

Capstone Project Report

Udacity Machine Learning Nanodegree

James Dellinger

June 6, 2018

1 Definition

1.1 Project Overview

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. [Home Credit Group's](#)[Gro] goal of providing a safe lending alternative to otherwise financially down-and-out folks is perfectly aligned with this desire of mine.

Home Credit's target demographic contains people who typically have 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.

Like any mainstream for-profit financial institution, Home Credit maximizes its own returns (in the form of interest payments) when it is able to lend money to as many reliable borrowers as its balance sheet allows. The only catch is that, unlike normal banks, Home Credit needs a way to gauge whether a loan applicant is the kind of person who will eventually repay the loan, all without relying on traditional measures of financial reliability like credit score. I believe that machine learning is just the right tool to help Home Credit turn the information that it does have about its loan applicants into sound lending decisions. Indeed, a good machine learning algorithm could enable Home Credit to expand its services to as many more worthy customers.

It was with this in mind that for my capstone project I decided to participate in the [Home Credit Default Risk competition](#) on Kaggle[Kagb]. The goal of the competition, which is sponsored by Home Credit, is to create an algorithm that accurately predicts the likelihood that an applicant will repay their loan.

The competition's dataset was provided by Home Credit Group's data scientists. It contains a wide breadth 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.

The dataset's 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 122 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:

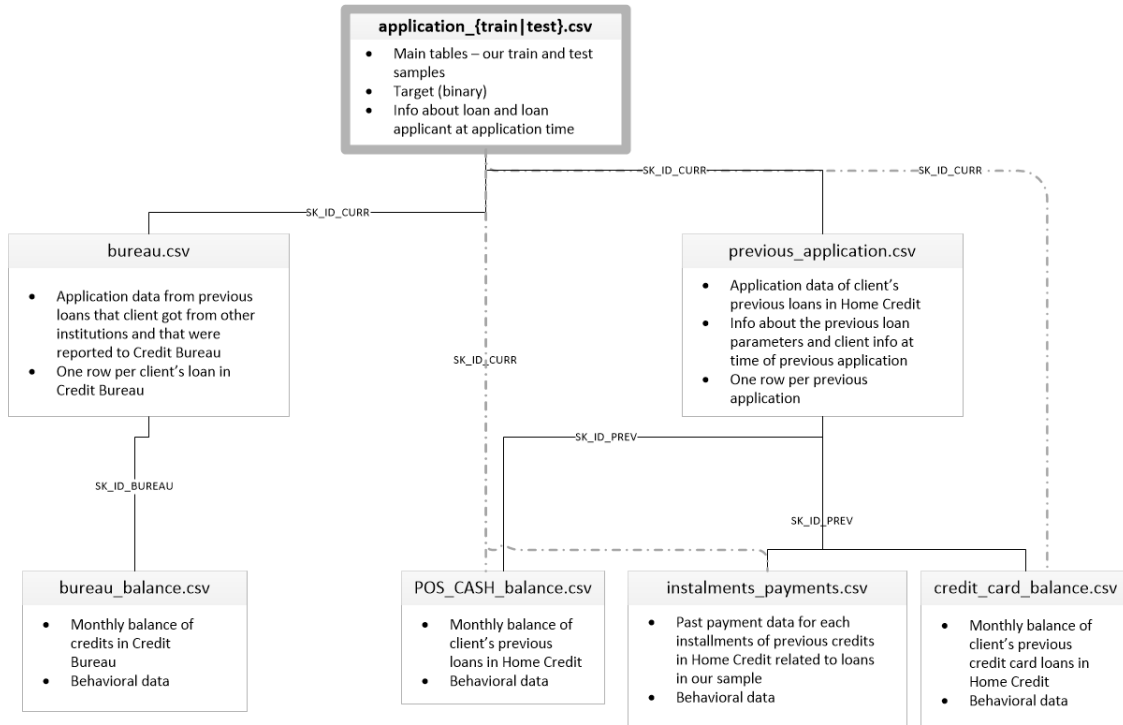


Figure 1: Home Credit Data Table Relationships[Kaga]

1.2 Problem Statement

Home Credit needs an algorithm that will take as inputs various personal and financial information originally contained in a loan applicant's profile and structured as features in

the seven data tables outlined above. The algorithm will then determine a probability of the applicant eventually 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 repay the loan and 0.0 indicates that there is zero chance that the applicant will eventually repay. The algorithm will be tested on a set of 48,744 individuals who previously borrowed from Home Credit. Competition participants must submit a CSV file 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 repaying their loan. The file must be formatted as follows:

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

To solve this competition I intend to try out several machine learning algorithms that can also return the probabilities of their predictions, such as logistical regression and multi-layer perceptron classifiers. For each classifier I build, I will experiment with tuning hyperparameters both manually and automatically. I will train my algorithms on the 307,511 borrower records that comprise the training segment. The seven data tables provided by Home Credit contain an awful lot of information that could take months to thoroughly investigate. So for the purposes of this project, I will focus on the features in the main data table. However, I will engineer at least one feature that is based on some portion of the data contained inside the other six tables. Finally, I will experiment with both manual and automatic feature selection.

Once I've finished my experimentation, I will choose classifier/hyperparameter/feature combination that scores highest on the test set and submit its test set probability predictions on the Kaggle competition's page. I am looking forward to seeing how my submission performs relative to other entries on the competition's leaderboard.

Home Credit knows which borrowers in the test set ultimately paid off their loans, and which ones eventually defaulted. The winning algorithm will therefore need to predict the highest probabilities of repayment for the largest fraction of borrowers who eventually repaid their loans, and will need to predict the lowest probabilities of repayment for the largest fraction of borrowers who ultimately defaulted on their loans. More formally, a winning algorithm will have the highest true positive rate, or sensitivity, (correctly guessing who will pay back the loans), while at the same time having a minimum false positive rate (as seldomly as possible making an incorrect guess that an individual would pay back their loan).

1.3 Metrics

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, or separating out, a particular target segment from 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 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).

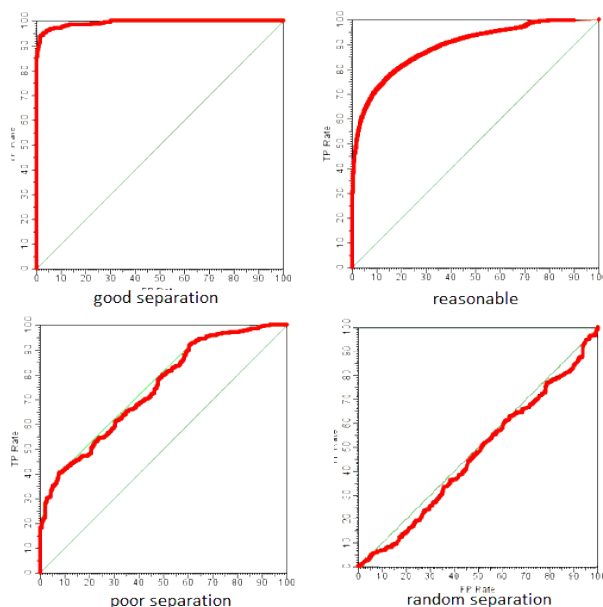


Figure 2: ROC Curve Examples[Wika]

Because the goal of this competition is to create an algorithm that correctly identifies the target segment of loan applicants who will pay back their loans, the area under the ROC curve is a perfectly appropriate measure of my algorithm's performance. 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. A competition-winning solution will be the algorithm that has an area under its ROC curve that is closest to 1. This algorithm will naturally be able to correctly predict more loan repayers than the other submissions (maximizing the true positive rate), while at the same time making the fewest incorrect predictions about which applicants are loan repayers (minimizing the false positive rate).

2 Analysis

2.1 Data Exploration

2.2 Exploratory Visualization

2.3 Algorithms and Techniques

2.4 Benchmark

Since Kaggle ranks submissions to this competition by the area under the ROC curve between their predicted probabilities of repayment and the actual observed targets (whether or not the borrower eventually repaid the loan), an ideal benchmark would be 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 the specifics of this model and its ROC curve area to be proprietary information, it is understandable that this score has not been disclosed. As a next best alternative, I will use as a benchmark the average ROC curve area of all entries on the competition's public leaderboard[Kaggle] as of June 2, 2018. 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. The reader may refer to the file `home-credit-default-risk-publicleaderboard(6-02-2018).csv`, located inside this project's Github repository[Del], to verify my calculation of this benchmark.

3 Methodology

3.1 Data Preprocessing

3.2 Implementation

3.3 Refinement

4 Results

4.1 Model Evaluation and Validation

4.2 Justification

5 Conclusion

5.1 Free-Form Visualization

5.2 Reflection

5.3 Improvement

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_BUREAU_ID**: 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_BUREAU_ID**: 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. **HOUR_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 Feature (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_INSTALLMENT_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_INSTALLMENT_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_INSTALLMENT_NUMBER**: On which installment we observe payment
5. **DAYS_INSTALLMENT**: 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_INSTALLMENT**: 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_INSTALMENT_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

- [Del] James Dellinger. *Machine Learning Nanodegree Capstone Project Repository*. URL: https://github.com/jamesdellinger/machine_learning_nanodegree_capstone_project. (accessed: 06.07.2018).
- [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. (accessed: 06.04.2018).
- [Wikb] Wikipedia. *Receiver operating characteristic*. URL: https://en.wikipedia.org/wiki/Receiver_operating_characteristic. (accessed: 06.07.2018).