

# Capstone Project

Machine Learning Engineer Nanodegree

Andrii Stelmashenko

February 13<sup>th</sup>, 2018

## Definition

### Project Overview

Loans are very popular in US. People take loans for education, property, car, etc. The definition of loan is:

A loan is money, property or other material goods that is given to another party in exchange for future repayment of the loan value amount along with interest or other finance charges [1]. There is a risk that a lender may not receive given credit back, it's called 'Credit Risk' [2]. Lenders, of course, want to minimize credit risks, for that they calculate 'Credit Score' [3]. A credit score is a number representing the creditworthiness of a person, the likelihood that person will pay his or her debts. Credit score is composed from several components related to a person, one of major components is 'Credit History'. A credit history is a record of a borrower's responsible repayment of debts.

Many people struggle to get loans due to insufficient or non-existent credit histories. And, unfortunately, this population is often taken advantage of by untrustworthy lenders. It's possible to use alternative data, including telco and transactional information, to predict their clients' repayment abilities, which can help people to get loans from trustworthy lenders.

### Problem Statement

There is a company Home Credit which works with the category of unbanked population. They have a set of data: client's previous credits provided by other financial institutions that were reported to Credit Bureau; Monthly balances of previous credits in Credit Bureau; monthly snapshots of credit card balances, cash loans applicant has/had with Home Credit; previous applications for Home Credit. Also they have a historical target score for existing customers, they want to predict for new applicants. The goal of this project is to predict probability that a person will pay it's debts based on related to the person financial information, some kind of alternative to 'Credit Score'.

Dataset is provided by competition organizer - HomeCredit company, and it is hosted on kaggle platform. There are seven different sources of data:

- application\_train/test: the main training and testing data with information about each loan application at Home Credit. Every loan has its own row and is identified by the feature SK\_ID\_CURR. The training application data comes with the TARGET indicating 0: the loan was repaid or 1: the loan was not repaid.
- bureau: data concerning client's previous credits from other financial institutions. Each previous credit has its own row in bureau, but one loan in the application data can have multiple previous credits.
- bureau\_balance: monthly data about the previous credits in bureau. Each row is one month of a previous credit, and a single previous credit can have multiple rows, one for each month of the credit length.
- previous\_application: previous applications for loans at Home Credit of clients who have loans in the application data. Each current loan in the application data can have multiple previous loans. Each previous application has one row and is identified by the feature SK\_ID\_PREV.
- POS\_CASH\_BALANCE: monthly data about previous point of sale or cash loans clients have had with Home Credit. Each row is one month of a previous point of sale or cash loan, and a single previous loan can have many rows.
- credit\_card\_balance: monthly data about previous credit cards clients have had with Home Credit. Each row is one month of a credit card balance, and a single credit card can have many rows.
- installments\_payment: payment history for previous loans at Home Credit. There is one row for every made payment and one row for every missed payment.

HomeCredit\_columns\_description.csv file contains descriptions for the columns in the various data files. There are 221 row in the file with human-readable description for each column.

## Metrics

Submissions are evaluated on area under the ROC curve [4] between the predicted probability and the observed target. The ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The evaluation logic is the next: submitted results are actually two probability distributions of two classes - 0: the loan was repaid or 1: the loan was not repaid. ROC curve will show how we did that split, it visualizes all possible classification thresholds. AUC represents the probability that a classifier will rank a randomly chosen positive observation higher than a randomly chosen negative observation.

The result will help HomeCredit business to decide if they want e.g. to minimize False Positive Rate or maximize True Positive Rate.

## Anasysis

### Data Exploration

Dataset is provided by competition organizer.

HomeCredit\_columns\_description.csv file contains descriptions for the columns in the various data files. There are 221 row in the file with human-readable description for each column. Relations between files is provided as a diagram:

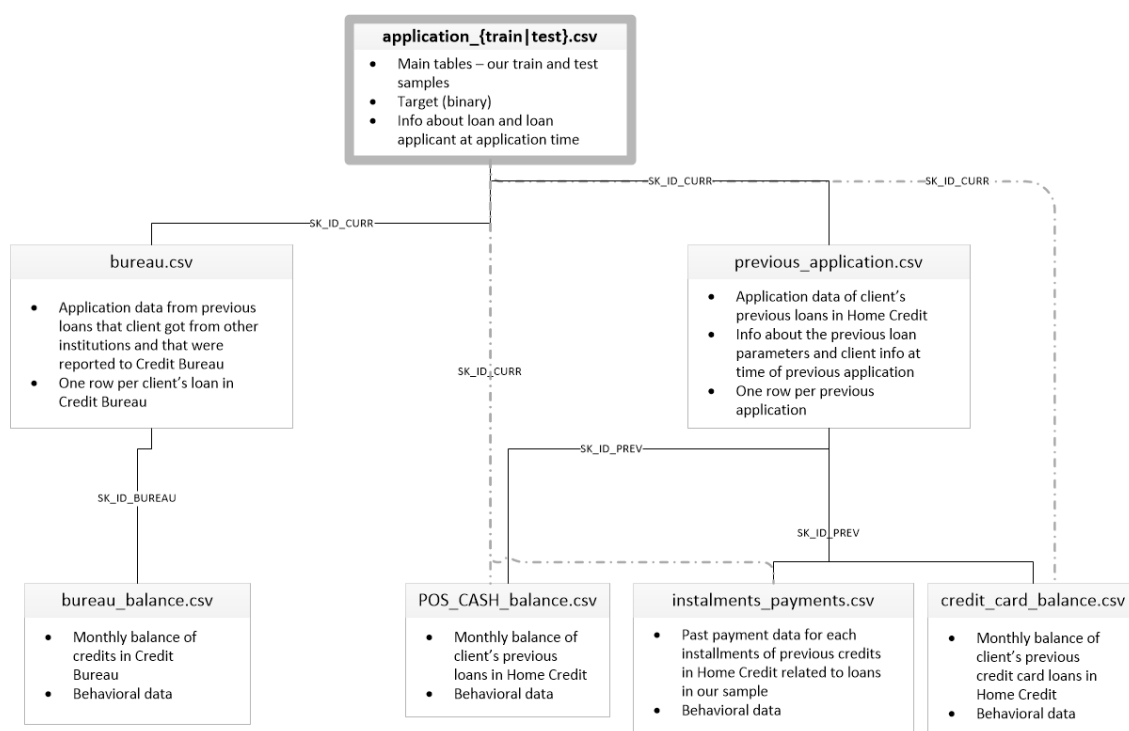


Figure 1 – Input data relations

Data files are split for convenience, later all these files should be concatenated using techniques like SQL Left Join[5] into single csv file with 221 column. Later in text I provide analysis only of some of columns which are discovered to be interesting, general analysis will be done only on one column as an example, other similar columns will be skipped to reduce document size.

All columns are of one of three data types: int, float, object. These types correspond to logical data type of each column: continuous, numerical discrete and categorical. Categorical columns are represented by object and int types. Let's count unique values:

```
app_train.select_dtypes('object').apply(pd.Series.nunique, axis = 0)
```

Partial output:

NAME_CONTRACT_TYPE	2
NAME_EDUCATION_TYPE	5
WEEKDAY_APPR_PROCESS_START	7
ORGANIZATION_TYPE	58

Some columns have only 2 unique values and some have more 58. Most algorithms work only with numerical values, so it is necessary to encode string values to numerical. Two most used techniques for that are Label Encoding and One-Hot-Encoding.

Columns of continuous type should be analyzed on anomalies. anomaly is the deviation in a quantity from its expected value[6]. Since it's often unclear what is 'expected value' there number of statistical methods how to detect anomaly or outlier, I decided to use Tukey's method[7]. It is based on interquartile range (IQR), the range is the next:

$$[Q1 - k * (Q3 - Q1), Q3 + k * (Q3 - Q1)]$$

Where  $k=1.5$  indicates an outlier,  $k=3$  means value is "far out".

Let's take AMT\_INCOME\_TOTAL column and apply IQR, the column has 3014 values which are far out from the rest. Executing outliers search for all continuous columns gives 3381 outlier when  $k=3$  and 25350 when  $k=1.5$ . AMT\_INCOME\_TOTAL contains most outliers. The column holds income of a client. These data points are very important because if client's income is big then most likely it will take big loan, if we look at such clients:

```
app_train.loc[income_total_outliers]['TARGET'].astype(int).value_counts()
0      2852
1       162
```

There is portion of rich clients which do not return loans back.

It depends on algorithm chosen if removing outliers is necessary, some algorithms are sensitive to outliers and some are not. In case of this projects outliers may be very important and it is better to choose algorithms which are not sensitive to outliers.

Missing values is another problem in real world projects. There are a lot of columns where missing values percentage goes up to 60:

Column Name	NaN %	NaN Count
COMMONAREA_AVG	69.872297	214865
NONLIVINGAPARTMENTS_MODE	69.432963	213514
NONLIVINGAPARTMENTS_MEDI	69.432963	213514
...		

Table 1 – Missing values

Having missing values in a dataset can cause errors with some machine learning algorithms. There are number of options what to do with missing values: drop rows or even whole columns where there are too many missing values, impute missing values, e.g. use mean or zeros.

Research on numerical columns using Five Number summary [9] showed that DAYS\_EMPLOYED column has group of abnormal value 365243 which is how many days before the application the person started current employment. This value means a person worked 1000 years before application, it seems like this number is sort of flag information is unknown or is not filled by some reason.

## Exploratory Visualization

Target column is binary, it takes either 0 or 1 values. Let's look at histogram plot.

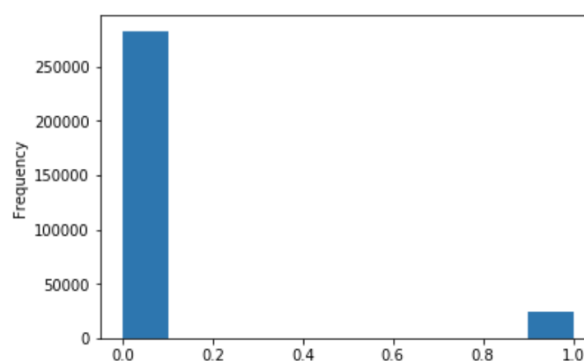


Figure 2 – Target distribution

From the fig. 2 it is seen that there is imbalanced class problem[8].

Missing values, was mentioned already above, can be visualized using data-dense matrix:

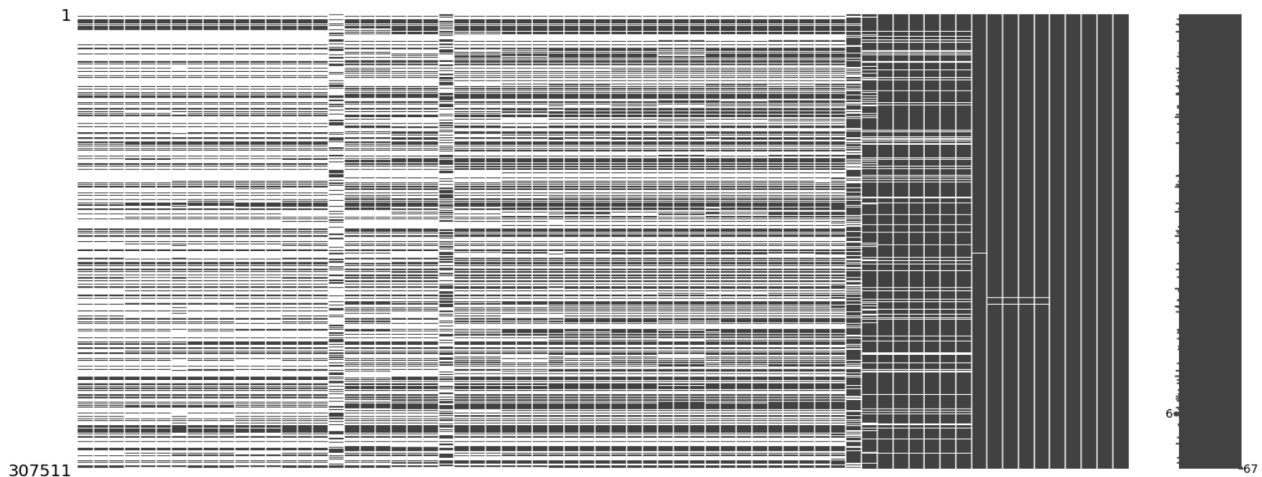


Figure 3 – Missing values matrix

It shows that most rows have missing values in a lot of columns at the same time (see horizontal blank lines).

Columns are sorted from the biggest portion of missing values to the smallest:

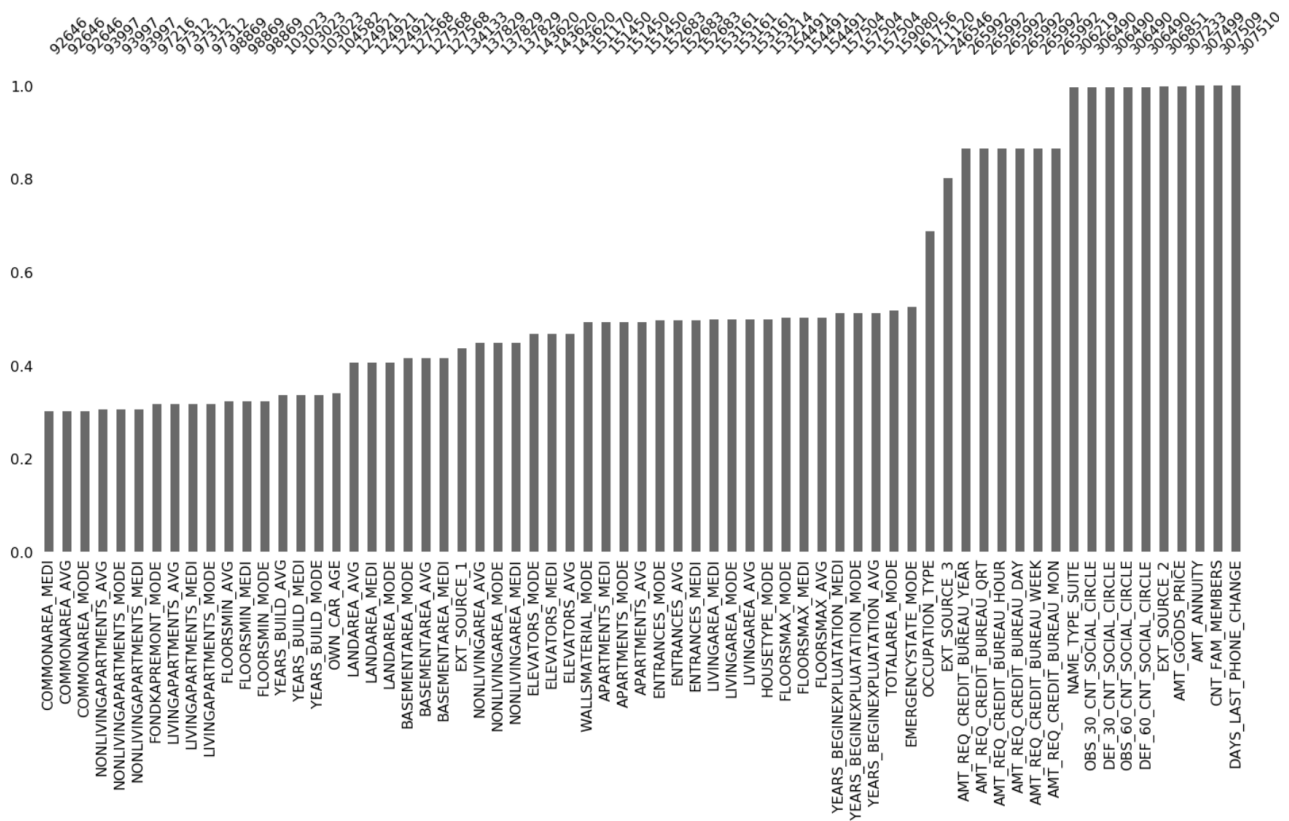


Figure 4 – Missing values barplot

The bar plot shows which columns have the biggest portion of missing values, which are mostly related to apartment/house. Certain algorithms require missing values to be filled in and some others can work with missing values. The decision what to do with missing values will be made after algorithm has been chosen.

There is certain amount of correlations exist in the data set (see Figure 5 below). On the below heatmap are shown only columns with correlations. There are columns describing apartment/house aspects and they are represented as \_AVG, \_MODE and \_MEDI, which are average, modus and median. Correlation value is close or equal to 1 which is very strong correlation type and may considered as duplicates. These columns are candidate for removal e.g. using principal components analysis technique or even, just dropping the columns \_MEDI and \_MODE leaving only \_AVG ones.

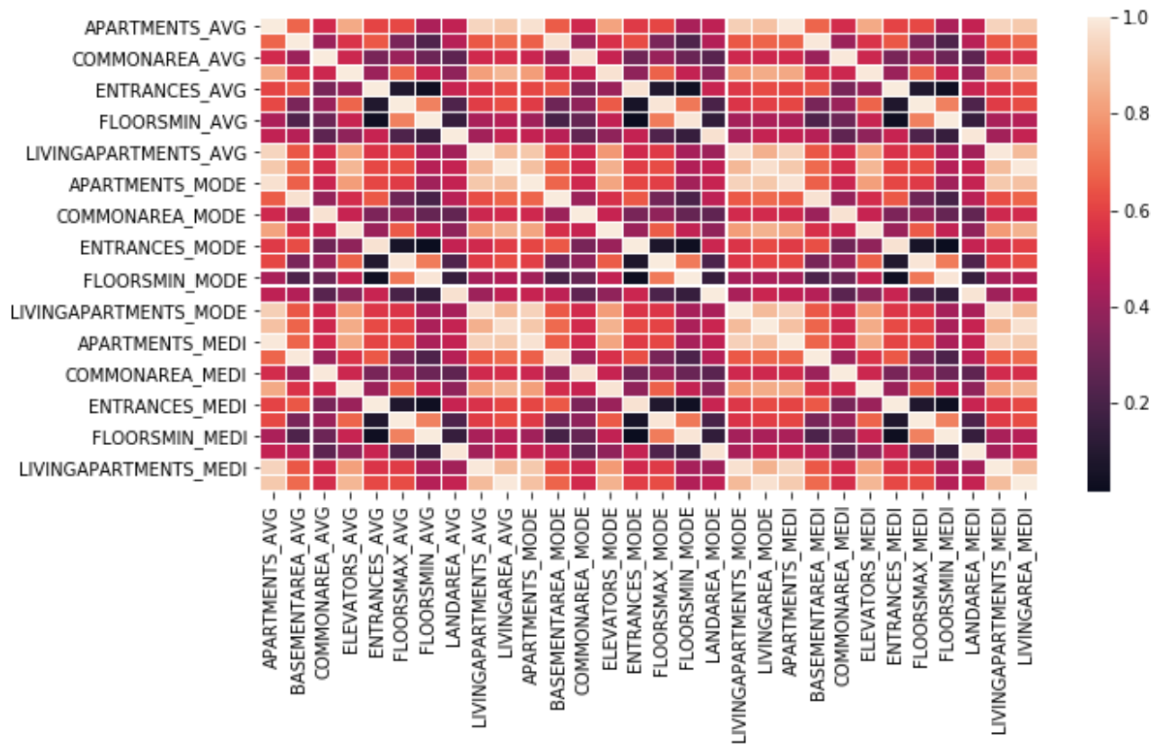


Figure 5 – Correlations

One more candidate for visualization is DAYS\_EMPLOYED column, it's worth to look at hist plot with and without outliers:

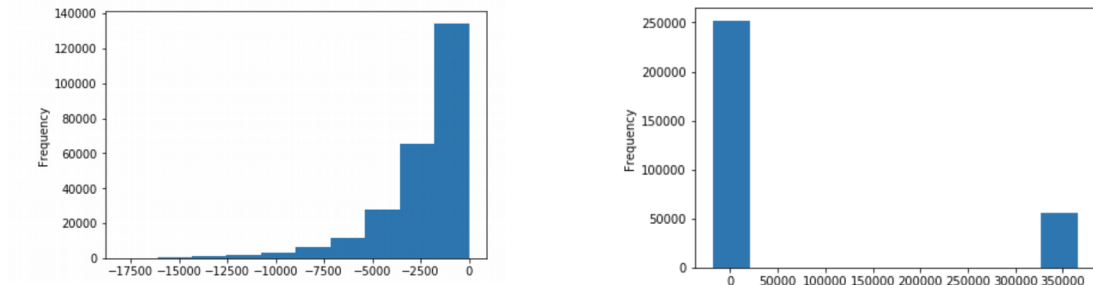


Figure 6 – Days employed column with and with out abnormalities

With out abnormal points distribution is more “gaussian”, many algorithms make assumption that data distribution of numerical columns is “gaussian”, which may affect model performance.

It's also useful to visualize outliers to see how far from other values they are, below is boxplot of OWN\_CAR\_AGE column:

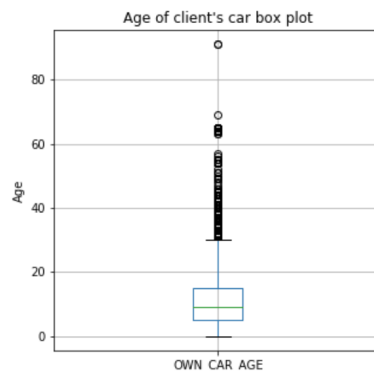


Figure 7 – Outliers visualization

From the plot it is possible to evaluate how many outliers are there outside of IQR.

## Algorithms and Techniques

Two types of classifiers will be used: Gradient Boosting Machine (LightGBM [10]) and Neural Networks. There are few reasons why these two among others.

LightGBM: gradient boosting machines proved to perform good on high dimensional datasets with mixed categorical and numerical features, it can natively handle categorical features and sparse datasets (see [10] section 4). It is very fast, it uses techniques Gradient-based One-Side Sampling and Exclusive Feature Bundling to deal with large number of data instances and large number of features respectively. It supports parallel execution out of the box [11]. Tree based models are interpretable [12], [13] which is strong side.

Continuous values handling is one of the weaknesses of tree based models, LGBM uses discretization technique which leads to information loss, it is not critical in most cases.

Neural Networks (NN) is very flexible method and allows to build accurate classifier using model of class function suitable for the problem. NN are prone to overfitting and require a large dataset (which is the case of chosen dataset). Another complexity in using NN is categorical features which gives sparse data, which leads to bad performance of NN due to the fact zeros make back propagation not working, in other words NN not learning, it does not change coefficients. There are solutions how to improve this, some of them are use Embeddings to reduce vector of categorical features [14] another one use Hashing [15] instead of One-Hot-Encoding.

At the end, if both models show good performance, it's possible to ensemble them and see if it makes performance better. For heterogeneous models stacking [16] is efficient approach.

## Benchmark

## Resources

1. Investopedia – Loan (<https://www.investopedia.com/terms/l/loan.asp>)
2. Wikipedia - Credit Risk ([https://en.wikipedia.org/wiki/Credit\\_risk](https://en.wikipedia.org/wiki/Credit_risk))
3. Wikipedia - Credit Score ([https://en.wikipedia.org/wiki/Credit\\_score\\_in\\_the\\_United\\_States](https://en.wikipedia.org/wiki/Credit_score_in_the_United_States))
4. Wikipedia - ROC curve ([https://en.wikipedia.org/wiki/Receiver\\_operating\\_characteristic](https://en.wikipedia.org/wiki/Receiver_operating_characteristic))
5. Wikipedia - LeftJoin ([https://en.wikipedia.org/wiki/Join\\_\(SQL\)](https://en.wikipedia.org/wiki/Join_(SQL)))
6. Wikipedia – Anomaly ([https://en.wikipedia.org/wiki/Anomaly\\_\(natural\\_sciences\)](https://en.wikipedia.org/wiki/Anomaly_(natural_sciences)))
7. Wikipedia – Tukey’s fences ([https://en.wikipedia.org/wiki/Outlier#Tukey's\\_fences](https://en.wikipedia.org/wiki/Outlier#Tukey's_fences))
8. Archiv.org – Imbalanced Class Problem (<https://arxiv.org/pdf/1305.1707.pdf>)
9. Wikipedia – Five Number summary ([https://en.wikipedia.org/wiki/Five-number\\_summary](https://en.wikipedia.org/wiki/Five-number_summary))
10. LightGBM (<https://papers.nips.cc/paper/6907-lightgbm-a-highly-efficient-gradient-boosting-decision-tree.pdf>)
11. LightGBM (<https://github.com/Microsoft/LightGBM/blob/master/docs/Features.rst#data-parallel-in-lightgbm>)
12. GBDT interpretability (<https://towardsdatascience.com/interpretable-machine-learning-with-xgboost-9ec80d148d27>)
13. LightGBM interpreter tool (<https://github.com/slundberg/shap>)
14. Embeddings (<https://developers.google.com/machine-learning/crash-course/embeddings/categorical-input-data>)
15. Hashing (<https://alex.smola.org/papers/2009/Weinbergeretal09.pdf>)
16. Wikipedia – Stacking Ensemble ([https://en.wikipedia.org/wiki/Ensemble\\_learning#Stacking](https://en.wikipedia.org/wiki/Ensemble_learning#Stacking))