

BUILDING CREDIT SCORING PROCESS IN VIETNAMESE COMMERCIAL BANKS USING MACHINE LEARNING

Ngày nhận bài: 07/11/2019

Ngày chấp nhận đăng: 17/01/2020

Dang Huong Giang, Nguyen Thi Phuong Dung

ABSTRACT

The industrial revolution 4.0 has affected the banking sector with the trend of transforming traditional banks into digital ones. Since the global financial crisis, risk management in banks has gained more prominence, and there has been a constant focus on how risks are being detected, measured, reported and managed. Recently, the world has seen a huge amount of data gathered within financial institutions (FIs). Crediting activities of banks must change for adapting with this trending. Current credit scoring system for individual clients of commercial banks mainly input the data and customer's information to provide the customer's credit score which helps the banks to make lending decision. Although the current system has accuracy, but it is considered as a rigid, inflexible method and still contains risks in measurement. Is there any method to increase the accuracy and inflexibility in this credit scoring system? How do we avoid missing good customers or prevent customers that are not reliable? Recently around the world, machine learning is widely considered in the financial services sector as a potential solution for delivering the analytical capability that FIs desire. Machine learning can impact every aspect of the FI's business model—improving client preferences, risk management, fraud detection, monitoring and client support automation. Therefore, this article aims to study the roles of machine learning and the application of machine learning in credit scoring systems of individual clients, building credit scoring process using machine learning in Vietnamese commercial banks.

Keywords: Machine Learning, credit scoring, Vietnamese commercial banks, Big Data, artificial intelligence.

1. Introduction

For financial institutions and the economy at large, the role of credit scoring in lending decisions cannot be overemphasized. An accurate and well-performing credit scorecard allows lenders to control their risk exposure through the selective allocation of credit based on the statistical analysis of historical customer data. The development in technology helps the commercial banks speeding up modernization process, changing banking services and activities from traditional aspect to electronic banking environment. Besides, data analytics and data management in banking sector will get more benefits from using Big Data. Collecting and analyzing big data will provide new knowledge, help making informed business decisions properly and faster, reduce operational cost, especially big data analytics

assist in statistical forecasting on banking operational activities. The modern technology makes it easier for commercial banks in collecting and developing database system, providing benefits in statistics, analytics and forecasting especially in lending and customer credit ratings.

The achievements of technology revolution 4.0 (or Industry 4.0) that impact finance and banking sector can be divided into two periods. The first period of this revolution (2008-2015) begins with cloud computing, open-source software system,

Dang Huong Giang, Department of Financial and Banking, University of Economics – Technology for Industries, Hanoi, Vietnam

Nguyen Thi Phuong Dung, School of Economics and Management, Hanoi University of Science and Technology, Hanoi, Vietnam

smart phones... The second period of this revolution are supposed to be from 2016 to 2020. At the moment, we are in the middle of second period with development of artificial intelligence, block-chain, data science, face recognition technology and biometric...With the development of the revolution, it is required that commercial banks must have strategies to actively seize opportunities, to improve their strengths in banking operations.

Artificial Intelligence (AI) can be used in credit rating system, which based on forecast models to predict and determine probability of repaying the loan of customer: whether they can pay on time, late or default. One of the most important benefits of credit rating is to help the banks to make better-informed decisions, to accept or reject providing loans/credits to customers, to increase or reduce the loan value, interest and loan's term.

With current credit score software system for individual clients of commercial banks, it is only set up to input data and customer's information into the system and the returned result is customer's credit score which help loan officers to make lending decision. However, this is a rigid, inflexible evaluation way. Although the current system has accuracy, it still has errors in measurement. What happens if a customer applies a loan at a bank and his/her loan application is rejected because of low credit score meanwhile it is approved by other banks and become a customer with good credit score. Conversely, a customer with good credit score is qualified for a loan but that loan is becoming a bad credit loan for the bank. These are 2 situations showing the failure of the current credit score software system at commercial banks. Is there any way to increase the accuracy in evaluating customers? How do we avoid missing potential customers or prevent customers that are not as good as they are

showing? The Machine Learning system with Big Data base has been considered as a solution for solving this problem. Therefore, the objective of this article is to study the roles of machine learning and the application of machine learning in credit scoring systems of individual clients in Vietnamese commercial banks.

2. Literature review

2.1. Overview of credit rating at commercial banks

Most of the banking profit comes from lending activities and providing credit/loans. Lending is a traditional banking activity that generates most of banking revenue and profit. Credit granting is an important part of banks' activities, as it may yield big profits. However, there is also a significant risk involved in making decisions in this area and the mistakes may be very costly for financial institutions (Zakrzewska, 2007). The main risk of lending for banks is the possibility of loss due to borrowers do not have abilities to repay the loan. In addition, the decision of whether or not to grant credit to customers who apply to them usually depends mainly on skills, knowledges as well as on experiences of loan officers (Thomas, 2000).

Credit rating system is an important tool to increase the objectivity, quality and efficiency of lending activity. Credit scoring model is a statistical analysis way performed by banks and financial institutions to evaluate a person's creditworthiness. It is a method that quantify risk levels based on credit scoring system. Factors used to evaluate a person's credit in credit-scoring models are different for each type of customers. Modern definition of credit scoring focuses on some main principles, including analyzing credit worthiness based on payment history, age, number of accounts, and credit card utilization, the borrower's

willingness to pay debt. Different types of loans may involve different credit factors specific to the loan characteristics; analyzing long-term risk that factor the influence of economic and business cycle as well as a tendency of ability to pay in the future; analyzing risk comprehensively based on credit scoring system.

It is necessary to use qualitative analysis to support quantitative analysis in credit scoring models. Quantitative analysis means to measure by quantity. When we do quantitative analysis, we are exploring facts, measures, numbers, and percentages, working with numbers, statistics, formula, and data. On the other hand, qualitative analysis allows you to interpret the information in non-mathematical ways. Analytic criteria may be changed to match with changes in technology and in accordance with risk management requirements. Collecting data used in credit scoring models need to be conducted objectively and flexibility. Using many different sources of information all at once to have a comprehensive analysis on financial situation of borrowers.

Scoring credit level for individual borrower of commercial banks is an internal method used by commercial banks to evaluate a customer's ability to pay off debt, risk level of loan, and based on that information, commercial banks will make decisions whether to approve or deny credit; manage risk; create appropriate policy for each type of borrower based on credit scoring results. Besides, credit scoring system is used also for classifying and supervising credit system. Classifying and supervising credit is applied for all customers and is conducted periodically; as well as when there are signs of inabilities to pay obligations.

One of traditional methods used to evaluate and approve credits or loans is relied

on some of rating criteria; however, some of them are very difficult to measure or evaluate correctly. For example, "5C's of credit", namely Character, Capacity, Capital, Collateral, and Conditions – a common method was used to consider when evaluating a consumer loan request (Abrahams & Zhang, 2008). Some of the criteria such as "Character" and "Capacity", that look at the ability of the borrower to repay the loan through income, are hard to evaluate. Moreover, credit scoring method based on "5C's of credit" standard has high cost. The breadth and depth of experiences are varied by loan officer, therefore, that led the potential for bias in individual decisions resulting inconsistent loan decisions. Due to these limitations, banks and financial institutions need to use credit scoring methods and assessment methods that are reliable, objective and low cost in order to help them decide whether or not to grant a credit for loan application (Akhavain, Frame, & White, 2005; Chye, Chin, & Peng, 2004). Moreover, according to Thomas and et al. (2002), banks need a credit scoring method that meets the following requirements: (1) cheap and easy to operate, (2) fast and stable, (3) make consistent decisions based on unbiased information which is independent from subjective feelings and emotions, and (4) the effectiveness of the credit scoring system can be easily checked and adjusted at any time to regulate promptly with changes in policies or conditions of the economy.

For credit classification and scoring, the traditional approach is purely based on statistical methods such as multiple regression (Meyer & Pifer, 1970), discriminant analysis (Altman, 1968, Banasik, Crook, & Thomas, 2003), and logistic regression (Desai, Crook, & Overstreet, 1996; Dimitras, Zanakakis, & Zopounidis, 1996; Elliott & Filinkov, 2008;

Lee, Chiu, Lu, & Chen, 2002). However, under requirements of the Basel Committee on Banking Supervision, banks and financial institutions are required to use credit scoring models which are more reliable in order to improve the efficiency of capital allocation. In order to meet these requirements, in recent years, there have been some new models of credit classification based on machine learning and artificial intelligence (AI) approaches. Unlike previous approaches, these new methods do not provide any strict assumptions in comparison to the traditional statistical approaches. Instead, these new approaches attempt to exploit and provide the knowledge, the output information based only on inputs that are observations and past information. For the credit classification problems, some machine learning models such as Artificial Neural Network (ANN), Support Vector Machines (SVMs), K Nearest Neighbors (KNN), Random Forest (RF), Decision Tree (DT), has proved to be superior in terms of accuracy as well as reliability compared to some traditional classification models (Chi et al., 2004, Huang et al., & Wang, 2007; Ince & Aktan, 2009; Martens et al., 2010).

2.2. Machine Learning and Machine Learning algorithm

Machine learning is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and

improve from experiences without being explicitly programmed. Machine learning focuses on the development of computer programs that can access data and use it learn for themselves.

Machine learning uses training, i.e., a learning and refinement process, to modify a model of the world. The objective of training is to optimize an algorithm's performance on a specific task so that the machine gains a new capability. Typically, large amounts of data are involved. The process of making use of this new capability is called inference. The trained machine-learning algorithm predicts properties of previously unseen data.

There are many different types of machine learning algorithms, with hundreds published each day, and they're typically grouped by either learning style (i.e. supervised learning, unsupervised learning, semi-supervised learning) or by similarity in form or function (i.e. classification, regression, decision tree, clustering, deep learning, etc.). Regardless of learning style or function, all combinations of machine learning algorithms consist of the following:

- Representation (a set of classifiers or the language that a computer understands)
- Evaluation (aka objective/scoring function)
- Optimization (search method; often the highest-scoring classifier, for example; there are both off-the-shelf and custom optimization methods used).

Table 1:

The three components of machine learning algorithms

Representation	Evaluation	Optimization
- Instances	Accuracy/Error rate	- Combinatorial optimization
K-nearest neighbor	Precision and recall	Greedy search
Support vector machine	Squared error	Beam search
- Hyperplanes	Likelihood	Branch-and-bound
Naive Bayes	Posterior probability	- Continuous optimization
		Unconstrained (Gradient descent,

Logistic regression	Information gain	Conjugate gradient, Quasi-Newton methods)
- Decision trees	K-L divergence	
- Sets of rules	Cost/Utility	Constrained (Linear programming, Quadratic programming)
Propositional rules	Margin	
Logic programs		
- Neural networks		
- Graphical models		
Bayesian networks		
Conditional random fields		

Machine learning emphasize on goals such as: (1) Teaching machine and computer to learn basic human skills such as listening, watching and understanding language, problem solving skill, programming... and (2) Assisting human beings in solving and finding solutions from a huge amount of information or big data that we have to face every day.

According to experimental researches: Machine learning algorithms along with data mining algorithms that based on new techniques and computation methods operate better for forecasting purpose. Machine learning algorithms are designed to learn from historical data to complete a task, or to make accurate predictions, or to behave intelligently.

Some of basic concepts in Machine Learning used in credit scoring:

Observation: symbol is x , which is input in algorithm. Observation can be a data point, row or sample in a data set. Observation usually represents as a vector $x = (x_1, x_2, \dots, x_n)$ which can be called as **feature vector** where each x_i is a **feature**. Feature vector is a list of features describing an observation with multiple attributes. (In Excel we call this a row). For example, we want to predict if a borrower can create a bad debt in the future or not based on calculation of a function in which Observation include features like biological sex, age, income, credit history...ect.

Label: symbol is y , output of calculation. Each observation will have an appropriate label to go with. In previous example, Label can be “overdue” or “on time”. Label can be described under many categories but they all can be converted into a number or a vector.

Model: are a function $f(x)$, A function assigns exactly one output to each input of a specified type. Input an observation x and return a label $y=f(x)$.

Parameter: Machine learning models are parameterized so that their behavior can be tuned for a given problem. These models can have many parameters and finding the best combination of parameters can be treated as a search problem. A model parameter is a configuration variable that is internal to the model and whose value can be estimated from the given data. For example, in a model of second degree polynomial function: $f(x)=ax_1+bx_2+c$, its parameters are set of (a,b,c) . However, there is a special parameter called hyperparameter

Parameter: all the model’s factors which are used for calculating the output. For example, the model is a quadratic polynomial function: $f(x) = ax_1 + bx_2 + c$, its parameter is a triad (a, b, c) .

Currently, there are many available machine learning algorithms, so the question is “Which algorithm is the best?”. There isn’t any clear answer for this question since the accuracy of each algorithm depends on input data and the structure of specific input data.

A general method to find a suitable model for a set of specific data is applying widely used and certified model.

Credit risk is still one of biggest challenges in banking system. Until now, commercial banks have not completely optimized forecast abilities of digital risk. A report from McKinsey shows that machine learning will be able to reduce credit deficit by 10%, with more than half of credit managers expect that time to process credit applications will reduce by 25% to 50%.

2.3. Experiences in applying Machine Learning in individual credit scoring at several commercial banks around the world

With machine learning, commercial banks and financial institutions have been able to apply sciences into their operations instead of prediction. A large number of commercial banks and financial institutions have been using AI to detect and prevent fraudulent transactions for several years around the world.

In 2017, JP Morgan Chase introduced COiN, a contract intelligence platform that using machine learning can review 12,000 annual commercial credit agreements in seconds. It would take staff around 360,000 hours per year to analyse the same amount of data.

AI-based scoring models combine customers' credit history and the power of big data, using a wider range of sources to improve credit decisions and often yielding better insights than a human analyst. Banks can analyze larger volumes of data – both financial and non-financial – by continuously running different combinations of variables and learning from that data to predict variable interactions.

In Germany, a recent Proof of Concept (PoC) model showing that running AI-based scoring models on Intel® Xeon® processors and using Intel® Performance Libraries can

help banks boost machine-learning and data analytics performance. Using Intel-optimized performance libraries in the Intel® Xeon® Gold 6128 processor helped machine-learning applications to make predictions faster when running a German credit data set of over 1,000 credit loan applicants

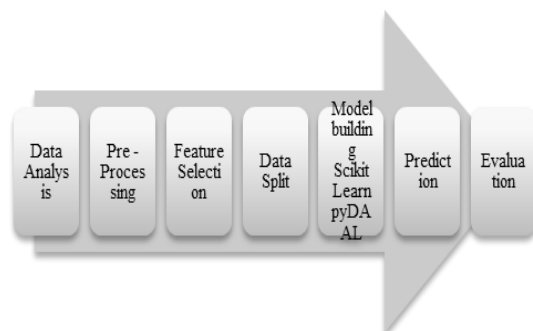


Figure 1: Proof of Concept (PoC) model

Source: Intel.com

Dataset analysis: This is the initial exploration of the data, including numerical and categorical variable analysis.

Pre-processing: Data pre-processing transforms the data before feeding it to the algorithm. In this case, it will involve converting the categorical variables to numerical variables using various techniques.

Feature Selection: In this step, the goal is to remove the irrelevant features which may cause an increase in run time, generate complex patterns, etc. This can be done either by using Random Forest or Xgboost algorithm.

Data split: The data is then split into train and test sets for further analysis.

Model Building: Machine-learning models are selected for training.

Prediction: During this stage, the trained model predicts the output for a given input based on its learning.

Evaluation: In order to measure performance, various evaluation metrics are available such as accuracy, precision, and recall.

3. Machine learning applications for vietnamese commercial banks

3.1. Credit scoring model for individual customers at Vietnam

In 2007, a research about "Credit Scoring for Vietnam's Retail Banking Market" by Dinh, T.T.H and Kleimeier, S. with credit scoring model for individual customers used at Vietnam's commercial banks includes a set of 22 variables such as age, income, education, occupation, time with employer, residential status, gender, marital status, loan type...ect. This model is used to determine the level of influence of these variables on credit risk and from the results collected to Table 2:

Variables included in the Vietnamese retail credit scoring model

Panel A: Variables considered in the first round of credit assessment

<u>Variable</u>	<u>Categories</u>
age	18-25, 26-40, 41-60, >60 (years)
education	postgraduate, graduate, high school, less than high school
occupation	professional, secretary, businessman, pensioner
total time in employment	<0.5, 0.5-1, 1-5, >5 (years)
time in current job	<0.5, 0.5-1, 1-5, >5 (years)
residential status	Owns home, rents, lives with parents, other
number of dependents	0, 1-3, 3-5, >5 (people)
applicant's annual income	<12, 12-36, 36-120, >120 (million VND)
family's annual income	<24, 24-72, 72-240, >240 (million VND)

Panel B: Variables considered in the second round of credit assessment

<u>Variable</u>	<u>Categories</u>
performance history with bank (short-term)	new customer, never delaid, payment delay less than 30 days, payment delay more than 30 days
performance history with bank (long-term)	new customer, never delaid, delay during 2 recent years, delay earlier than 2 recent years
total outstanding loan value	<100, 100-500, 500-1000, >1000 (million VND)
other services used	
average balance in saving account during previous year	savings account, credit card, savings account and credit card, none
	<20, 20-100, 100-500, >500 (million VND)

create an individual credit scoring model applied for Vietnam's retail banks.

Individual credit scoring model consist of 2 components which are borrower's personal characteristics score as well as ability to repay the debt; the borrower's banking relationship score (as shown in Table 1). Based on total scores, banks and financial institutions classify risk levels into 10 different classes from Aaa to D. In order to apply this model, it is required that commercial banks have to create score system for each variable that is suitable with its current status and its individual customer database system.

Panel C: Loan decision

<u>Applicant's scoring</u>	<u>Score</u>	<u>Loan decision</u>
	≥ 400	Lend as much as requested by borrower
Aaa	351-400	Lend as much as requested by borrower
Aa	301-350	Lend as much as requested by borrower
a	251-300	Loan amount depends on the type of collateral
Bbb	201-250	Loan amount depend on the type of collateral with
Bb	151-200	assessment
b	101-150	Loan application requires further assessment
Ccc	51-100	Reject loan application
Cc	0-50	Reject loan application
c	0	Reject loan application
d		Reject loan application

Source: Dinh, T.T.H and Kleimeier, S (2007)

Table 3:

The credit scoring model's variables and estimated coefficients

(Note that the variables are selected based on the stepwise method. In this table the included variables are ranked by absolute value of the coefficients.)

included variables	estimated coefficient	standard error	significance level
time with bank	-1.774	0.121	0.0%
gender	-1.557	0.222	1.0%
number of loans	-0.938	0.051	1.4%
loan duration	-0.845	0.080	3.7%
deposit account	-0.750	0.104	3.1%
region	-0.652	0.030	13.6%
residential status	-0.551	0.278	44.6%
current account	-0.492	0.208	10.4%
collateral value	-0.402	0.096	9.8%
number of dependants	-0.356	0.096	9.9%
	-0.285	0.054	2.5%
time at present address	-0.233	0.101	68.1%
	-0.190	0.057	53.0%
marital status	-0.181	0.047	3.4%
collateral type	-0.156	0.067	60.3%
home phone	-0.125	0.054	3.3%
education			
loan purpose	-3.176	0.058	4.6%
constant			

In addition, commercial banks also use Fico model to rate credit score for retail customers.

The most widely adopted credit scores are FICO Scores created by Fair Isaac Corporation. 90% of top lenders use FICO Scores to help them make billions of credit-related decisions every year. FICO Scores are calculated solely based on information in consumer credit reports maintained at the credit reporting agencies.

By comparing this information to the patterns in hundreds of thousands of past credit reports, FICO Scores estimate your level of future credit risk.

Base FICO Scores have a 300-850 score range. The higher the score, the lower the risk. But no score says whether a specific individual will be a "good" or "bad" customer.

While many lenders use FICO Scores to help them make lending decisions, each lender has its own strategy, including the level of risk it finds acceptable for a given credit product.

Table 4:

FICO's 5 credit score components

Proportion	Components
35%	Payment history: The first thing any lender wants to know is whether you've paid past credit accounts on time. This helps a lender figure out the amount of risk it will take on when extending credit. This is one of the most important factors in a FICO® Score. Be sure to keep your accounts in good standing to build a healthy history.
30%	Amount owed: Having credit accounts and owing money on them does not necessarily mean you are a high-risk borrower with a low FICO® Score. However, if you are using a lot of your available credit, this may indicate that you are overextended-and banks can interpret this to mean that you are at a higher risk of defaulting.
15%	Length of credit history: In general, a longer credit history will increase your FICO® Scores. However, even people who haven't been using credit long may have high FICO Scores, depending on how the rest of their credit report looks. Your FICO® Scores take into account: <ul style="list-style-type: none"> - How long your credit accounts have been established, including the age of your oldest account, the age of your newest account and an average age of all your accounts. - How long specific credit accounts have been established. - How long it has been since you used certain accounts.
10%	Credit mix: FICO® Scores will consider your mix of credit cards, retail accounts, installment loans, finance company accounts and mortgage loans. Don't worry, it's not necessary to have one of each.
10%	New credit: Research shows that opening several credit accounts in a short period of time represents a greater risk-especially for people who don't have a long credit history. If you can avoid it, try not to open too many accounts too rapidly.

Source: Fico.com.

FICO credit scoring model is used when banks have ability to review and check customers' credit history easily. Credit data is recorded and updated from credit institutions. According to FICO's credit score model, borrowers who have scores at and above 700 are considered "good", individuals who have credit score less than 620 are considered risky borrowers and banks will be afraid to grant loans for them

3.2. Applying machine learning in individual credit scoring at Vietnamese commercial banks

Over the years, some modeling techniques to implement credit ratings have developed, including parametric or non-parametric, statistics or Machine Learning, Supervised or unsupervised algorithms, Artificial Neural. Recent techniques include very sophisticated approaches, using hundreds of different models, different models of testing methods, combining a variety of algorithms to achieve high accuracy. However, the most outstanding model building technique called Credit Scorecard is widely applied by many banks in the world. (ex. Commonwealth Bank of Australia, Standard Chartered Bank...) is Standard Scorecard, it's based on (Logistic Regression Model).

Credit card model is simple, easy to understand, easy to deploy and run fast. Combining statistics and Machine Learning, the accuracy of this method is equivalent to sophisticated techniques. Its output score can be directly applied to assess the probability of bad debt, thereby providing inputs to the valuation of bad debt based on the risk. This is very important for lenders who need to comply with the Basel II.

The credit card model can be described as: attributes input from customers, customer characteristics (For example, age, income, occupation, etc.), their past credit

information (For example, information collected from the National Credit Information Center – CIC, with other credit information that the bank has..., based on model calculations, each attribute will be assigned a certain coefficient, Their sum is equal to the output score. Based on the output score, it can be identified bad debt probability (PD – Probability of Default). This probability makes it easy to calculate the value of credit risk, so, the bank quickly determine minimum amount of capital for credit risk in accordance with Basel II standards. This is the reason that Credit Scoring Engine is based on this model researched and applied by Hyperlogy for the customers.

Therefore, credit scoring system and customer ratings by a scorecard, created by Machine Learning technology, Logistic regression model application, not only assessment ability to perform financial obligations of a customer to a bank such as pay interest and repay the loan principal when due, but it is also a tool of bank support the bank in controlling compliance with Basel II. From the theory of the credit card model, the authors propose Machine Learning application process to Credit rating at Vietnamese commercial banks is as follows:

3.2.1. Choosing machine learning algorithms

Traditional models usually focus on the strengths of the borrower's finances and abilities to repay the loan. They classify borrowers based on their credit history, quality of collateral, payment history and other considerations. That makes it easier for banks when it comes to clarify the relationships between consumer's behavior and credit score.

However, the way which consumers spend their money on saving and lending are

changing, as well as the technology. Many financial institutions are using credit scoring model to reduce risks in credit scoring and in granting, credits. Credit scoring models based on traditional statistical theory have been used widely at present. However, these traditional models cannot be used when there is a lot of input data. Since big data have an influence on the accuracy of model-based forecast. Machine learning can be used in credit scoring in order to reach higher accuracy level from analyzing a large amount of big data.

A typical business procedure in providing loan services is to receive loan application, to determine credit risk, to make decision on granting a loan and to supervise the repayment of interest and principal. During mentioned above process, many things can happen, such as: how we can accelerate credit analysis and underwriting process; how we can supervise repaying process and how we can timely intervene when there is a chance of default. To solve both problems, we can create a two-stage of credit scoring model.

Establishing process:

All applicants for a loan need to be checked. The model can be used to analyze and learn from historical application data, thereby determine whether a new applicant is credible enough to grant a loan or not and whether specific criteria of the applicant are provided, such as income, marital status, age, credit history (whether or not had bad debt in the past) etc.

Supervising process:

The system will check database of borrowers who have been approved a loan. By using the repayment historical data and the status of customers who completed the entire loan process, we can train another model to make a forecast of whether this

customer have a high probability of insolvency. By observing the applicant's repayment profile for the first few payback periods and changing the characteristics, this model will help to make new adjustments based on updated information. This process is more efficient in processing time as well as it is more accurate than the traditional way.

Machine learning algorithms such as: Decision tree (DT), Support vector machine (SVM), Genetic algorithm (GA), artificial neural network (ANN) have many advantages for statistical models and optimize credit risk assessment techniques. Machine learning algorithms are evaluated as having an advantage over other statistical methods in assessing corporate credit risk, especially for nonlinear regression model. In machine learning, the hybrid approach method is a prospective field of study to improve classification or predict performance, to rate credit score. Combined method provides better performance than single evaluation methods.

Several commercial banks in the world have used a "hybrid" model that uses machine learning algorithms such as support vector machine, support vector machine, artificial neural network and decision tree to calculate credit scores in order to support credit decision making. Using the proposed hybrid model and experimental results show that this model has higher accuracy in classification when compared to other credit scoring methods.

There are 3 classes of evaluation criteria are used to compare testing methods: Accuracy Class 1, Class 2, Overall accuracy are calculated by the following formula:

Accuracy Class 1 = Number of classification and number of "bad" observations / Total number of bad observations

Accuracy type 2 = Number of classification and number of "good" observations / Total number of good observations

Overall accuracy = Number of correct classification / Total number of observations

When there is a need of rating customers' credit, the hybrid model using SVM technique has higher accuracy and better performance than other techniques in credit scoring model. This model combines both classification method and clustering method, some of machine learning techniques such as SVM, Decision Trees, Artificial Neural Networks are used for classification, Fuzzy c-means (FCM) clustering method is used for clustering. Therefore, with the proposed hybrid model along with SVM techniques using the method of aggregating the level of the members gives us best accurate results and the best performance when calculating credit scores in order to limit credit risks in banking operations.

3.2.2. Proposing the process of applying Machine Learning in credit scoring of individual customers at Vietnamese commercial banks

Developing a Machine Learning model requires many steps and it will have some similarities/steps with developing a common technology software. Frameworks for Approaching the Machine Learning Process (Regina Esi Turkson, 2016):

- (1) Data collection and data preparation
- (2) Choose a Model
- (3) Train the Model
- (4) Evaluate the Model
- (5) Parameter Tuning
- (6) Make Predictions

Machine Learning is a modern branch of software development and computer science. As in conventional software, pre-packaged

solutions are less costly but are not suitable with the specific needs. While on the other hand, building Machine Learning systems can lead to solutions that are more flexible with universal variability. Based on research on current individual customer credit scoring system and the goal of applying Artificial Intelligence with Machine Learning technology and Big Data platform (nguồn trích dẫn), the authors propose the process of applying Artificial Intelligence in credit scoring of individual customers for commercial banks including the following steps:

- a. Determining the goals of the system
 1. Assist in credit evaluation and credit approval and in credit risk management
 2. Set up and develop a customer policy
 3. Overcome weaknesses of previous credit scoring system.
 4. Meet the requirements of State Bank.
- b. Building the foundation for customer database (Big Data)

Commercial banks need to standardize individual customer credit history data at their banking system with customer information and scoring criteria. The bigger data, the better. Suggesting the minimum period to keep and store digital customer records is 10 years. The customer information fields are organized in scientific way with below information:

1. Regarding personal identity;
2. Regarding transaction history with bank;
3. Regarding loan information;
4. Regarding collateral assets;
5. Regarding the customer's debt repayment history.

....

c. Building an individual customer credit scoring model with foundation of Big Data and Machine Learning technology

Building an individual customer credit scoring model with Big Data and Machine Learning is implemented through following steps:

1. Choose an appropriate credit scoring model;
2. Program the Machine Learning system based on the set goals, selected credit scoring model and existing Big Data system.
- d. Operating system experimentally, evaluating and using system officially.

Using a new model with Machine Learning technology to evaluate and forecast credit scoring results are based on previous customer database. Comparing the prediction results of Machine Learning with the actual results that happened with this customer group and comparing with the forecast results of the current individual customer credit scoring system. Through analyzing old customers' data with existing customer database and comparing accuracy level with old technologies to evaluate new credit scoring models based on Machine Learning technology, assess the strengths and weaknesses of the new model to make adjustments when needed Applying on a new customer database and then evaluate the accuracy level based on actual results.

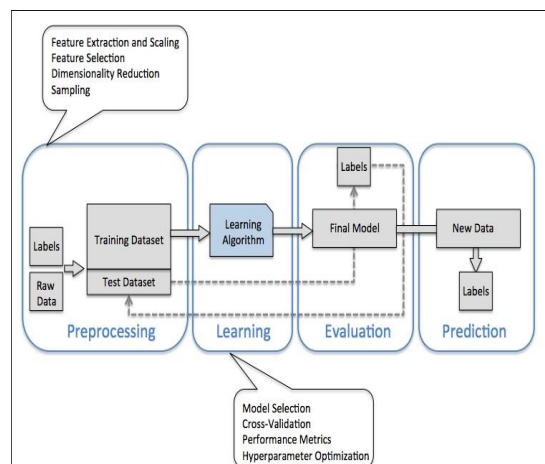


Figure 3. A roadmap for building machine learning systems in personal credit scorings

4. Conclusions

In the two scoring processes mentioned above, both look the same, but they have different models. The repayment supervising process looks like the loan granting process, but it is learned and drawn from various historical data, particularly the data from former customers who have completed repayment which include records from payment history and customer's characteristics.

Today, the most used machine learning algorithms can be classified as single or full classification. Representatives of the single classification algorithms are Classification and Regression Tree (CART), Naïve Bayes, Support Vector Machines (SVM), Logistic Regression. The modification of single classification with multiple learning models to solve the same problem, is widely used, such as Random Forests, CART-Adaboost, etc.

Basically, machine learning is like teaching financial knowledge to a new student. He will learn from historical data to determine the quality of a loan based on

several indicators, after that he will have experiences to make his own decisions whether or not to grant a credit in the following cases.

In banking and insurance industry, many commercial banks develop its own applications based on Machine Learning, including Credit Scoring, Risk Analysis, Fraud Detection, Cross-Selling.

Machine Learning can still make mistakes. Algorithms are created by human; therefore they are still affected by human, and like all other areas of data analysis, there will be times when the data collected is good or usable, but sometimes data collected is not good and should be ignored.

Machine Learning also has some limits on transparency, especially when it involves some "black boxes" that are an essential part of the Neural Network

However, Machine Learning is a great tool that plays an increasingly important role in the evolution of technology, helping artificial intelligence (AI) reach many more users.

REFERENCES

- Altman I. E. (1968), "Financial ratios, discriminant analysis and the prediction of corporate bankruptcy," *The journal of finance*, vol. 23.
- Arezzo, M.F.; Guagnano, G. (2018), "Response-Based Sampling for Binary Choice Models with Sample Selection", *Econometrics*, 6, 12.
- Baesens, B., Van Gestel, T., Viaene, S., M. Stepanova, J. Suykens, and Vanthienen. J (2003), "Benchmarking state-of-art classification algorithm for credit scoring," *Journal of operational research society*, 627-635.
- Cho, S., H. Hong and Ha, B. (2010), "A hybrid approach based on the combination of variable selection next term using previous term decision trees and case-based reasoning next term using the Mahalanobis distance: For bankruptcy prediction," *Expert systems with applications*, Elsevier journals, 3482-3488.
- Dinh, T.T.H and Kleimeier, S. (2007). Credit Scoring for Vietnam's Retail Banking Market, *Int. Rev. Financ*, 16, 471-495

- Gestel, T.V, B. Baesens, J. A. Suykens, D. Van den Poel, D.-E. Baestaens, and Willekens, B. (2006), “Bayesian kernel-based classification for financial distress detection,” *European journal of operational research*.
- Hsieh, N.C. (2005), “Hybrid mining approach in design of credit scoring model,” *Expert systems with applications*.
- Intel (2018), “Machine Learning-Based Advanced Analytics Using Intel® Technology”, Reference Architecture.
- Louzada, F.; Ara, A.; Fernandes, G.B. (2016), “Classification methods applied to credit scoring: Systematic review and overall comparison”, *Comput. Oper. Res.*, 21, 117–134.
- McKinsey (2017), “Smartening up with Artificial Intelligence (AI) - What’s in it for Germany and its Industrial Sector?”, McKinsey report.
- Munkhdalai, L., Namsrai, O.E., Lee, Y.J. and Ryu, K.H (2019), “An Empirical Comparison of Machine-Learning Methods on Bank Client Credit Assessments”, *Sustainability journal*, 11, 699.
- Oreski, S.; Oreski, G. (2014), “Genetic algorithm-based heuristic for feature selection in credit risk assessment” *Expert Syst. Appl*, 41, 2052–2064.
- Regina Esi Turkson; Edward Yeallakuor Baagyere ; Gideon Evans Wenya (2016), “A machine learning approach for predicting bank credit worthiness”, *Third International Conference on Artificial Intelligence and Pattern Recognition (AIPR)*.
- Shi, J.; Xu, B. (2016), “Credit scoring by fuzzy support vector machines with a novel membership function” *J. Risk Financ. Manag*, 9, 13.
- Schebesch, B. and R. Stecking, (2005), “Support vector machine for classifying and describing credit applicants: Detecting typical and critical regions,” *Journal of the operational research society*.
- Tsai, C.F. and M. L. Chen,” Credit rating by hybrid Machine Learning techniques,” *Applied soft computing* 2010
- Wang, G., J. Hao, J. Ma, and H. Jiang,” A comparative assessment of ensemble learning for credit scoring,” *Expert systems with applications*, 2011.
- Xia, Y.; Liu, C.; Li, Y.; Liu, N. A (2017), “Boosted decision tree approach using Bayesian hyper-parameter optimization for credit scoring”, *Expert Syst. Appl*, 78, 225–241..