

# Capstone Project Proposal

Udacity Machine Learning Nanodegree

James Dellinger

June 22, 2018

## 1 Domain Background

For my capstone project I will participate in the [Home Credit Default Risk competition](#) on Kaggle[Kagb]. The goal of the competition is to be able to predict the likelihood that an applicant will experience difficulty in repaying their loan. The competition is sponsored by [Home Credit](#)[Gro], whose mission is to provide a positive and safe borrowing experience to groups of people that traditional, mainstream banks and financial institutions typically refuse to serve. The competition's dataset is available for download directly from 'Data' section of its [Kaggle webpage](#)[Kaga]. There has been extensive research in using machine learning algorithms to predict loan repayment outcomes, such as a 2017 paper published by Xiaojiao Yu that explores using an XGBoost model to predict online lending risk: [Machine learning application in online lending risk prediction](#)[Yu17].

Home Credit targets a demographic that typically has no recourse but to deal with shady characters such as loan sharks when borrowing money. Many of these unbanked individuals are hard-working, well-intentioned folks who, either due to circumstances beyond their control or past mistakes, have fallen through the financial system's cracks.

I want to live in a world where a second chance is available to anyone who would make a good faith effort to make the most of it. This is why Home Credit's goal of providing a safe lending alternative to otherwise financially down-and-out folks resonates with me. I want to help Home Credit have the ability to expand its services to as many deserving applicants as possible.

## 2 Problem Statement

Home Credit needs an algorithm that will take as inputs various personal and alternative financial information originally taken from a loan applicant's profile, and then determine a probability of the applicant having at least one late payment when repaying their loan. This probability will be in the range  $[0.0, 1.0]$ , where 1.0 represents a 100% certainty that the applicant will have at least one delinquent repayment and 0.0 indicates that there is

zero chance that the applicant will ever be delinquent. The algorithm will be tested on a set of 48,744 individuals who previously borrowed from Home Credit. A CSV file must be produced that contains one header row, and 48,744 prediction rows, where each prediction row contains both a user ID, the `SKI_ID_CURR` column, and the probability, the `TARGET` column, of that user being delinquent. The file must be formatted as follows:

```
SK_ID_CURR , TARGET
100001 , 0.1
100005 , 0.9
100013 , 0.2
etc .
```

Home Credit knows which borrowers ultimately paid off their loans, and which ones had one or more late payments. A good algorithm will need to predict a high probability of delinquency for the majority of borrowers who did actually make one or more late payments. This algorithm will also need to predict a low probability of delinquency for the majority of borrowers who eventually did successfully repay their loans with no late payments.

### 3 Datasets and Inputs

This competition’s dataset was provided by Home Credit Group’s data scientists. It contains a wide variety of personal and financial information belonging to 356,255 individuals who had previously been recipients of loans from Home Credit. These individuals are divided into training and testing sets. The training group contains 307,511 individuals’ records. The test group contains 48,744 records. The dataset is anonymized, with each individual represented by their loan ID. Any personally indentifying infomation, such as name, phone number, or address, has been omitted.

Because Home Credit targets the unbanked population, it is unable to rely on traditional measures, such as a credit score, that mainstream financial institutions use when making lending decisions. Home Credit works around this obstacle by looking at an extensive and diverse array of personal and financial information for each of its applicants. These features range from common characteristics, such as marital status, age, type of housing, region of residence, job type, and education level, to some incredibly niche characteristics, such as how many elevators are in an applicant’s apartment building. Home Credit also looks at nitty gritty aspects of applicants’ financial backgrounds, including month-by-month payment performance on any loans or credit card balances that the applicant has previously had with Home Credit, as well as the amount and monthly repayment balances of any loans that the applicant may have received from other lenders.

All of these features are spread across seven data tables. The main data table (application\_{train | test}.csv) contains 120 features that comprise applicants’ personal background information. The other six data tables contain applicants’ previous loan and credit card balance payment histories. Detailed descriptions of each data table can be found in [Appendix A](#). Explanations of all features in each data table are in [Appendix B](#). The following diagram

provides a brief summary:

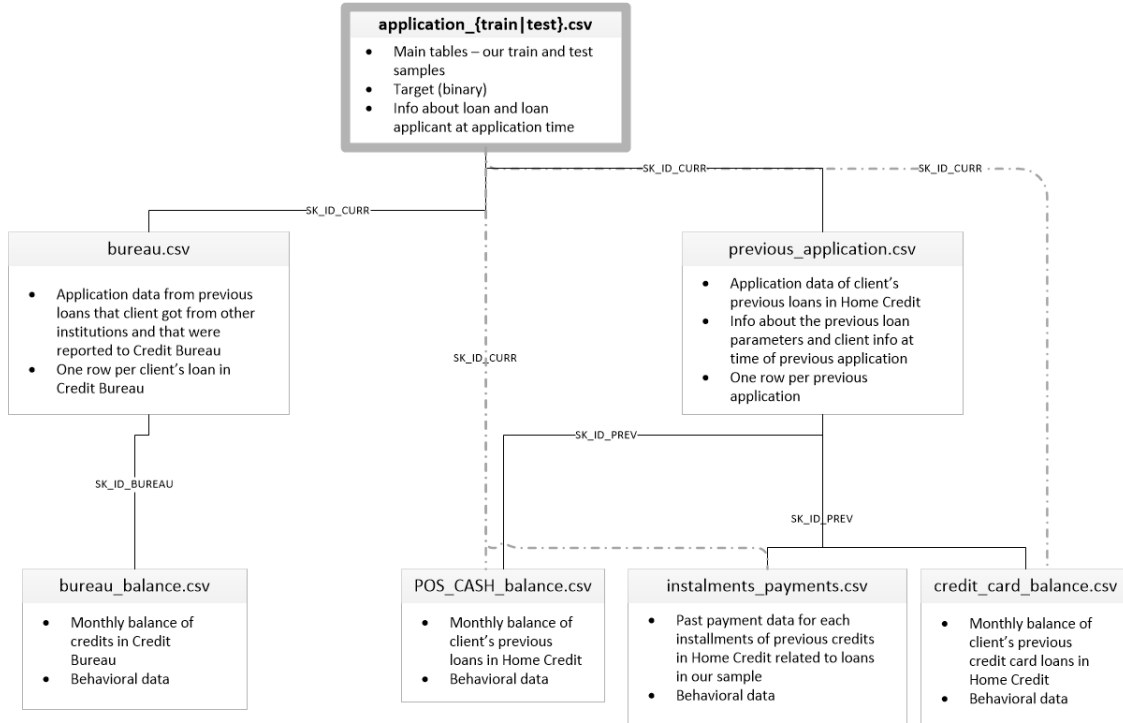


Figure 1: Home Credit Data Table Relationships[Kaga]

## 4 Solution Statement

The solution to this competition will be an algorithm that computes delinquency probabilities (in the range  $[0.0, 1.0]$ ) for each of the 48,744 Home Credit borrowers in the dataset's test segment. These predictions will be stored in a CSV file that adheres to the format described in the Problem Statement above in Section 2. This algorithm will take as inputs features from the seven data tables outlined above. It will be trained on the 307,511 Home Credit borrower records that comprise the training segment. The winning algorithm will predict the highest probabilities of delinquency for the borrowers who had at least one late payment, and will predict the lowest probabilities of delinquency for the borrowers who ultimately repaid their loans with no late payments. More formally, a winning algorithm will have the highest true positive rate, or sensitivity, (correctly guessing who will be delinquent), while at the same time having a minimum false positive rate (as seldom as possible making an incorrect guess that an individual will have difficulty paying back their loan).

## 5 Benchmark Model

Solutions to this competition will be ranked on Kaggle by the area under the ROC curve between their predicted probabilities and the observed targets (whether or not the borrower was actually delinquent).

I will train a Gaussian Naive Bayes classifier on a fully preprocessed dataset and use its ROC area under curve score as my primary benchmark. Gaussian Naive Bayes is ideal for this because it needs no hyperparameter input nor tuning. It also trains faster than most other models and I expect I should have no difficulty training it on the dataset's 300,000+ rows locally on my own machine. I am optimistic that with appropriate model selection and hyperparameter tuning I will eventually be able to build a learning algorithm that outperforms this out-of-the-box Gaussian Naive Bayes model.

In order to gauge how my solution compares to others in the Kaggle community, I will use the average ROC curve area of all entries on the competition's public leaderboard[Kage] as of June 2, 2018 as a secondary, personal benchmark. As of this date, there were 5,398 total entries from 1,515 participants. The average area under the ROC curves of all entries was 0.7379. I hope to design an algorithm that will meet or exceed this score.

It would also be nice to know the area under ROC curve of the prediction model employed by Home Credit's data scientists prior to this competition's debut on Kaggle. However, since Home Credit possibly considers this model and its ROC curve area to be proprietary information, it is understandable that it was not disclosed.

## 6 Evaluation Metric

The area under the ROC (receiver operating characteristic) curve[Wikb] will be the evaluation metric for my solution to this competition. This area can range from a minimum value of 0, or 0%, to a maximum of 1, or 100%. The size of the area under an ROC curve indicates how good a job a classifier does of identifying a certain target segment in a dataset. An area of 1 indicates that the classifier is perfect – that it can find every true positive (exhibiting perfect sensitivity, or recall) without making any mistakes (not accidentally labelling some true negatives as false positives), thus also giving it perfect specificity, which leads to a false positive rate (1 - specificity), of zero. An area of 0 indicates that the classifier isn't able to find and properly label any of the true positives. An area of 0.5 is the level of performance we'd expect from a classifier that randomly labels each point in the dataset – on the whole, this classifier would make just as many mistakes (false positives) as it makes correct predictions (true positives).

Because the goal of this competition is to create an algorithm that correctly identifies applicants who were delinquent in repaying their loans, area under the ROC curve is a perfectly appropriate measure of my algorithm's performance. While my target benchmark is the score of 0.7379 that I mentioned above, at a bare minimum, my classifier must achieve a score of better than 0.5 – anything less would mean that my classifier performs only as good as, or worse than, a random classifier.

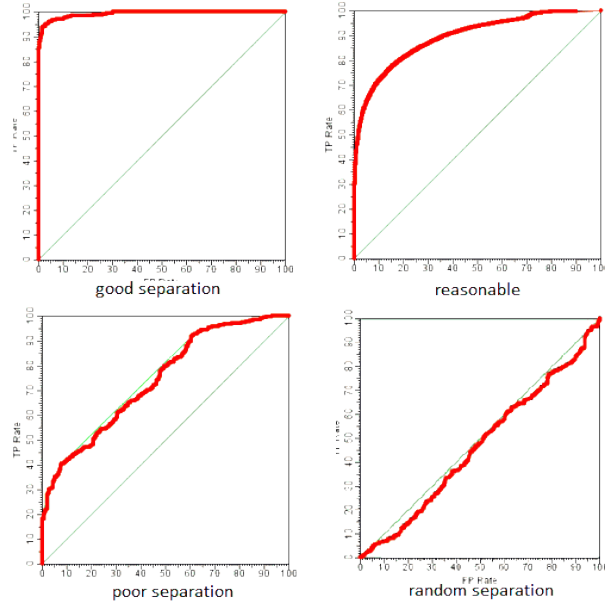


Figure 2: ROC Curve Examples[Wika]

## 7 Project Design

Although Home Credit provides seven data tables for this competition, for the purposes of this project I will focus most of my attention on the main data table. However, at a minimum, I will engineer at least one feature from the remaining six data tables. I may dive deeper into these other six data tables if time permits. I will structure my approach to creating a high-performing classification algorithm according to the steps in the following provisional workflow:

### I Data Exploration

- Compile lists and short descriptions of each feature in all seven of the competition's data tables.
- Display number of entries (rows) in each data table.
- Display number of features in each data table.
- Generate and display 5-row samples of each data table.
- Generate and display statistical descriptions of the numerical features of each data table that will be used by my algorithm.
- Investigate whether any features in any of the data tables have 'NaN' entries that will need to be handled, and whether some features have sparse enough data such that the entire feature should be removed.

- (g) Also investigate whether any borrowers (rows) in the main data table have ‘NaN’ for the majority of their feature entries. Some of these datapoints may also need to be removed.
- (h) Review the statistical descriptions of each numerical feature, and ensure that there are no other unexpected numerical entries in these features, eg. a feature having some negative numbers among its entries, where only positive numbers would be expected based on how the feature was defined by Home Credit.
- (i) List the numerical features in the main data table that are not already identified as being normally distributed (they may eventually need to be log-normalized).
- (j) List the categorical features in the main data table that will need to be one-hot encoded.
- (k) List the categorical features in the data table that are already one-hot encoded.
- (l) Define at least one feature that could be engineered from the data contained in the other six data tables outside of the main data table.

## II Exploratory Visualization

- (a) Plot histograms of the numerical features in the main data table, in order to confirm which ones are indeed already normally distributed.
- (b) Plot histograms of the numerical features in the main data table that are supposedly not already normally distributed.
- (c) Identify any features containing outliers that will need to be addressed.

## III Algorithms and Techniques to Consider

### (a) Dimensionality Reduction:

- i. Feature transformation and dimensionality reduction using Principle Component Analysis (PCA)
- ii. Dimensionality reduction using Sparse Random Projection (in case using PCA is too taxing on my compute resources)

### (b) Learning Algorithms:

Any algorithm that I use will need to not only be able to simply segment the data points into the “delinquent repayers” and “on-time repayers” categories, but must also be able to compute the probability that each data point in the test set belongs to the “delinquent repayers” segment. These probabilities will hopefully be higher for points that the algorithm labels as “delinquent repayers,” and will be lower for the points that are labeled as “on-time repayers.” This requirement automatically eliminates algorithms such as SVM, K-Means, and Hierarchical Clustering from consideration.

Below are several algorithms that meet the above criteria:

- i. Classifier Algorithms:
  - Multi-layer Perceptron Classifier
  - Logistic Regression Classifier
  - Gaussian Naive Bayes Classifier
- ii. Ensemble Methods:
  - AdaBoost Classifier
  - Bagging Classifier
  - Extra Trees Classifier
  - Gradient Boosting Classifier
- (c) **Other Techniques:**
  - i. Gridsearch CV for hyperparameter tuning.
  - ii. Breaking the training set up into an 80%-20% train-test split, or use K-Fold Cross-Validation, for model selection.
  - iii. SelectKBest for feature selection.

#### IV Data Preprocessing

- (a) Separate out training and testing target data into separate variables. Drop the targets column from both training and testing sets of the main data table.
- (b) Apply one-hot encoding to all non-numerical features in the main data table.
- (c) Apply non-linear feature scaling using the natural logarithm to any numerical features in the main data table that aren't already normally distributed.
- (d) Scale all numerical features are to the range  $[0.0,1.0]$ .
- (e) Impute 'NaN' values for numerical features.
- (f) Replace 'NaN' with 0 for all binary categorical features.
- (g) **Engineer** at least one feature based on some of the data contained in the six other data tables and append it to the main data table.
- (h) Build a preprocessing pipeline that can be used to apply the above steps in order to preprocess test sets before a classifier makes predictions on them.

#### V Implementation

- (a) **First and foremost**, it is absolutely necessary to set aside a validation segment from the training set in order to compare the different algorithms I attempt and to avoid overfitting. This should be about 20% of the size of the original training set. Alternatively, K-Fold Cross-Validation could also be used.
- (b) Run PCA or Sparse Random Projection on the features in the main data table.
- (c) Plot the feature weights and display the explained variances of the dimensions outputted by the chosen dimensionality reduction algorithm, in order to get an idea of how many dimensions to use when classifying the data points.

- (d) Create different classifiers, using some or all of the different classifier algorithms listed above.
- (e) Create a training and prediction pipeline that outputs probability predictions for each data point in test segment to a CSV file, and also calculates and outputs the area under the ROC curve of the predictions.

## VI Refinement

- (a) Experiment with adjusting hyperparameters on the best performing classifier.
- (b) Try feature selection methods such as SelectKBest.

## VII Model Evaluation and Validation

- (a) Generate a table that ranks the test set scores (area under ROC curve) of all the different algorithms and featuresets that I experiment with.
- (b) Plot ROC curves of the results from various classifier algorithms
- (c) Choose the algorithm/featureset combo that had the highest score on the testing set.
- (d) Conduct sensitivity analysis by manipulating some of the training data in order to ascertain the solution's robustness.
- (e) Save the testing set predictions to a CSV file, and submit on Kaggle.
- (f) Compare my final algorithm's score to the benchmark score, and observe my submission's public rank on Kaggle.

# Appendices

## A Data Table Descriptions

From the [data description page](#) on the Home Credit Kaggle competition's website[Kaga].

### 1 application\_{train | test}.csv

- This is the main table, broken into two files for Train (with TARGET) and Test (without TARGET).
- Static data for all applications. One row represents one loan in our data sample.

### 2 bureau.csv

- All client's previous credits provided by other financial institutions that were reported to Credit Bureau (for clients who have a loan in our sample).



- For every loan in our sample, there are as many rows as number of credits the client had in Credit Bureau before the application date.

### 3 **bureau\_balance.csv**

- Monthly balances of previous credits in Credit Bureau.
- This table has one row for each month of history of every previous credit reported to Credit Bureau – i.e the table has (#loans in sample # of relative previous credits # of months where we have some history observable for the previous credits) rows.

### 4 **previous\_application.csv**

- All previous applications for Home Credit loans of clients who have loans in our sample.
- There is one row for each previous application related to loans in our data sample.

### 5 **POS\_CASH\_balance.csv**

- Monthly balance snapshots of previous POS (point of sales) and cash loans that the applicant had with Home Credit.
- This table has one row for each month of history of every previous credit in Home Credit (consumer credit and cash loans) related to loans in our sample – i.e. the table has (#loans in sample # of relative previous credits # of months in which we have some history observable for the previous credits) rows.

### 6 **installments\_payments.csv**

- Repayment history for the previously disbursed credits in Home Credit related to the loans in our sample.
- There is a) one row for every payment that was made plus b) one row each for missed payment.
- One row is equivalent to one payment of one installment OR one installment corresponding to one payment of one previous Home Credit credit related to loans in our sample.

### 7 **credit\_card\_balance.csv**

- Monthly balance snapshots of previous credit cards that the applicant has with Home Credit.
- This table has one row for each month of history of every previous credit in Home Credit (consumer credit and cash loans) related to loans in our sample – i.e. the table has (#loans in sample # of relative previous credit cards # of months where we have some history observable for the previous credit card) rows.

## B Data Table Features

From the HomeCredit\_columns\_description.csv file on the [data description page](#) on the Home Credit Kaggle competition's website[Kaga].

### B.1 Main Data Table Features (application\_{train | test}.csv)

- 1 **SK\_ID\_CURR**: ID of loan in our sample
- 2 **TARGET**: Target variable (1 - client with payment difficulties: he/she had late payment more than X days on at least one of the first Y installments of the loan in our sample, 0 - all other cases)
- 3 **NAME\_CONTRACT\_TYPE**: Identification if loan is cash or revolving
- 4 **CODE\_GENDER**: Gender of the client
- 5 **FLAG\_OWN\_CAR**: Flag if the client owns a car
- 6 **FLAG\_OWN\_REALTY**: Flag if client owns a house or flat
- 7 **CNT\_CHILDREN**: Number of children the client has
- 8 **AMT\_INCOME\_TOTAL**: Income of the client
- 9 **AMT\_CREDIT**: Credit amount of the loan
- 10 **AMT\_ANNUITY**: Loan annuity
- 11 **AMT\_GOODS\_PRICE**: For consumer loans it is the price of the goods for which the loan is given
- 12 **NAME\_TYPE\_SUITE**: Who was accompanying client when he was applying for the loan
- 13 **NAME\_INCOME\_TYPE**: Clients income type (businessman, working, maternity leave,Ö)
- 14 **NAME\_EDUCATION\_TYPE**: Level of highest education the client achieved
- 15 **NAME\_FAMILY\_STATUS**: Family status of the client
- 16 **NAME\_HOUSING\_TYPE**: What is the housing situation of the client (renting, living with parents, ...)
- 17 **REGION\_POPULATION\_RELATIVE**: Normalized population of region where client lives (higher number means the client lives in more populated region) – normalized

- 18 **DAYS\_BIRTH**: Client's age in days at the time of application – time only relative to the application
- 19 **DAYS\_EMPLOYED**: How many days before the application the person started current employment – time only relative to the application
- 20 **DAYS\_REGISTRATION**: How many days before the application did client change his registration – time only relative to the application
- 21 **DAYS\_ID\_PUBLISH**: How many days before the application did client change the identity document with which he applied for the loan – time only relative to the application
- 22 **OWN\_CAR\_AGE**: Age of client's car
- 23 **FLAG\_MOBIL**: Did client provide mobile phone (1=YES, 0=NO)
- 24 **FLAG\_EMP\_PHONE**: Did client provide work phone (1=YES, 0=NO)
- 25 **FLAG\_WORK\_PHONE**: Did client provide home phone (1=YES, 0=NO)
- 26 **FLAG\_CONT\_MOBILE**: Was mobile phone reachable (1=YES, 0=NO)
- 27 **FLAG\_PHONE**: Did client provide home phone (1=YES, 0=NO)
- 28 **FLAG\_EMAIL**: Did client provide email (1=YES, 0=NO)
- 29 **OCCUPATION\_TYPE**: What kind of occupation does the client have
- 30 **CNT\_FAM\_MEMBERS**: How many family members does client have
- 31 **REGION\_RATING\_CLIENT**: Our rating of the region where client lives (1,2,3)
- 32 **REGION\_RATING\_CLIENT\_W\_CITY**: Our rating of the region where client lives with taking city into account (1,2,3)
- 33 **WEEKDAY\_APPR\_PROCESS\_START**: On which day of the week did the client apply for the loan
- 34 **HOUR\_APPR\_PROCESS\_START**: Approximately at what hour did the client apply for the loan rounded
- 35 **REG\_REGION\_NOT\_LIVE\_REGION**: Flag if client's permanent address does not match contact address (1=different, 0=same, at region level)
- 36 **REG\_REGION\_NOT\_WORK\_REGION**: Flag if client's permanent address does not match work address (1=different, 0=same, at region level)

- 37 **LIVE\_REGION\_NOT\_WORK\_REGION**: Flag if client's contact address does not match work address (1=different, 0=same, at region level)
- 38 **REG\_CITY\_NOT\_LIVE\_CITY**: Flag if client's permanent address does not match contact address (1=different, 0=same, at city level)
- 39 **REG\_CITY\_NOT\_WORK\_CITY**: Flag if client's permanent address does not match work address (1=different, 0=same, at city level)
- 40 **LIVE\_CITY\_NOT\_WORK\_CITY**: Flag if client's contact address does not match work address (1=different, 0=same, at city level)
- 41 **ORGANIZATION\_TYPE**: Type of organization where client works
- 42 **EXT\_SOURCE\_1**: Normalized score from external data source – normalized
- 43 **EXT\_SOURCE\_2**: Normalized score from external data source – normalized
- 44 **EXT\_SOURCE\_3**: Normalized score from external data source – normalized
- 45 **APARTMENTS\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 46 **BASEMENTAREA\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 47 **YEARS\_BEGINEXPLUATATION\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 48 **YEARS\_BUILD\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 49 **COMMONAREA\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized

- 50 **ELEVATORS\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 51 **ENTRANCES\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 52 **FLOORSMAX\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 53 **FLOORSMIN\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 54 **LANDAREA\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 55 **LIVINGAPARTMENTS\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 56 **LIVINGAREA\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 57 **NONLIVINGAPARTMENTS\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 58 **NONLIVINGAREA\_AVG**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized

- 59 **APARTMENTS\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 60 **BASEMENTAREA\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 61 **YEARS\_BEGINEXPLUATATION\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 62 **YEARS\_BUILD\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 63 **COMMONAREA\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 64 **ELEVATORS\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 65 **ENTRANCES\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 66 **FLOORSMAX\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 67 **FLOORSMIN\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized

- 68 **LANDAREA\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 69 **LIVINGAPARTMENTS\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 70 **LIVINGAREA\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 71 **NONLIVINGAPARTMENTS\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 72 **NONLIVINGAREA\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 73 **APARTMENTS\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 74 **BASEMENTAREA\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 75 **YEARS\_BEGINEXPLUATATION\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 76 **YEARS\_BUILD\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized

- 77 **COMMONAREA\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 78 **ELEVATORS\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 79 **ENTRANCES\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 80 **FLOORSMAX\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 81 **FLOORSMIN\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 82 **LANDAREA\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 83 **LIVINGAPARTMENTS\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 84 **LIVINGAREA\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 85 **NONLIVINGAPARTMENTS\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized



- 86 **NONLIVINGAREA\_MEDI**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 87 **FONDKAPREMONT\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 88 **HOUSETYPE\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 89 **TOTALAREA\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 90 **WALLSMATERIAL\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 91 **EMERGENCYSTATE\_MODE**: Normalized information about building where the client lives, What is average (\_AVG suffix), modus (\_MODE suffix), median (\_MEDI suffix) apartment size, common area, living area, age of building, number of elevators, number of entrances, state of the building, number of floor – normalized
- 92 **OBS\_30\_CNT\_SOCIAL\_CIRCLE**: How many observation of client's social surroundings with observable 30 DPD (days past due) default
- 93 **DEF\_30\_CNT\_SOCIAL\_CIRCLE**: How many observation of client's social surroundings defaulted on 30 DPD (days past due)
- 94 **OBS\_60\_CNT\_SOCIAL\_CIRCLE**: How many observation of client's social surroundings with observable 60 DPD (days past due) default
- 95 **DEF\_60\_CNT\_SOCIAL\_CIRCLE**: How many observation of client's social surroundings defaulted on 60 (days past due) DPD
- 96 **DAYS\_LAST\_PHONE\_CHANGE**: How many days before application did client change phone
- 97 **FLAG\_DOCUMENT\_2**: Did client provide document 2

- 98 **FLAG\_DOCUMENT\_3**: Did client provide document 3
- 99 **FLAG\_DOCUMENT\_4**: Did client provide document 4
- 100 **FLAG\_DOCUMENT\_5**: Did client provide document 5
- 101 **FLAG\_DOCUMENT\_6**: Did client provide document 6
- 102 **FLAG\_DOCUMENT\_7**: Did client provide document 7
- 103 **FLAG\_DOCUMENT\_8**: Did client provide document 8
- 104 **FLAG\_DOCUMENT\_9**: Did client provide document 9
- 105 **FLAG\_DOCUMENT\_10**: Did client provide document 10
- 106 **FLAG\_DOCUMENT\_11**: Did client provide document 11
- 107 **FLAG\_DOCUMENT\_12**: Did client provide document 12
- 108 **FLAG\_DOCUMENT\_13**: Did client provide document 13
- 109 **FLAG\_DOCUMENT\_14**: Did client provide document 14
- 110 **FLAG\_DOCUMENT\_15**: Did client provide document 15
- 111 **FLAG\_DOCUMENT\_16**: Did client provide document 16
- 112 **FLAG\_DOCUMENT\_17**: Did client provide document 17
- 113 **FLAG\_DOCUMENT\_18**: Did client provide document 18
- 114 **FLAG\_DOCUMENT\_19**: Did client provide document 19
- 115 **FLAG\_DOCUMENT\_20**: Did client provide document 20
- 116 **FLAG\_DOCUMENT\_21**: Did client provide document 21
- 117 **AMT\_REQ\_CREDIT\_BUREAU\_HOUR**: Number of enquiries to Credit Bureau about the client one hour before application
- 118 **AMT\_REQ\_CREDIT\_BUREAU\_DAY**: Number of enquiries to Credit Bureau about the client one day before application (excluding one hour before application)
- 119 **AMT\_REQ\_CREDIT\_BUREAU\_WEEK**: Number of enquiries to Credit Bureau about the client one week before application (excluding one day before application)
- 120 **AMT\_REQ\_CREDIT\_BUREAU\_MON**: Number of enquiries to Credit Bureau about the client one month before application (excluding one week before application)

- 121 **AMT\_REQ\_CREDIT\_BUREAU\_QRT**: Number of enquiries to Credit Bureau about the client 3 month before application (excluding one month before application)
- 122 **AMT\_REQ\_CREDIT\_BUREAU\_YEAR**: Number of enquiries to Credit Bureau about the client one day year (excluding last 3 months before application)

## B.2 Bureau Data Table Features (**bureau.csv**)

- 1 **SK\_ID\_CURR**: ID of loan in our sample - one loan in our sample can have 0,1,2 or more related previous credits in credit bureau – hashed
- 2 **SK\_ID\_BUREAU**: Recoded ID of previous Credit Bureau credit related to our loan (unique coding for each loan application) – hashed
- 3 **CREDIT\_ACTIVE**: Status of the Credit Bureau (CB) reported credits
- 4 **CREDIT\_CURRENCY**: Recoded currency of the Credit Bureau credit – recoded
- 5 **DAYS\_CREDIT**: How many days before current application did client apply for Credit Bureau credit – time only relative to the application
- 6 **CREDIT\_DAY\_OVERDUE**: Number of days past due on CB credit at the time of application for related loan in our sample
- 7 **DAYS\_CREDIT\_ENDDATE**: Remaining duration of CB credit (in days) at the time of application in Home Credit – time only relative to the application
- 8 **DAYS\_ENDDATE\_FACT**: Days since CB credit ended at the time of application in Home Credit (only for closed credit) – time only relative to the application
- 9 **AMT\_CREDIT\_MAX\_OVERDUE**: Maximal amount overdue on the Credit Bureau credit so far (at application date of loan in our sample)
- 10 **CNT\_CREDIT\_PROLONG**: How many times was the Credit Bureau credit prolonged
- 11 **AMT\_CREDIT\_SUM**: Current credit amount for the Credit Bureau credit
- 12 **AMT\_CREDIT\_SUM\_DEBT**: Current debt on Credit Bureau credit
- 13 **AMT\_CREDIT\_SUM\_LIMIT**: Current credit limit of credit card reported in Credit Bureau
- 14 **AMT\_CREDIT\_SUM\_OVERDUE**: Current amount overdue on Credit Bureau credit
- 15 **CREDIT\_TYPE**: Type of Credit Bureau credit (Car, cash,...)

- 16 **DAYS\_CREDIT\_UPDATE**: How many days before loan application did last information about the Credit Bureau credit come – time only relative to the application
- 17 **AMT\_ANNUITY**: Annuity of the Credit Bureau credit

### B.3 Bureau Balance Data Table Features (**bureau\_balance.csv**)

- 1 **SK\_ID\_BUREAU**: Recoded ID of Credit Bureau credit (unique coding for each application) - use this to join to CREDIT\_BUREAU table – hashed
- 2 **MONTHS\_BALANCE**: Month of balance relative to application date (-1 means the freshest balance date) – time only relative to the application
- 3 **STATUS**: Status of Credit Bureau loan during the month (active, closed, DPD0-30,Ö [C means closed, X means status unknown, 0 means no DPD, 1 means maximal did during month between 1-30, 2 means DPD 31-60,Ö 5 means DPD 120+ or sold or written off ])

### B.4 Previous Application Data Table Features (**previous\_application.csv**)

- 1 **SK\_ID\_PREV**: ID of previous credit in Home credit related to loan in our sample. (One loan in our sample can have 0,1,2 or more previous loan applications in Home Credit, previous application could, but not necessarily have to lead to credit) – hashed
- 2 **SK\_ID\_CURR**: ID of loan in our sample – hashed
- 3 **NAME\_CONTRACT\_TYPE**: Contract product type (Cash loan, consumer loan [POS] ,...) of the previous application
- 4 **AMT\_ANNUITY**: Annuity of previous application
- 5 **AMT\_APPLICATION**: For how much credit did client ask on the previous application
- 6 **AMT\_CREDIT**: Final credit amount on the previous application. This differs from AMT\_APPLICATION in a way that the AMT\_APPLICATION is the amount for which the client initially applied for, but during our approval process he could have received different amount - AMT\_CREDIT
- 7 **AMT\_DOWN\_PAYMENT**: Down payment on the previous application
- 8 **AMT\_GOODS\_PRICE**: Goods price of good that client asked for (if applicable) on the previous application

- 9 **WEEKDAY\_APPR\_PROCESS\_START**: On which day of the week did the client apply for previous application
- 10 **HOURLY\_APPR\_PROCESS\_START**: Approximately at what day hour did the client apply for the previous application – rounded
- 11 **FLAG\_LAST\_APPL\_PER\_CONTRACT**: Flag if it was last application for the previous contract. Sometimes by mistake of client or our clerk there could be more applications for one single contract
- 12 **NFLAG\_LAST\_APPL\_IN\_DAY**: Flag if the application was the last application per day of the client. Sometimes clients apply for more applications a day. Rarely it could also be error in our system that one application is in the database twice
- 13 **RATE\_DOWN\_PAYMENT**: Down payment rate normalized on previous credit – normalized
- 14 **RATE\_INTEREST\_PRIMARY**: Interest rate normalized on previous credit – normalized
- 15 **RATE\_INTEREST\_PRIVILEGED**: Interest rate normalized on previous credit – normalized
- 16 **NAME\_CASH\_LOAN\_PURPOSE**: Purpose of the cash loan
- 17 **NAME\_CONTRACT\_STATUS**: Contract status (approved, cancelled, ...) of previous application
- 18 **DAYS\_DECISION**: Relative to current application when was the decision about previous application made time only relative to the application
- 19 **NAME\_PAYMENT\_TYPE**: Payment method that client chose to pay for the previous application
- 20 **CODE\_REJECT\_REASON**: Why was the previous application rejected
- 21 **NAME\_TYPE\_SUITE**: Who accompanied client when applying for the previous application
- 22 **NAME\_CLIENT\_TYPE**: Was the client old or new client when applying for the previous application
- 23 **NAME\_GOODS\_CATEGORY**: What kind of goods did the client apply for in the previous application
- 24 **NAME\_PORTFOLIO**: Was the previous application for CASH, POS, CAR, Ö
- 25 **NAME\_PRODUCT\_TYPE**: Was the previous application x-sell o walk-in

- 26 **CHANNEL\_TYPE**: Through which channel we acquired the client on the previous application
- 27 **SELLERPLACE\_AREA**: Selling area of seller place of the previous application
- 28 **NAME\_SELLER\_INDUSTRY**: The industry of the seller
- 29 **CNT\_PAYMENT**: Term of previous credit at application of the previous application
- 30 **NAME\_YIELD\_GROUP**: Grouped interest rate into small medium and high of the previous application – grouped
- 31 **PRODUCT\_COMBINATION**: Detailed product combination of the previous application
- 32 **DAYS\_FIRST\_DRAWING**: Relative to application date of current application when was the first disbursement of the previous application – time only relative to the application
- 33 **DAYS\_FIRST\_DUE**: Relative to application date of current application when was the first due supposed to be of the previous application – time only relative to the application
- 34 **DAYS\_LAST\_DUE\_1ST\_VERSION**: Relative to application date of current application when was the first due of the previous application – time only relative to the application
- 35 **DAYS\_LAST\_DUE**: Relative to application date of current application when was the last due date of the previous application – time only relative to the application
- 36 **DAYS\_TERMINATION**: Relative to application date of current application when was the expected termination of the previous application – time only relative to the application
- 37 **NFLAG\_INSURED\_ON\_APPROVAL**: Did the client requested insurance during the previous application

## **B.5 POS CASH Balance Data Table Features (POS\_CASH\_balance.csv)**

- 1 **SK\_ID\_PREV**: ID of previous credit in Home Credit related to loan in our sample. (One loan in our sample can have 0,1,2 or more previous loans in Home Credit)
- 2 **SK\_ID\_CURR**: ID of loan in our sample

- 3 **MONTHS\_BALANCE**: Month of balance relative to application date (-1 means the information to the freshest monthly snapshot, 0 means the information at application - often it will be the same as -1 as many banks are not updating the information to Credit Bureau regularly ) – time only relative to the application
- 4 **CNT\_INSTALMENT**: Term of previous credit (can change over time)
- 5 **CNT\_INSTALMENT\_FUTUR**: Installments left to pay on the previous credit
- 6 **NAME\_CONTRACT\_STATUS**: Contract status during the month
- 7 **SK\_DPD**: DPD (days past due) during the month of previous credit
- 8 **SK\_DPD\_DEF**: DPD during the month with tolerance (debts with low loan amounts are ignored) of the previous credit

## **B.6 Installments Payments Data Table Features (installments-payments.csv)**

- 1 **SK\_ID\_PREV**: ID of previous credit in Home credit related to loan in our sample. (One loan in our sample can have 0,1,2 or more previous loans in Home Credit) – hashed
- 2 **SK\_ID\_CURR**: ID of loan in our sample – hashed
- 3 **NUM\_INSTALMENT\_VERSION**: Version of installment calendar (0 is for credit card) of previous credit. Change of installment version from month to month signifies that some parameter of payment calendar has changed
- 4 **NUM\_INSTALMENT\_NUMBER**: On which installment we observe payment
- 5 **DAYS\_INSTALMENT**: When the installment of previous credit was supposed to be paid (relative to application date of current loan) – time only relative to the application
- 6 **DAYS\_ENTRY\_PAYMENT**: When was the installments of previous credit paid actually (relative to application date of current loan) – time only relative to the application
- 7 **AMT\_INSTALMENT**: What was the prescribed installment amount of previous credit on this installment
- 8 **AMT\_PAYMENT**: What the client actually paid on previous credit on this installment

## B.7 Credit Card Balance Data Table Features (`credit_card_balance.csv`)

- 1 **SK\_ID\_PREV**: ID of previous credit in Home credit related to loan in our sample. (One loan in our sample can have 0,1,2 or more previous loans in Home Credit) – hashed
- 2 **SK\_ID\_CURR**: ID of loan in our sample – hashed
- 3 **MONTHS\_BALANCE**: Month of balance relative to application date (-1 means the freshest balance date) – time only relative to the application
- 4 **AMT\_BALANCE**: Balance during the month of previous credit
- 5 **AMT\_CREDIT\_LIMIT\_ACTUAL**: Credit card limit during the month of the previous credit
- 6 **AMT\_DRAWINGS\_ATM\_CURRENT**: Amount drawing at ATM during the month of the previous credit
- 7 **AMT\_DRAWINGS\_CURRENT**: Amount drawing during the month of the previous credit
- 8 **AMT\_DRAWINGS\_OTHER\_CURRENT**: Amount of other drawings during the month of the previous credit
- 9 **AMT\_DRAWINGS\_POS\_CURRENT**: Amount drawing or buying goods during the month of the previous credit
- 10 **AMT\_INST\_MIN\_REGULARITY**: Minimal installment for this month of the previous credit
- 11 **AMT\_PAYMENT\_CURRENT**: How much did the client pay during the month on the previous credit
- 12 **AMT\_PAYMENT\_TOTAL\_CURRENT**: How much did the client pay during the month in total on the previous credit
- 13 **AMT\_RECEIVABLE\_PRINCIPAL**: Amount receivable for principal on the previous credit
- 14 **AMT\_RECIVABLE**: Amount receivable on the previous credit
- 15 **AMT\_TOTAL\_RECEIVABLE**: Total amount receivable on the previous credit
- 16 **CNT\_DRAWINGS\_ATM\_CURRENT**: Number of drawings at ATM during this month on the previous credit



- 17 **CNT\_DRAWINGS\_CURRENT**: Number of drawings during this month on the previous credit
- 18 **CNT\_DRAWINGS\_OTHER\_CURRENT**: Number of other drawings during this month on the previous credit
- 19 **CNT\_DRAWINGS\_POS\_CURRENT**: Number of drawings for goods during this month on the previous credit
- 20 **CNT\_INSTALLMENT\_MATURE\_CUM**: Number of paid installments on the previous credit
- 21 **NAME\_CONTRACT\_STATUS**: Contract status (active signed,...) on the previous credit
- 22 **SK\_DPD**: DPD (Days past due) during the month on the previous credit
- 23 **SK\_DPD\_DEF**: DPD (Days past due) during the month with tolerance (debts with low loan amounts are ignored) of the previous credit

## References

- [Yu17] Xiaojiao Yu. “Machine learning application in online lending risk prediction”. In: *ArXiv e-prints* (July 2017). arXiv: [1707.04831](https://arxiv.org/abs/1707.04831) [q-fin.RM].
- [Gro] Home Credit Group. *Home Credit Group Website*. URL: <http://www.homecredit.net/>. (accessed: 06.01.2018).
- [Kaga] Kaggle. *Home Credit Default Risk Competition Data*. URL: <https://www.kaggle.com/c/home-credit-default-risk/data>. (accessed: 06.01.2018).
- [Kagb] Kaggle. *Home Credit Default Risk Competition Overview*. URL: <https://www.kaggle.com/c/home-credit-default-risk/>. (accessed: 06.01.2018).
- [Kagc] Kaggle. *Home Credit Default Risk Competition Public Leaderboard*. URL: <https://www.kaggle.com/c/home-credit-default-risk/leaderboard>. (accessed: 06.02.2018).
- [Wika] ML Wiki. *ML Wiki ROC Analysis*. URL: [http://mlwiki.org/index.php/ROC\\_Analysis](http://mlwiki.org/index.php/ROC_Analysis). (accessed: 06.04.2018).
- [Wikb] Wikipedia. *Receiver operating characteristic*. URL: [https://en.wikipedia.org/wiki/Receiver\\_operating\\_characteristic](https://en.wikipedia.org/wiki/Receiver_operating_characteristic). (accessed: 06.07.2018).