# Learning Style Prediction Using Students' E-textbook Reading Behaviors Data

**Meijun GU[a]\* Bo JIANG[b] & Chengjiu YIN**

[a]*College of Educational Science and Technology, Zhejiang University of Technology, China*
[b]*Department of Educational Information and Technology (Shanghai Engineering Research Center of Digital Education Equipment), East China Normal University, China*
[c]*Information Science and Technology Center, Kobe University, Kobe, Japan*
*\* bjiang@deit.ecnu.edu.cn*

**Abstract:** Adaptivity is one of the most prominent features of intelligent textbooks in the 21st century. Learning style is a personality characteristic of learners, which is used to describe learners' preference for processing information in a certain way. Learning style was often measured by questionnaires, which were easily influenced by learners' subjective cognition and external interference. This study proposes a data-driven approach to automatically detect learning style of learners. In the learning environment of e-textbook, 234 students' reading data was collected, and a learner model is constructed using machine learning technology. The results show that the proposed model achieves a promising performance in prediction learning style. This will help measure learning style more accurately and provide support for personalization. The learner model applied to e-textbook can promptly and dynamically monitor the changes of students' learning behavior in the online environment, and adaptively intervene, remedy or enhance.

**Keywords:** intelligent textbook, adaptivity, machine learning, learning style

## 1. Introduction

E-textbooks are believed to play an important role in future learning, but most of existing e-textbooks have not considered readers' personalities. Recently, there are increasing interests from educational technology community in designing intelligent e-textbook. Intelligent e-textbook is essentially an adaptive learning system that can provide learners personalized learning service. Usually, an adaptive learning system includes four parts: content model, learner model, instructional model and adaptive engine (Ritter et al., 2020).

Learning style (LS) is the personality characteristic of learners, which can be a part of leaner model. Learners with different LS have different learning preferences. That is to say, in the face of a certain theme or contents presented by e-textbook, students tend to select individual reading resources and adopt specific reading strategies (Gomede, Barros, & Mendes, 2020). Some studies have found that LS can guide the e-textbook to improve students' learning process (Truong, 2016). Therefore, the design of e-textbooks should consider the differences of learners' LS and optimize the learning process.

Many researchers classify LS according to different standards, hoping to recommend suitable reading resources and find effective reading strategies for learners. Compared with other LS models, Felder-Silverman learning styles model (FSLSM) has a more detailed classification of LS. FSLSM summarizes learners' learning preferences from four dimensions of information processing, perceiving, inputting and understanding (Felder, & Silverman, 1988). The first dimension of FSLSM is divided learners into active and reflective according to whether he prefers cooperative inquiry or independent reasoning. The second dimension is divided into sensing and intuitive learners according to whether he prefers touching things to learn or observing things to learn. The third dimension is divided into visual and verbal learners according to whether he prefers to see charts, tables, and figures or words and texts. The fourth dimension can be divided into sequential and global learners according to whether he prefers to acquire information step by step or overall grasp.

As we all know, the traditional method of LS measurement is based on questionnaires or scales. It has two major problems. On the one hand, it takes a lot of time to fill in the questionnaire and process

data (Aissaoui, Madani, Oughdir, & Allioui, 2019). On the other hand, this method is easy to be interfered by external factors, for example, students' comprehensive deviation will affect results of the measurement (Bernard, Chang, Popescu, & Graf, 2017). To solve these issues, this study proposes a data-driven approach for measuring students' LS. The aim of this study is to propose a easy-to-use method that detect e-textbook readers' LS automatically and accurately.


## 2. Related Work

### 2.1 Adaptive Intelligent Textbook

Intelligent textbooks are considered as family members of adaptive systems. The adaptive system usually includes four parts: content model, learner model, instructional model and adaptive engine. Among them, the content model covers various knowledge components and prerequisite relationships. The learner model includes personalized features of learners and interactive behaviors with e-textbook. The instructional model recommends appropriate teaching resources and skills. The adaptive engine is responsible for recommending the learning strategies and adjusting learning materials according to the adaptive rules. At present, the research on these four parts is not completely mature (Boulanger & Kumar, 2019).

(Huang, Yudelson, Han, He, & Brusilovsky, 2016) summarizes the latest development of intelligent textbooks in the past. Previous studies mainly focused on content model and learner model. The first generation of adaptive textbooks focuses on tracking learners' knowledge status and uses adaptive navigation technology to recommend students to read the most relevant pages (Thaker, Brusilovsky, & He, 2019). (Bommanapally, Subramaniam, Parakh, Chundi, & Puppala, 2020) constructs a knowledge repository of course learning objects in order to automatically generate personalized e-textbook. Among them, the knowledge content model is relatively complex, and the learner model rarely considers the students' personalities, which is relatively simple.

In the future, more personalities of learners should be considered into intelligent textbooks. And intelligent textbooks should not be regarded as a repository for learning and teaching activities (Ritter et al., 2020) or a tool of collecting data (He, 2014; Yin, Ren, Polyzou, & Wang, 2019), but as an adaptive system which make decisions by adaptive technologies.

### 2.2 Automatic Detection of LS Based on Data Driven Approach

Learners' learning style is not unchangeable. Spending a lot of time on the questionnaire will reduce learners' learning motivation (Dorça, Araújo, Carvalho, Resende, & Cattelan, 2016) Compared with the conventional ones, the data-driven approach has the following advantages. Firstly, it is more objective. Data-driven methods are based on data mining and machine learning technology (Aissaoui, Madani, Oughdir, & Allioui, 2019). The result is not affected by learners' subjective comprehension. Secondly, data-driven methods are dynamic. The e-textbooks will provide timely feedback on changes in students' reading behaviors. Thirdly, data-driven methods are more accurate. That is because the prediction results of LS are based on a large amount of students' reading behavior data.

According to the relevant literature, it can be found that researchers have been always looking for a data-driven mechanism for automatic detection of LS. For example, (Truong, 2016) classified data-driven detection of LS into three sub problems: consideration of learners' personal traits, selection of LS models and classification algorithm selection of learner model. Sub question 1 is the first step of learner model. (Normadhi et al., 2019) divides learners' personal characteristics into four categories: the mixture of cognition, affection, behavior or psychomotor and mix. For sub problem 2, previous studies have proved that FSLSM is the most suitable for adaptive e-textbook to detect LS compared with other LS models (Pham & Adina, 2013; Bernard, Chang, Popescu, & Graf, 2017) . For sub problem 3, many researchers use different classifiers to implement the automatic prediction of LS (Bernard, Chang, Popescu, & Graf, 2017; Garcia, Amandi, Schiaffino, & Campo, 2007; Sheeba & Krishnan, 2019). Because the data-driven approach needs sufficient training data to achieve accurate identification of

personal traits, in this study, a mixture method of questionnaire and data-driven approach is used to detect learners' learning styles.


## 3. Data

With the help of e-textbook system, a system established by (Yin, Ren, Polyzou, & Wang, 2019), data was collected from a total of 234 students were collected. Each record contains students' ID, gender, scores of LS and reading behaviors data. There are 16 behavioral characteristic variables generated by students, as shown in Table 1. We add a feature variable SumMarkerC, which is used to represent the total number of students taking notes.

Table 1. *Seventeen Behavioral Characteristic Variables*

| Variable Name: Description | Variable Name: Description |
|---|---|
| *PCC*: times of login using PC terminal | *ReadPages*: reading pages |
| *TabletC*: times of login using tablet terminal | *ReadTime*: reading time |
| *MobileC*: times of login using mobile terminal | *PrevC*: times of turning pages back |
| *HighLightC*: times of marking highlight | *NextC*: times of turning pages forward |
| *UnderLineC*: times of marking underline | *BacktrackRate = PreC / NextC* |
| *BookMarkC*: times of marking bookmark | *Pretest*: pretest score |
| *MemoC*: times of marking memo | *Middle*: middle score |
| *MarkerC: HighLightC + UnderLineC* | *GPA*:grade point average |
| *SumMarkerC*: total times of *HighLightC*, *UnderLineC*, *BookMarkC* and *MemoC* | |

The scores of LS are measured by FSLSM, which is used for model training and performance evaluation. FSLSM consists of 44 questions, which are divided into four dimensions to describe learners' preference of processing, perceiving, inputting and understanding information. Each dimension has 11 questions. And each question has two options, "a" and "b". The number of "a" minus the number of "b" equals the score of each dimension. It is worth noting that the score can be only restricted to odd numbers between - 11 and 11. Therefore, 12 types of LS are finally formed, as shown in Table 2.

Table 2. *Classification of 12 types of LS*

| Dimension | Variable Name | Description | Learners' LS |
|---|---|---|---|
| D1 | | *activescore*∈[3,5,7,9,11] | reflective |
| | *activescore* | *activescore*∈[-1,1] | balance |
| | | *activescore*∈ [-11,-9,-7,-5,-3] | active |
| D2 | | *sensingscore*∈[3,5,7,9,11] | intuitive |
| | *sensingscore* | *sensingscore*∈[-1,1] | balance |
| | | *sensingscore*∈[-11,-9,-7,-5,-3] | sensing |
| D3 | | *visualscore*∈[3,5,7,9,11] | verbal |
| | *visualscore* | *visualscore*∈[-1,1] | balance |
| | | *visualscore*∈[-11,-9,-7,-5,-3] | visual |
| D4 | | *sequentialscore*∈[3,5,7,9,11] | global |
| | *sequentialscore* | *sequentialscore*∈[-1,1] | balance |
| | | *sequentialscore*∈[-11,-9,-7,-5,-3] | sequential |

The personalized characteristics of the balance are relatively not obvious while the other two types are obvious and the differences are relatively large. Therefore, this study first eliminates learners'

data of the balance, and transforms the multi classification prediction problem into a binary classification problem.

## 4. Methods

Learning style prediction process is divided into five steps, including data exploration, determination of prediction target, dimension reduction, construction of learner model and evaluation of model performance, as shown in Figure 1.
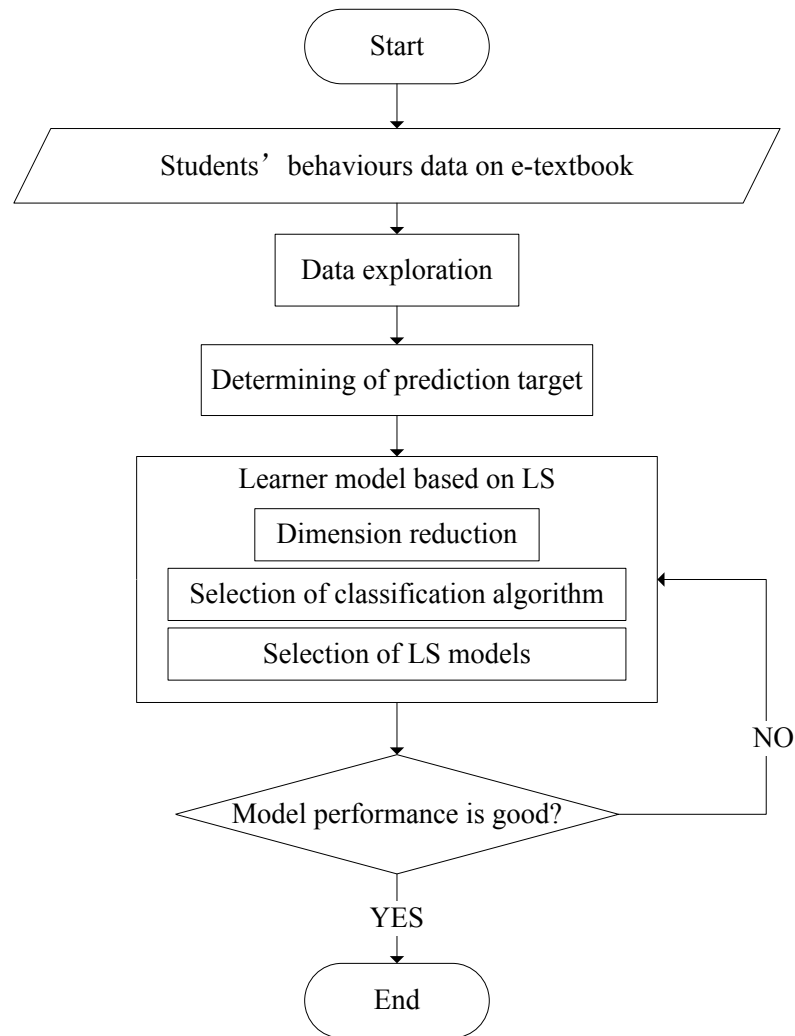


*Figure 1.* Flow chart of LS prediction process

The data exploration mainly focuses on the distribution of LS and the differences of reading behaviors among LS. The distribution will affect the choice of the selection of evaluation methods and evaluation indicators. If there exists the problem of data imbalance among categories, the three methods can be used before next step: over-sampling, under-sampling or mixed-sampling.

Feature selection and feature extraction are usually used to achieve dimensionality reduction. Feature selection is to directly select a few dimension data from high-dimensional data to represent the whole; feature extraction reduces the dimension of feature matrix by establishing mapping relationship between high-dimensional data and low-dimensional data. In this study, PCA (Jolliffe, 1986) and Lasso (Tibshirani, 1996) were used. PCA is a feature extraction method and Lasso is a feature selection method.

The construction of the learner model is mainly divided into the following parts: determining the LS model, selecting the classifier, dividing the proportion of training set and test set, and determining the super parameters. In this study, we would like to attempt to use four types of classifiers, including LR (Cox, 1958), NB (John & Langley, 1995), DT (Quinlan, 1992) and SVM (Cortes & Vapnik, 1995).

In addition, it is necessary to consider the evaluation methods and indicators of model performance. In order to reduce over fitting in a certain extent and obtain as much effective information as possible from limited data, 10 fold cross validation (Golub, Heath, & Wahba, 1979) is introduced. Considering the imbalance of data, it is not rigorous to evaluate the performance of the model only with the accuracy and F1 measurement is selected.

## 5. Results

### 5.1 Data Exploration and Analysis

The distribution of 234 students' LS is shown in Figure 2. The conclusion can be drawn that the distribution of D1 and D4 is relatively balanced, while that of dimensions D2 and D3 is very unbalanced. Therefore, we use the mixed sampling method to preprocess the data for D2 and D3.
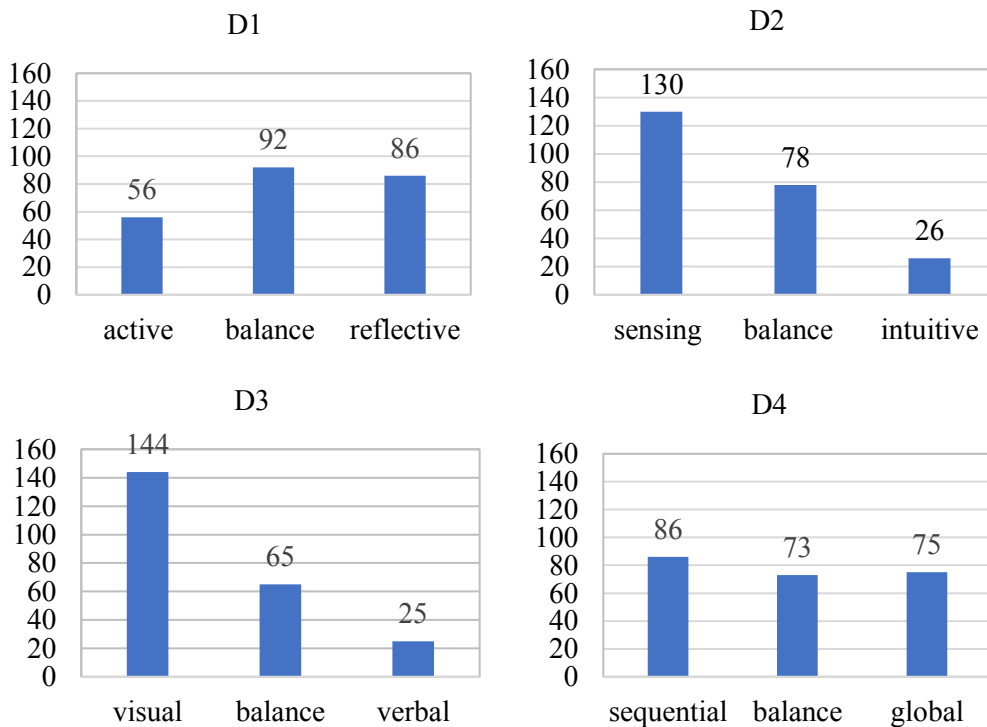


*Figure 2*. Distribution of 12 Types of LS in Four Dimensions

This study aims at e-textbook identifying students' LS automatically and more accurately in terms of students' reading behaviors data. At first, we should learn about whether there are significant differences among groups with different LS and what the differences are. One-way ANOVA is adopted. Before that, independence, normality and homogeneity of variance within the group passed the test.

Table 3. *Variables with significant difference in D3*

| Reading behaviors | visual (N=144) | balance (N=65) | verbal (N=25) | $F$ | $p$ |
|---|---|---|---|---|---|

|  | M | SD | MD | M | SD | MD | M | SD | MD |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|
| *BookMarkC* | 1.93 | 4.54 | 1.00 | 2.71 | 4.83 | 1.00 | 3.72 | 9.41 | 0.00 | 3.08 | 0.048 |
| *HighLightC* | 14.17 | 17.92 | 6.00 | 22.15 | 27.25 | 15.00 | 25.32 | 26.53 | 26.00 | 3.69 | 0.027 |
| *MarkerC* | 15.45 | 19.12 | 7.00 | 26.22 | 33.57 | 16.00 | 26.24 | 27.04 | 29.00 | 3.94 | 0.021 |
| *SumMarkerC* | 27.96 | 32.53 | 15.00 | 43.23 | 48.14 | 29.00 | 51.00 | 50.38 | 58.00 | 4.57 | 0.011 |

The results show that the first dimension has significant difference in *MobileC* ($p= 0.045$), marginal significant difference in *HighlightC* ($p= 0.065$), *MarkerC* ($p= 0.077$) and *SumMarkerC* ($p= 0.095$). The second dimension has marginal significant difference in *ReadPages* ($p= 0.065$) and *NextC* ($p= 0.066$). The fourth dimension has marginal significant difference in *GPA* ($p= 0.065$).

The third dimension has significant difference in *BookMarkC* ($p= 0.048$), *HighLightC* ($p= 0.065$), *MarkerC* ($p= 0.027$) and *SumMarkerC* ($p= 0.021$), marginal significant difference in *MemoC* ($p= 0.055$), *Middle* ($p=0.074$) and *GPA* ($p= 0.057$), as shown in the Table 4.



*Figure 3.* Radar Chart of Students' Differences in Taking Notes in D3

It can be obviously concluded in Figure 3, verbal learners make notes the most of *BookMarkerC, HighLightC, MarkerC* and *SumMarkerC*, while visual learners do the least. Comparing the median and the average of the four reading behaviors, we can find that more than 50% of the verbal learners are above average. More than 50% of the visual learners are below average and the balanced learners are in the middle. In other words, most verbal learners have the habit of taking notes, while most visual learners do not.

## 5.2 Evaluation of Learner Model

We select the popular F1 measure as the evaluation indicator of model performance. According to calculation formulas of F1, Table 4 lists results of the performance evaluation. For D2 and D3, DT classifier have the best prediction performance compared with other classifiers. The maximum F1 measurement of D2 can reach 83.33% and that of D3 can reach 95.63%. For D1 and D4, the four classifiers have poor prediction performance, and F1 measurement ranges from 49.59% to 66.00%. Generally, there is no significant difference between the two methods of dimension reduction.

Table 4. *Results of the Model Performance*

| F1 | D1 | | D2 | | D3 | | D4 | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | PCA | LASSO | PCA | LASSO | PCA | LASSO | PCA | LASSO |
| LR | 63.49% | 66.00% | 71.43% | 75.06% | 70.00% | 61.25% | 54.71% | 55.88% |
| NB | 60.23% | 55.35% | 72.86% | 70.48% | 78.33% | 73.12% | 57.06% | 49.96% |
| DT | 58.37% | 62.00% | 81.43% | 79.40% | 83.33% | 95.63% | 49.59% | 53.31% |
| SVM | 61.33% | 61.40% | 78.57% | 71.96% | 81.67% | 66.25% | 52.94% | 57.76% |

## 6. Discussion

In this study, the data collected from e-textbook is first mixed-sampled to solve the problem of imbalance. Then a learner model for predicting LS is successfully established. The model is evaluated by confusion matrix method. The results show that the model performs well in the two dimensions of the FSLSM. It can be observed from the results obtained that the learner model may be used to identify students' LS based on their reading behaviors data in e-Textbook. The data-driven approach is used to automatically identify students' LS in this study. This method avoids the external interference and subjective understanding bias caused by the traditional LS measurement method based on the questionnaire. It has the advantages of dynamic adaptiveness, objective feedback and higher accuracy.

Thus, by identifying students' LS, the model can not only guide the development of intelligent textbooks, but also recommend learning materials suitable for learners with different LS in order to improve the learning process. For example, more than half of verbal learners have the habit of taking notes while visual learners do not in this sample. For verbal learners, intelligent textbooks can first recommend more text-based learning resources or additional supplementary materials; for visual learners, visual learning resources such as pictures are presented first. In addition, it has a positive impact on learners. Learners can perceive their own learning preferences and obtain personalized learning materials, which can reduce their cognitive load and improve their self-confidence (Durak & Saritepeci, 2018).

As a family member of adaptive systems, intelligent textbooks provide personalized feedback and support through learners' autonomous learning. Instructional model, learner model, domain model and adaptive engine are four parts of the adaptive system, which have become valuable researches and are worth breaking through all the time. Learners' personal traits belong to the content of learner model, and most studies use LS to simulate learners' personalities. It is important to take advantage of LS automatic detection to design adaptive system so as to provide better personalized service (Boulanger & Kumar, 2019). In the future work, we will classify the learning resources and mark the key pages based on the results of this study (Deligiannis, Panagiotopoulos, Patsilinakos, Raftopoulou, & Symvonis, 2019). This will help provide learners with selective references of peers, who has the same LS, so as to read the key content back correctly and effectively for the purpose of providing support for personalization (Pursel, Ramsay, Dave, Liang, & Giles, 2019; Ritter et al., 2020).

**References**

Aissaoui, O. E., Madani, Y. E., Oughdir, L., &amp; Allioui, Y. E. (2019). Combining supervised and unsupervised machine learning algorithms to predict the learners' learning styles. Procedia Computer Science, 148, 87-96. doi:10.1016/j.procs.2019.01.012

Bernard, J., Chang, T., Popescu, E., & Graf, S. (2017). Learning style Identifier: Improving the precision of learning style identification through computational intelligence algorithms. *Expert Systems with Applications, 75*, 94-108. doi:10.1016/j.eswa.2017.01.021

Bommanapally, V., Subramaniam, M., Parakh, A., Chundi, P., & Puppala, V. M. (2020). Learning Objects Based Adaptive Textbooks with Dynamic Traversal for Quantum Cryptography. *Second Workshop on Intelligent Textbooks The 21th International Conference on Artificial Intelligence in Education (AIED'2020).*

Boulanger, D., & Kumar, V. (2019). An Overview of Recent Developments in Intelligent e- Textbooks and Reading Analytics. *First Workshop on Intelligent Textbooks The 20th International Conference on Artificial Intelligence in Education (AIED'2019), 2384*, 44-56.

Cortes, C., & Vapnik, V. (1995). Support-Vector Networks. *Machine Learning, 20*(3), 273-297. doi:10.1023/A:1022627411411

Cox, D. R. (1958). The Regression Analysis of Binary Sequences. *Journal of the Royal Statistical Society: Series B (Methodological), 20*(2), 215-232. doi:10.1111/j.2517-6161.1958.tb00292.x

Deligiannis, N., Panagiotopoulos, D., Patsilinakos, P., Raftopoulou, C., & Symvonis, A. (2019). Interactive and Personalized Activity eBooks for Learning to Read: The iRead Case. *Proceedings of the First Workshop on Intelligent Textbooks Co-located with 20th International Conference on Artificial Intelligence in Education (AIED 2019), 2384*, 57-69.

Dorça, F. A., Araújo, R. D., Carvalho, V. C., Resende, D. T., &amp; Cattelan, R. G. (2016). An Automatic and Dynamic Approach for Personalized Recommendation of Learning Objects Considering Students Learning Styles: An Experimental Analysis. Informatics in Education, 15(1), 45-62. doi:10.15388/infedu.2016.03

Durak, H. Y., & Saritepeci, M. (2018). Analysis of the relation between computational thinking skills and various variables with the structural equation model. *Computers & Education, 116*, 191-202. doi:10.1016/j.compedu.2017.09.004

Felder, R. M., & Silverman, L. K. (1988). Learning and teaching styles in engineering education.

Garcia, P., Amandi, A., Schiaffino, S., & Campo, M. (2007). Evaluating Bayesian networks' precision for detecting students' learning styles. *Computers in Education, 49*(3), 794-808. doi:10.1016/j.compedu.2005.11.017

Golub, G. H., Heath, M. T., & Wahba, G. (1979). Generalized Cross-Validation as a Method for Choosing a Good Ridge Parameter. *Technometrics, 21*(2), 215-223. doi:10.1080/00401706.1979.10489751

Gomede, E., Barros, R. M., & Mendes, L. D. (2020). Use of Deep Multi-Target Prediction to Identify Learning Styles. *Applied Sciences, 10*(5). doi:10.15417/1881

He, L. (2014). The Adaptive Teaching in the Setting of Big Data. *Proceedings of the 2014 International Conference on Education, Management and Computing Technology*. doi:10.2991/icemct-14.2014.79

Huang, Y., Yudelson, M., Han, S., He, D., & Brusilovsky, P. (2016). A Framework for Dynamic Knowledge Modeling in Textbook-Based Learning. *Proceedings of the 2016 Conference on User Modeling Adaptation and Personalization - UMAP '16, 141-150.* doi:10.1145/2930238.2930258

John, G. H., & Langley, P. (1995). Estimating continuous distributions in Bayesian classifiers. *Uncertainty in Artificial Intelligence,* 338-345.

Jolliffe, I. T. (1986). Principal Component Analysis. *Springer Series in Statistics*. doi:10.1007/978-1-4757-1904-8

Normadhi, N. B., Shuib, L., Nasir, H. N., Bimba, A., Idris, N., &amp; Balakrishnan, V. (2019). Identification of personal traits in adaptive learning environment: Systematic literature review. Computers &amp; Education, 130, 168-190. doi:10.1016/j.compedu.2018.11.005

Pham, Q. D., &amp; Adina, M. F. (2013). Adaptation to Learners' Learning Styles in a Multi-Agent E-Learning System. Internet Learning. doi:10.18278/il.2.1.2

Pursel, B., Ramsay, C., Dave, N., Liang, C., & Giles, C. L. (2019). BBookX: Creating Semi-Automated Textbooks to Support Student Learning and Decrease Student Costs. *Proceedings of the First Workshop on Intelligent Textbooks Co-located with 20th International Conference on Artificial Intelligence in Education (AIED 2019), 2384*, 81-86.

Quinlan, J. R. (1992). Programs for machine learning. San Francisco.

Ritter, S., Fisher, J., Lewis, A., Finocchi, S. B., Hausmann, B., & Fancsali, S. (2019). What's a Textbook? Envisioning the 21st Century K-12 Text. *First Workshop on Intelligent Textbooks The 20th International Conference on Artificial Intelligence in Education (AIED'2019).*

Rout, N., Mishra, D., & Mallick, M. K. (2017). Handling Imbalanced Data: A Survey. *Advances in Intelligent Systems and Computing International Proceedings on Advances in Soft Computing, Intelligent Systems and Applications,* 431-443. doi:10.1007/978-981-10-5272-9_39

Sheeba, T., & Krishnan, R. (2019). Automatic Detection of Students Learning Style in Learning Management System. *Smart Technologies and Innovation for a Sustainable Future Advances in Science, Technology & Innovation,* 45-53. doi:10.1007/978-3-030-01659-3_7

Thaker, K., Brusilovsky, P., & He, D. (2019). Student Modeling with Automatic Knowledge Component Extraction for Adaptive Textbooks. *Proceedings of the First Workshop on Intelligent Textbooks Co-located with 20th International Conference on Artificial Intelligence in Education (AIED 2019), 2384*, 95-102.

Tibshirani, R. (1996). Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society: Series B (Methodological), 58*(1), 267-288. doi:10.1111/j.2517-6161.1996.tb02080.x

Truong, H. M. (2016). Integrating learning styles and adaptive e-learning system: Current developments, problems and opportunities. Computers in Human Behavior, 55, 1185-1193. doi:10.1016/j.chb.2015.02.014

Yin, C., Ren, Z., Polyzou, A., & Wang, Y. (2019). Learning Behavioral Pattern Analysis Based on Digital Textbook Reading Logs. *International Conference on Human-computer Interaction,* 471-480. doi:10.1007/978-3-030-21935-2_36