_MODE': nan,  
'NONLIVINGAPARTMENTS\_MODE': nan,  
'NONLIVINGAREA\_MODE': nan,  
'APARTMENTS\_MEDI': nan,  
'BASEMENTAREA\_MEDI': nan,  
'YEARS\_BEGINEXPLUATATION\_MEDI': nan,  
'YEARS\_BUILD\_MEDI': nan,  
'COMMONAREA\_MEDI': nan,  
'ELEVATORS\_MEDI': nan,  
'ENTRANCES\_MEDI': nan,  
'FLOORSMAX\_MEDI': nan,  
'FLOORSMIN\_MEDI': nan,  
'LANDAREA\_MEDI': nan,  
'LIVINGAPARTMENTS\_MEDI': nan,  
'LIVINGAREA\_MEDI': nan,  
'NONLIVINGAPARTMENTS\_MEDI': nan,  
'NONLIVINGAREA\_MEDI': nan,  
'FONDKAPREMONT\_MODE': nan,  
'HOUSETYPE\_MODE': nan,  
'TOTALAREA\_MODE': nan,  
'WALLSMATERIAL\_MODE': nan,  
'EMERGENCYSTATE\_MODE': nan,

'OBS\_30\_CNT\_SOCIAL\_CIRCLE': 0.0,  
'DEF\_30\_CNT\_SOCIAL\_CIRCLE': 0.0, '  
OBS\_60\_CNT\_SOCIAL\_CIRCLE': 0.0,  
'DEF\_60\_CNT\_SOCIAL\_CIRCLE': 0.0,  
'DAYS\_LAST\_PHONE\_CHANGE': 0.0,  
'FLAG\_DOCUMENT\_2': 0, 'FLAG\_DOCUMENT\_3': 1,  
'FLAG\_DOCUMENT\_4': 0, '  
FLAG\_DOCUMENT\_5': 0, 'FLAG\_DOCUMENT\_6': 0,  
'FLAG\_DOCUMENT\_7': 0,  
'FLAG\_DOCUMENT\_8': 0, 'FLAG\_DOCUMENT\_9': 0,  
'FLAG\_DOCUMENT\_10': 0,  
'FLAG\_DOCUMENT\_11': 0, 'FLAG\_DOCUMENT\_12': 0,  
'FLAG\_DOCUMENT\_13': 0,  
'FLAG\_DOCUMENT\_14': 0, 'FLAG\_DOCUMENT\_15': 0,  
'FLAG\_DOCUMENT\_16': 0, '  
FLAG\_DOCUMENT\_17': 0, 'FLAG\_DOCUMENT\_18': 0,  
'FLAG\_DOCUMENT\_19': 0,  
'FLAG\_DOCUMENT\_20': 0, 'FLAG\_DOCUMENT\_21': 0,  
'AMT\_REQ\_CREDIT\_BUREAU\_HOUR': 0.0,  
'AMT\_REQ\_CREDIT\_BUREAU\_DAY': 0.0,  
'AMT\_REQ\_CREDIT\_BUREAU\_WEEK': 0.0,  
'AMT\_REQ\_CREDIT\_BUREAU\_MON': 0.0,  
'AMT\_REQ\_CREDIT\_BUREAU\_QRT': 0.0,  
'AMT\_REQ\_CREDIT\_BUREAU\_YEAR': 3.0}





## BENEFITS

- Banks Risk factor reduced
- Deserving Candidates will not suffer
- One click solution
- Saves time





