

Online learning behavior analysis based on machine learning

Online learning
behavior
analysis

Ning Yan

Shanghai Open University, Shanghai, China, and

Oliver Tat-Sheung Au

Open University of Hong Kong, Kowloon, Hong Kong

97

Received 15 August 2019
Revised 10 September 2019
13 September 2019
Accepted 16 September 2019

Abstract

Purpose – The purpose of this paper is to make a correlation analysis between students' online learning behavior features and course grade, and to attempt to build some effective prediction model based on limited data.

Design/methodology/approach – The prediction label in this paper is the course grade of students, and the eigenvalues available are student age, student gender, connection time, hits count and days of access. The machine learning model used in this paper is the classical three-layer feedforward neural networks, and the scaled conjugate gradient algorithm is adopted. Pearson correlation analysis method is used to find the relationships between course grade and the student eigenvalues.

Findings – Days of access has the highest correlation with course grade, followed by hits count, and connection time is less relevant to students' course grade. Student age and gender have the lowest correlation with course grade. Binary classification models have much higher prediction accuracy than multi-class classification models. Data normalization and data discretization can effectively improve the prediction accuracy of machine learning models, such as ANN model in this paper.

Originality/value – This paper may help teachers to find some clue to identify students with learning difficulties in advance and give timely help through the online learning behavior data. It shows that acceptable prediction models based on machine learning can be built using a small and limited data set. However, introducing external data into machine learning models to improve its prediction accuracy is still a valuable and hard issue.

Keywords Machine learning, Online learning, Learning behaviour analysis

Paper type Research paper

1. Introduction

In a traditional learning environment, teachers can understand the learning effect of students through the examination of each course. It is difficult to obtain the specific learning process of each student during the teaching process, and it is also difficult to predict the students' learning performance from the students features and learning behavior data. Therefore, it is also difficult for teachers to master the situation of students and provide timely intervention and help.

With the rapid development of online learning, especially the rise of MOOC in recent years, a large amount of data has gradually accumulated in various online learning platforms. The online behavior data are the data generated by the learners' interaction with the operating platform during the learning process, and is mainly recorded and stored in real time by the learning platform database and other tools. The classification of a large amount of learning behavior data has certain significance for the collection of data.

© Ning Yan and Oliver Tat-Sheung Au. Published in *Asian Association of Open Universities Journal*. Published by Emerald Publishing Limited. This article is published under the Creative Commons Attribution (CC BY 4.0) licence. Anyone may reproduce, distribute, translate and create derivative works of this article (for both commercial and non-commercial purposes), subject to full attribution to the original publication and authors. The full terms of this licence may be seen at <http://creativecommons.org/licenses/by/4.0/legalcode>

This work is supported by the Croucher Chinese Visiting Scholar Project (2018–2019). The authors would like to thank the Croucher Foundation and the Open University of Hong Kong for providing the valuable research opportunity.



The researcher determines the relevant indicators according to the objectives of the analysis, so as to determine the type of learning behavior from which the data are collected based on the results of the behavior classification, making the data collection process easier and more convenient, and making the results of the behavior description more accurate. According to the theory of behavioral science, learners conduct learning activities through online learning platforms based on their own needs and generate different learning behaviors. The login behavior is generated based on the requirements of the course plan, the contents of the guide, the contents of the evaluation and the academic achievements, etc. The resource access behavior is generated based on the demand for the resources. The retrieve behavior is generated based on the requirements of the query learning content, the subject tool, the course content, and the like. The forum interaction behavior is generated based on the need to communicate and collaborate with others. The feedback and reflection behavior are generated based on reflecting learners themselves, evaluation curriculum and peers' needs. And the learning task-related behavior is generated based on the requirements of completing the learning task.

According to Burgos *et al.* (2018), with the increase in the amount of data on the online learning platform, researchers are looking for ways to make those data understandable and meaningful. In order to analyze and dig out more potential educational laws, researchers delve into the theory, framework, tools and practices of learning analysis. In recent years, people have been studying more and more research on learning behavior analysis, especially on predicting students' learning performance. Valiente *et al.* (2015) point out that with the continuous development of learning analysis research, researchers have begun to study learning predictions using machine learning methods and tools.

However, in recent years, with the widespread use of the internet, especially the rise of the mobile internet, teachers and students, are more inclined to use mobile smart terminals and social applications such as WeChat, WhatsApp, Twitter and Facebook to communicate with each other, discuss learning problems, even submit homework on tablets or mobile phones, rather than traditional tools on LMS or CMS. Because the above data are difficult to collect, researchers have much lower effective data available from online learning platform than before. Therefore, how to analyze online learning behavior and build an effective learning performance prediction platform is a valuable issue under the premise of limited data.

2. The overview of online learning behavior analysis

With the rise of educational data and the rise of big data-related technology research, online learning behavior analysis has attracted more and more researchers' attention. Compared with the traditional offline learning analysis, because online learning behavior analysis can obtain various recorded data of learners' online learning, instead of obtaining subjectively strong data through questionnaires, it is more objective.

Learning is a series of activities carried out by learners for their own needs. In the learning process, the learner's expectations of the learning objectives are achieved through learning tools during a certain learning time, while interacting with learning content, learning environment and learning partners. Finally, the learning achievement is formed. Therefore, a complete learning process should include: learners, learning partners, learning content, learning tools, learning objectives, learning achievement, learning environment and learning time. Whether it is offline learning or online learning, the learning process should include the above elements.

The learning behavior is a series of actions that learners produce during the learning process, including reading books, answering questions, watching videos, viewing courseware, browsing forums, uploading resources, accessing learning platforms, discussing and communicating with others, and so on. Combined with behavioral science theory, learning behavior can be considered as a process in which learners interact based on

a two-way interaction between a learning goal and a learning environment. Think of each learning behavior as a system, then it should include subjects of behavior, objects of behavior, operations of behavior, environments of behavior and outcomes of behavior. Each behavior activity is a systematic process in which the actor of behavior interacts with the object of behavior in a certain learning environment and time, performs operations and produces certain results.

The occurrence of each learning behavior is guided by subjective thinking and is also limited by the environment. When a learner acquires knowledge from a video resource, there may be operations such as fast forward and pause; when the learner is engaged in collaborative learning, it may be necessary to express his own opinion on the problem. The operational behaviors in online learning are diversified and vary with different operating objects. Online learning behaviors can be divided into the following four categories.

2.1 Trajectory behaviors

The learner's behaviors such as login, exit and jump between web pages in the system are trajectory behaviors. Trajectory behaviors are the most extensive type of behavior and the most basic type of learning behaviors. When a learner retrieves a certain type of resource, this behavior belongs to both the resource learning behavior and the trajectory behavior. The difference between the trajectory behavior and the resource learning behavior is that the resource learning behavior contains the body of the search, and the track behavior is only to retrieve the action, not including the main content of the retrieval.

2.2 Social behaviors

In the traditional offline learning environment, learners, teachers and other learners can communicate and collaborate with each other face to face, while people in online learning environment (OLE) rely on the network, including real-time communication and discussion-based communication. When the learners collaborate to complete a certain task, they can express their own opinions separately; when the learning is in doubt, they can consult with the teacher or other learners; when they have some insights in the study, they can post in the forum or in the individual learning space, and so on. In the social behaviors of online learning, learners generally submit text-based content that reflects the learner's internal knowledge structure and learner's personal characteristics.

2.3 Resource learning behaviors

Resource learning behaviors are the subject behaviors in learning behavior, and resources are an important part of learning. The online learning platform provides a variety of learning resources, such as multimedia types, text types and the like. When the learner browses the multimedia type of learning resources, it will produce fast-forward, pause and loop-playing behaviors; when the learner reads the text-type learning resources, it will generate jump positioning, pause preview and the like.

2.4 Evaluation and reflection behaviors

Most of the evaluations of learners focus on tests, tests, assignments, etc., in order to test the learning achievement of learners at a certain stage. An online learning process often includes records of learner test scores and submission of answers, but if only the final answer or grade is recorded, the learner's process information in the answer process will be lost. Therefore, the learner's answers to the test are modified several times, and this process should be recorded to reflect the learner's actual learning situation more comprehensively.

3. Relevant studies in online learning behavior prediction

Learning behavior prediction is an important research content and application goal of learning behavior analysis. In the past few years, with the continuous development of learning analysis technology, more and more learning prediction research work has emerged. Dietzuhler and Hurn point out that predicting learners' performance by analyzing learning behaviors helps developers to evaluate online learning systems more effectively, continuously improve system availability, and expand system functions to visualize learners' behaviors and the future trends of learning behavior; help teachers to understand the trends of learning behavior; help teachers to engage in appropriate human interventions for learners at the appropriate time; help teachers to continuously improve the curriculum and improve the quality of teaching. At the same time, the system can also help teachers to give timely help to learners with poor performance at the appropriate time to improve their learning performance (Dietzuhler and Hurn, 2013).

First, from the perspective of predicting objectives (labels), they can be mainly divided into three categories: dropout rate prediction, passing rate prediction and learning performance prediction. Balakrishnan and Coetzee (2013) use four behavioral eigenvalues and hidden Markov models to predict the possibility of dropping out of school on the MOOC platform. Kloft *et al.* use click data and SVM models to predict learner dropout rates on the MOOC platform (Kloft *et al.*, 2014). Jiang *et al.* choose two logistic regression models to predict the performance of a week's learners on the MOOC platform (Jiang *et al.*, 2014). Qiu *et al.* develop a unified predictive model that can be used to predict learners' academic performance and certificate acquisition (Qiu *et al.*, 2016). Sorour and Mine (2016) use predictive trees and random forests in machine learning to predict learner performance.

Second, from the perspective of the selection of eigenvalues, many researchers carry out a variety of theoretical research, which covers a variety of eigenvalues that may be related to learning effects from different perspectives. For example, Brown summarizes three categories of predicting eigenvalues: learner characteristics, learning behavior feature and student work. For different types of eigenvalues, he discusses related predictive models and cases (Brown, 2012). Berry proposes three indicators that affect academic achievement: academic factors, demographic factors, and cultural and social factors. In actual prediction, limited to various objective factors, not all feature values related to learning effects can be collected, so there are often many manual trade-offs in the operation (Berry, 2017). For example, Balakrishnan and Coetzee (2013) use the cumulative percentage of video lectures that can be viewed, the number of posts in the forum, the number of replies in the forum and the number of course progress views as predictors. Romero *et al.* directly predict learner performance from the forum's participation. The predictive indicators used include: the number of information sent by the learner, the number of new topics created by the learner, the number of posts read by the learner, the time the learner stayed on the forum, the concentration of the learner, and the adherence of the learner (Romero *et al.*, 2013).

Third, the prediction algorithms can generally be divided into two categories: white box models and black box models. Villagr -Arnedo *et al.* (2017) point out that it is generally thought that the predictions of the white-box model are "interpretable," and the predictions of black-box learning are "unexplained." In machine learning area, the decision tree model in machine learning is a typical white box model, and the artificial neural networks (ANNs) model is a typical black box model. In the related research on learning behavior prediction, it is generally believed that black box prediction has higher accuracy, especially when dealing with complex relationships. The white box prediction has a higher degree of interpretation, that is, a specific explanation can be given for the predicted result. Of course, high interpretation and high accuracy are often contradictory. When the degree of interpretation is high, the accuracy of prediction may be reduced, and vice versa.

4. The overview of machine learning and ANN

In general, machine learning is usually considered to be a branch of Artificial Intelligence (AI). The concept of “Machine Learning” was first proposed by Arthur Samuel in 1959. Tom Mitchell gives a more formal definition: “If a computer program is measured by P for a certain type of task T based on experience E , then we call this computer program learning from experience E , for a certain type of task T , its performance is measured by P ” (Mitchell, 2003).

Machine learning tasks can be divided into many categories. It mainly includes supervised learning, semi-supervised learning and unsupervised learning. In supervised learning, a machine learning algorithm builds a mathematical model from data sets which contain input and expected output. Semi-supervised learning builds mathematical models from incomplete training data, some of which are unlabeled. In unsupervised learning, machine learning algorithms build mathematical models from data sets that contain only input values but no expected output labels. Unsupervised learning can be used to discover structures in data, such as grouping and clustering of data points. Unsupervised learning can discover patterns in data and group inputs into different classes, such as feature learning.

Supervised learning algorithms mainly include classification algorithms and regression algorithms. In simple terms, Classification is mainly used to predict discrete or nominal labels, while regression is mainly used to predict continuous or ordered labels.

The ANN is generally considered to be a branch of supervised machine learning algorithm. An ANN is a mathematical model or computational model that mimics the structure and function of a biological neural network for estimating or approximating functions. Multi-layer feedforward neural networks are generally trained using the Backpropagation (BP) algorithm. The BP algorithm performs learning by iteratively processing a set of training samples, comparing the network prediction of each sample with the known class labels. For each training sample, the weights are modified to minimize the mean square error between the network prediction and the actual class. Because the modification is done “reverse,” the algorithm is also called “Backpropagation” (BP).

There have been some research efforts to apply ANN models to student performance predictions. Zacharis points out that students’ performance can be accurately predicted in blended learning environments by training the ANN model (Zacharis, 2016). Mason *et al.* (2018) use three different ANN models to predict the retention rate of engineering students and compared the output of the three models. The ANN algorithm can also be used to predict whether a student can successfully complete a course (Daud *et al.*, 2017). ANN models can be used not only to predict students’ performance, but also to predict the performance of school teachers (Osman, 2016).

5. Course grade prediction using ANN

The student data of the ELEC S305F Autumn 2018 course of OUHK with a total of 192 student records are used in the experiment. The collected student eigenvalues include: student age, student gender, connection time, hits count, and days of access. The predicting objectives (labels) of the model are the course grade. All the data used in the experiment are acquired from the OLE and Registry database of OUHK. The course grade range is (A, A-, B+, B, B-, C+, C, F). All data were preprocessed before the experiment, and the connection time unit was converted to minutes.

The ANN used in the experiment is classical three-layer feedforward neural networks. The scaled conjugate gradient BP algorithm is adopted here. The ratio of training data set, verification data set and test data set is set as 3:1:1. The overall prediction accuracy is 30.8 percent in the first experiment.

In general, data normalization and eigenvalue discretization can effectively improve the prediction accuracy of the machine learning model, so we try to normalize and discretize the data to improve the prediction accuracy of the model.

5.1 Prediction model using data normalization

In the field of machine learning, different evaluation eigenvalues usually have different dimensions and dimensional units, which often affect the results of data analysis. In order to eliminate the dimensional impact between indicators, data standardization needs to be done to resolve the comparability between data indicators. After the original data are processed by data standardization, the indicators are in the same order of magnitude, which is suitable for comprehensive comparative evaluation. The most typical data standardization process is the normalization of data. In summary, the purpose of normalization is to make the preprocessed data limited to a certain range, such as $[0, 1]$ or $[-1, 1]$, thus eliminating the adverse effects caused by singular sample data.

The normalization methods commonly used in machine learning include linear proportional transformation, range conversion and Z-score. The range conversion method is used in this experiment:

$$y_i = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (1)$$

After the data normalization using the range conversion method, the experimental result shows that the overall prediction accuracy of ANN model is improved to 39.6 percent.

5.2 Prediction model using eigenvalue discretization

In the experiment, the three eigenvalues which are connection time, hits count and days of access are equally distanced. The discretization rule is shown in Table I.

Using the eigenvalue discretization method, the experimental result shows that the overall prediction accuracy of ANN model increases from 30.8 to 38.5 percent.

5.3 Prediction model of binary classification (good/not good)

The above section is about multi-value classification. To improve the predictive accuracy of ANN model, we may decrease the number of predictive labels and transform the multi-value classification problem into the binary classification problem. Then the course grades which is equal to or higher than B (which is: A, A-, B or B+) is transformed to 1, and the course grade which is lower than B (which is: B-, C, C+ or F) is transformed to 0. Then it is called a "Good/Not good" binary classification problem. The result of the ANN model prediction (classification) for "Good/Not Good" is shown in Table II.

The experimental result shows that although the data set and the number of eigenvalues is small, the ANN model for the course grade binary classification using data normalization (Good/Not Good) has acceptable prediction accuracy. The overall prediction accuracy of test data set is 78.4 percent. It is obvious that the binary classification model has much higher prediction accuracy than the multi-class classification model.

Connection time (min)	Hits count	Days of access	Discretization value
0-1	0-10	0-10	1
2-5	11-20	11-20	2
6-10	21-30	21-30	3
11-15	31-40	31-40	4
16-20	41-50	41-50	5
21 or above	51 or above	51 or above	6

Table I.
Eigenvalue
discretization rule

5.4 Prediction model of binary classification (Pass/Fail)

As in Section 5.3, we can propose and deal with the “Pass/Fail” binary classification problem. The result of the ANN model prediction (classification) for “Pass/Fail” is shown in Table III.

The experimental result shows that the ANN model for the course grade binary classification using data normalization (Pass/Fail) has a low prediction accuracy rate for the Label Fail. This can be attributed the small total number of data samples and the low proportion of Fail in the student’s grades. The neural network model can easily treat such data as noise processing, resulting in a lower prediction accuracy rate.

6. Correlation analysis of the student data

Students’ course grades are converted to numbers (GPA) before conducting correlation analysis, and the correspondence between GPA and course grade is shown in Table IV.

The result of the Pearson correlation analysis of student data is shown in Table V.

From the correlation analysis data of Table V, it can be found that days of access and course grade have the highest correlation, Hits count is second, and the connection time has relatively low correlation with course grade. Age and gender are the least relevant to course grade.

Data set	Actual label	Predicted label		Correct percentage
Train		0	1	
	0	7	27	20.6
	1	6	72	92.3
	Overall percentage	11.6	88.4	70.5
Verify	0	2	11	15.4
	1	3	26	89.7
	Overall percentage	11.9	88.1	66.7
Test	0	2	8	20.0
	1	0	27	100.0
	Overall percentage	5.4	94.6	78.4

Table II.
Experimental data of
the binary
classification model
using normalization I

Data set	Actual label	Predicted label		Correct percentage
Train		0	1	
	0	1	13	7.1
	1	0	95	100.0
	Overall percentage	0.9	99.1	88.1
Verify	0	0	4	0.0
	1	0	38	100.0
	Overall percentage	0.0	100.0	90.5
Test	0	0	3	0.0
	1	0	37	100.0
	Overall percentage	0.0	100.0	92.5

Table III.
Experimental data of
the binary
classification model
using normalization II

Course grade	A	A-	B+	B	B-	C+	C	F
GPA	4.0	3.7	3.3	3.0	2.7	2.3	2.0	0.0

Table IV.
Correspondence
between GPA and
course grade

	Age	Gender	Connection time	Hits count	Days of access	Course grade
<i>Age</i>						
Pearson correlation	1	0.068	0.049	0.002	-0.034	0.067
Sig.		0.348	0.500	0.973	0.643	0.359
<i>Gender</i>						
Pearson correlation	0.068	1	0.086	0.150*	0.119	0.092
Sig.	0.348		0.237	0.038	0.102	0.203
<i>Connection time</i>						
Pearson correlation	0.049	0.086	1	0.665**	0.425**	0.107
Sig.	0.500	0.237		0.000	0.000	0.141
<i>Hits count</i>						
Pearson correlation	0.002	0.150*	0.665**	1	0.896**	0.250**
Sig.	0.973	0.038	0.000		0.000	0.000
<i>Days of access</i>						
Pearson correlation	-0.034	0.119	0.425**	0.896**	1	0.294**
Sig.	0.643	0.102	0.000	0.000		0.000
<i>Course grade</i>						
Pearson correlation	0.067	0.092	0.107	0.250**	0.294**	1
Sig.	0.359	0.203	0.141	0.000	0.000	1

Note: *,**Correlation is significant at the 0.05 and 0.01 levels, respectively (two-tailed)

Table V.
Pearson correlation
analysis of
student data

The above results are also compatible with common sense, because days of access reflects the student's persistence in the course learning, and hits count reflects the student's concentration on the course learning. At first glance, it seems weird that connection time has low correlation with course grade. However, it is reasonable in fact, because while students are in the OLE, maybe they just open the learning page but are engaged in other activities, so the connection time does not accurately reflect the students' learning performance, and its correlation with the course grade is relatively low. From Table V, it is also found that student age and gender have the lowest correlation with course grade. Days of access and hits count have the highest correlation coefficient with each other, so it shows that the persistence and concentration of a student is highly correlated.

7. Conclusions and outlook

In summary, the following inferences can be obtained from this paper.

First, days of access has the highest correlation with course grade, followed by hits count, and connection time is much less relevant to students' learning performance, i.e. course grade. Student age and gender have the lowest correlation with course grade. Because days of access reflects the student's learning persistence, and hits count reflects the students' concentration on the course. Connection time is difficult to reflect the persistence and concentration of students' learning, so the correlation between online time and student achievement is relatively low. The correlation between age or gender and course grade is very low, almost negligible. However, since the students in the experimental data set are all full-time students, the variance of student age is small, so the inference is still open to question. If the age variance is large enough, the result may be different.

Second, based on the above inferences, this paper provides some way for teachers to roughly predict students' performance by examining online learning data, and to provide timely intervention and guidance for students with learning difficulties. If the teacher finds that some students have much lower number of days of access and hits count than the

average number on the learning platform, he or she can pay appropriate attention to them and give them effective learning support to prevent them from failing the exam.

Third, Data normalization and eigenvalue discretization can effectively improve the prediction accuracy of the machine learning model, so data preprocessing is important to improve the prediction accuracy of the machine learning model. Although the experimental data set and the number of eigenvalues is small, we can still get acceptable predictive models, especially for binary classification problems. However, in the future, we will work with other institutions to obtain larger student data sets and more diverse eigenvalues from OLEs. At the same time, questionnaire data, etc. will be added to the machine learning training set. Then an ANN model with higher precision to predict students' course grades more accurately would be obtained.

However, it is obvious that the accuracy of the prediction (classification) model is not so high because the experimental data set is small and limited. In fact, even for days of access, which has the highest correlation with course grade, the Pearson correlation coefficient for both is only close to 0.3, which is just moderately correlated, let alone other eigenvalues. If we can get a bigger and better student data set, including more effective personality and learning behavior data, such as: the number of homework submitted, the data of student grades of their homework, the data about social interaction among students and teachers, the data about student posting and messaging on the forum, more complex models could be built, and the prediction accuracy of the machine learning model would be improved significantly. In fact, more useful data exist on mainstream social applications. It would be very valuable to extract the data from social applications for the training of machine learning models, but the data on social platforms are difficult for researchers to get and to use. It is envisaged that in the future more external data should be introduced for train machine learning models to improve their prediction accuracy. This practice will inevitably encounter thorny issues like how to deal with student privacy, and how to deal with those issues is a big challenge for us, so we should think comprehensively, cautiously and deeply about them.

References

- Balakrishnan, G.K. and Coetzee, D. (2013), "Predicting student retention in massive open online courses using hidden Markov models", Technical Report No. EECS-2013-109, UC Berkeley, available at: www2.eecs.berkeley.edu/Pubs/TechRpts/2013/ (accessed September 3, 2019).
- Berry, L.J. (2017), "Using learning analytics to predict academic success in online and face-to-face learning environments", Dissertation for the Degree of Doctor of Education, Boise State University, Boise, ID, March 6, available at: <https://scholarworks.boisestate.edu/cgi/viewcontent.cgi?article=2317&context=td> (accessed September 3, 2019).
- Brown, M. (2012), "Learning analytics: moving from concept to practice", EDUCAUSE Learning Initiative, July, available at: <https://library.educause.edu/-/media/files/library/2012/7/elib1203-pdf> (accessed September 3, 2019).
- Burgos, C., Campanario, M.L., Peña, D., Lara, J.A., Lizcano, D. and Martínez, M.A. (2018), "Data mining for modeling students' performance: a tutoring action plan to prevent academic dropout", *Computers & Electrical Engineering*, Vol. 66, pp. 541-556.
- Daud, A., Aljohani, N.R., Abbasi, R.A., Lytras, M.D., Abbas, F. and Alowibdi, J.S. (2017), "Predicting student performance using advanced learning analytics", *Proceedings of the 26th International Conference on World Wide Web Companion, Perth*, pp. 415-421.
- Dietzuhler, B. and Hurn, J.E. (2013), "Using learning analytics to predict (and improve) student success: a faculty perspective", *Journal of Interactive Online Learning*, Vol. 12 No. 2, pp. 17-26.
- Jiang, S., Williams, A.E., Schenke, K., Warschauer, M. and O'Dowd, D.K. (2014), "Predicting MOOC performance with week 1 behavior", *Proceedings of the 7th International Conference on Educational Data Mining, International Educational Data Mining Society, London*, pp. 273-275.

- Kloft, M., Stiehler, F., Zheng, Z. and Pinkwart, N. (2014), "Predicting MOOC dropout over weeks using machine learning methods", *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP) in Doha, Association for Computational Linguistics*, pp. 60-65.
- Mason, C., Twomey, J., Wright, D. and Whitman, L. (2018), "Predicting engineering student attrition risk using a probabilistic neural network and comparing results with a backpropagation neural network and logistic regression", *Research in Higher Education*, Vol. 59 No. 3, pp. 382-400.
- Mitchell, M. (2003), *Machine Learning* (translated by Zeng huayin *et al.*), China Machine Press, Beijing.
- Osman, A.H. (2016), "An evaluation model of teaching assistant using artificial neural network", *VAWKUM Transactions on Computer Sciences*, Vol. 11 No. 2, pp. 10-14.
- Qiu, J., Tang, J., Liu, T.X., Gong, J., Zhang, C., Zhang, Q. and Xue, Y. (2016), "Modeling and predicting learning behavior in MOOCs", *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*, ACM Press, San Francisco, CA, pp. 93-102.
- Romero, C., López, M.I., Luna, J.M. and Ventura, S. (2013), "Predicting students' final performance from participation in on-line discussion forums", *Computers & Education*, Vol. 68 No. C, pp. 458-472.
- Sorour, S.E. and Mine, T. (2016), "Building an interpretable model of predicting student performance using comment data mining", *5th IIAI International Congress on Advanced Applied Informatics (IIAI-AAI)*, Kumamoto, IEEE Press, pp. 285-291.
- Valiente, J.A.R., Merino, P.J.M., Leony, D. and Kloos, C.D. (2015), "ALAS-KA: a learning analytics extension for better understanding the learning process in the khan academy platform", *Computers in Human Behavior*, Vol. 47, pp. 139-148.
- Villagrà-Arnedo, C.J., Gallego-Durán, F.J., Llorens-Largo, F., Compañ-Rosique, P., Satorre-Cuerda, R. and Molina-Carmona, R. (2017), "Improving the expressiveness of black-box models for predicting student performance", *Computers in Human Behavior*, Vol. 72, pp. 621-631.
- Zacharis, N.Z. (2016), "Predicting student academic performance in blended learning using artificial neural networks", *International Journal of Artificial Intelligence and Applications*, Vol. 7 No. 5, pp. 17-29.

Corresponding author

Ning Yan can be contacted at: ynphd@hotmail.com