

# A CASE STUDY FOR ANALYSIS AND PREDICTION IN LEARNING USING A PEER LEARNING SYSTEM

Adelina Aleksieva-Petrova<sup>1</sup> and Milen Petrov<sup>2</sup>

<sup>1</sup>*Technical University of Sofia, Faculty of Computer Systems and Technologies, Bulgaria*

<sup>2</sup>*Sofia University, Faculty of Mathematics and Informatics, Bulgaria*

## ABSTRACT

Using data from different learning systems and tools to analyze and predict student achievements and behaviors is a method to increase learning effectiveness. In this paper, the peer learning system is presented and used to simulate analysis and prediction processes in a learning environment. The aggregation and sorting methods enforce the analysis process. The Decision Tree classification method is applied and evaluated to predict the result.

## KEYWORDS

Learning Environment, Prediction, Student Behavior, Technology-Enhanced Learning

## 1. INTRODUCTION

The COVID-19 pandemic has significantly changed the world and influenced every sector of society. Technology Enhanced Learning (TEL) became extremely popular during the global epidemic as different platforms and tools must be used to transform all learning activities to remote or mixed mode. Using learning management systems (LMS), tools for videoconference, systems for assessments, and many others help to support the learning process and to deliver learning content to students in an effective way. These systems define the learning environment as a technology model which could be used by both students and teachers.

Most of the systems and tools are web-based, leading to the challenge of web browsers used by students and teachers. On the other hand, some tools are desktop applications, that require installation and configuration depending on the operating system. Therefore, knowing the preferences of learners in using OS or browsers is a challenge to setting up an appropriate learning environment.

Those systems also generate and collect data that could use for different analyses of student achievements or problems and predict student needs and behaviors which teachers could be used to improve the learning process and outcomes. This process is known as Educational Data Mining which “converts raw data coming from educational systems into useful information that could potentially have a great impact on educational research and practice” (Kaur, 2015). This data could give information about learner achievements and results and also help with selecting and configuring the learning environment.

All these challenges the APTITUDE project tries to resolve, specifically to propose a framework based on learning analytics for improving adaptation and recommendation of the learning process. The project has designed and developed the platform which is a middleware layer between different learning systems and tools to deploy the learning environment. Using data for student behaviors from external systems and tools the APTITUDE platform has analyzed and predicted the learning process and content.

This research aims to investigate the possibility of using systems learning data to be used in two directions as regards the analysis of learner preferences of different things and to predict different outcomes which could be used in the learning environment. To address the research goal, we investigate the data provided by a web-based peer learning system in the Web Technologies course at Sofia University.

This paper is organized as follows. After the introduction, Section 2 reviews the literature regarding related work. A brief conceptual description of the APTITUDE framework, main components, systems, subsystems, and its relations are given in Section 3. Section 4 presents the case study for an internal peer learning system. The analysis and prediction process using system data are described in Section 5 and the results are reported and interpreted. Section 6 concludes the paper.

## 2. RELATED WORKS

Different approaches are used to analyze learning data to predict student achievements and behavior. For example, the approach based on Bayesian Networks for modeling the behavior of the students is proposed in cases where the assessment or prediction should take into consideration very large amounts of data from a variety of sources. The results are promising as regards both predict of student behaviors, based on modeled past experience, and assessment (Xenos, 2004).

Using classification-based algorithms, such as multi-layer perception, Naïve Bayes, SMO, J48, and REPTree supported by the WEKA tool, the slow learners among students are identified and predicted. The results show that multi-layer perception performs best with 75% accuracy and proves to be a potentially effective and efficient classifier algorithm (Kaur, 2015). Machine learning algorithm Nave Bayes is also investigated to concentrate student efforts in a specific area to improve their academic achievement (Pallathadka, 2021). It is compared with the other three algorithms such as ID3, C4.5, and SVM. In this research SVM is the most accurate technique for classifying a data set of student performance.

A combination of Random Matrix Theory, a Community Detection algorithm, and statistical hypothesis tests is proposed to detect groups of students who have similar learning behaviors and outcomes (Mai, 2022). The experiments show that the higher performing groups are more active in practical-related activities throughout the course rather than groups that engage more with lecture notes.

The next research point is to investigate the learning environment and the use of different approaches to providing effective learning to students. For example, the introduction of educational games (gaming) and related virtual reality significantly improves the perception of educational material and the results obtained in all studied disciplines in the field of engineering and IT sciences (Getova, 2021). For that reason, several software instruments are created for analyzing and evaluating both the learning and playing outcomes of the player with different knowledge, learning goals, preferences, and learning styles (Dankov, 2021). In other research implementation of automated testing and evaluation is present which is involved in a learning environment (Chenchev, 2018).

The technology-enriched classrooms including videoconferences and video streaming lessons are investigated using a neuro-fuzzy model of Quality of Experience prediction for students. The results show that the perceived Quality of Experience is mostly influenced by the person's personality traits and learning style (Vasileva-Stojanovska, 2015).

## 3. A BRIEF INTRODUCTION TO THE APTITUDE PROJECT

The Aptitude project proposes a flexible framework aimed to recommend and adapt learning content and activities from/to different learning platforms, systems, and tools based on learning analytics (figure 1). For the conceptual system design of the framework, the main software components, systems, subsystems, tools, and its relations are given here. For the components on the diagram, we use a similar notation to the UML class diagram, but instead of classes, we represent software components. As the proposed overview is at a high level, UML component diagrams are not appropriate, as they represent too many specifics on the implementation of the components. The central system component here is the Aptitude server.

The Publisher and Consumer Registers components contain a list of systems and tools which are registered into the Aptitude platform. The Publisher Register is where the source of the data from systems log files and event streams is used in learning analytics. When considering the learning content and activities in different learning systems, we define the micro and macro level of the adaptation and the recommendation processes. The micro-level is related to learning resources and activities inside the system itself, while the macro level is related to learning resources and activities that are outside the system. Conversely, the Consumer Register is a component that has a role as external consumers of services offered by the Aptitude framework.

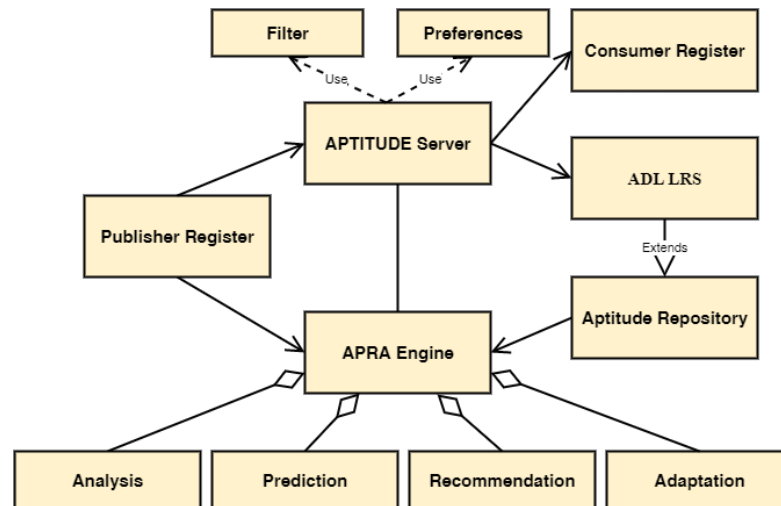


Figure 1. Meta model of Aptitude framework

Filter and Preferences components help to set the configurations and to establish a connection between the registered learning system or tools and the APTITUDE Server. The students are identified through ADL Learning Record Store based on the xAPI standard and stored in Aptitude Repository. Also, reference to learning content and activities is collected into the repository too. APRA engine is a core component, working behind the scenes, and is a complex component, that aggregates four subsystems for Analysis, Prediction, Recommendation, and Adaptation. It is outside the scope of the current paper, but in general, it manages its sub-components and has bi-directional interaction with the APTITUDE server: from one side it collects data from the server and its state, feeding four subcomponents, and from the other, it supports the work of the server for adaptation and recommendation processes.

#### 4. CASE STUDY DESCRIPTION

The presented framework was applied in the last semester of the academic year 2021/2022 in the Web technologies course at Sofia University for a bachelor's degree in Software Engineering. Publisher Register is an internal peer learning system that is deployed and set up. The system is a web-based application that provides peer review and peer assessment as a subset of the peer learning process (Petrov, 2022). The students can see and choose a topic in the subject of Web technologies and they should write an essay based on paper bibliographies which provide a flexible and non-invasive, in which one can quantify actual student behavior involving library products and resources (Hovde, 2000). In such a way they would need to locate appropriate resources for the assigned research paper.

When the student essay is uploaded into the system the peer review process is started. The students could anonymously comment on individual parts of a submitted essay or write a general comment and assign a grade (peer assignment). All students see the comments and grades, but without knowing which peer student wrote the comment.

To create a model of analysis of learning data and prediction, we have simulated two components from the aforementioned APRA Engine using the system log file in two directions:

- Analysis process – learning preferences analysis according to learning content, and according to used browsers;
- Prediction process – here prediction is defined conceptually, as most of the information about students in the current study is confidential, so we use only prediction of operating system usage. Even though this didn't give much insight into the learning it positions the prediction system component as a viable one in the whole process.

The system log file with a total of 21767 data instances and 9 variables are taken for the analysis of this research. The log file includes information for the type of learner activity, started date of activity, learner device type, operating system, browser, description, anonymized ID (identifier), and IP address of the user.

In this paper, the selected potential variables are id, browser, device\_type, os, type, IP\_ANON, and ANON\_ID. The variables started\_at and description are used as meta attributes.

## 5. RESULTS

### 5.1 Analysis Process

For data processing, in this section, we use Talend Open Studio, which is open source and allows us to “execute simple extract, transform and load (ETL) and data integration tasks of data, and manage files”<sup>1</sup>. The workflow in the Talend software tool is implemented to create data analysis (figure 2). To identify the most interesting essay topics and most using OS and browser data aggregation and sorting are provided.

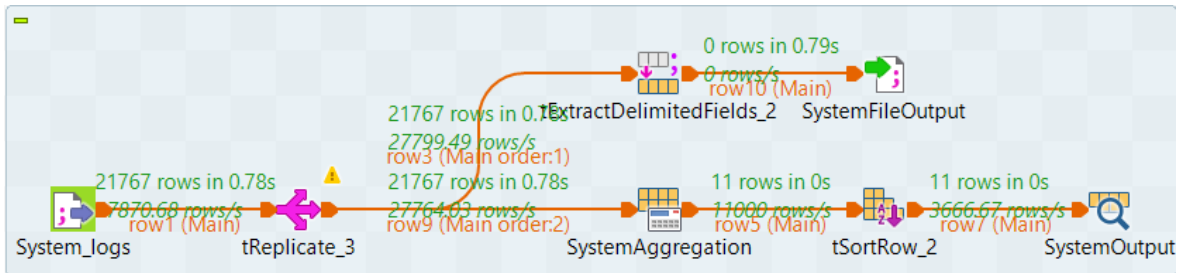


Figure 2. Workflow for the Analysis process

Figure 3 shows the statistics for the most used browsers and most reading essays. The results show that most students use Chrome browser, followed by Firefox, Opera, Edge, and Safari. There is a significant number of users, browsing the system via Facebook mobile application.

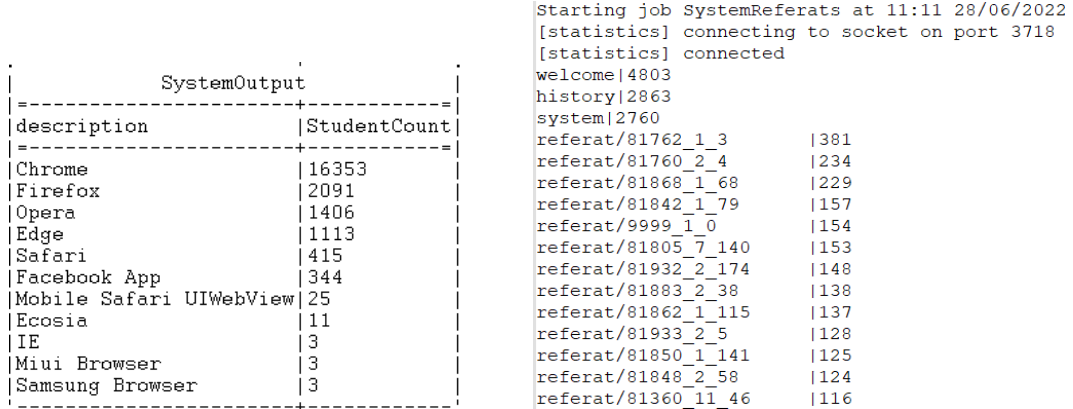


Figure 3. Experimental data results on left, and top reading results on the right

The statistics for the most visited pages are the welcome page (landing page of the system, after successful authentication); history page (with the last results of the submitted works), system pages, and follows (anonymously) ranking on the reading results. The ranking on the reading results shows which essay is most browsed and read. A test account (with id 9999\_1\_0) has published an example of an essay for students which is first published in the system, and it ranks in the top 5 readings of the students. The results can help in ranking the essays as content and quality, ranking the authors of the essays, as the most credible reviewers, and using the results as input for classifications. It can be used as pointing out good practices, and a good example for other students.

<sup>1</sup> <https://www.talend.com/products/talend-open-studio/>

## 5.2 Prediction Process

To elevate student satisfaction with a given technology-enhanced educational setup, we used the classification technique as a Decision Tree to predict which OS is used by students. First, we used a random sample with 75 % of the data to create a model based on the given training data. This model is used to predict the response for the remaining 5441 instances.

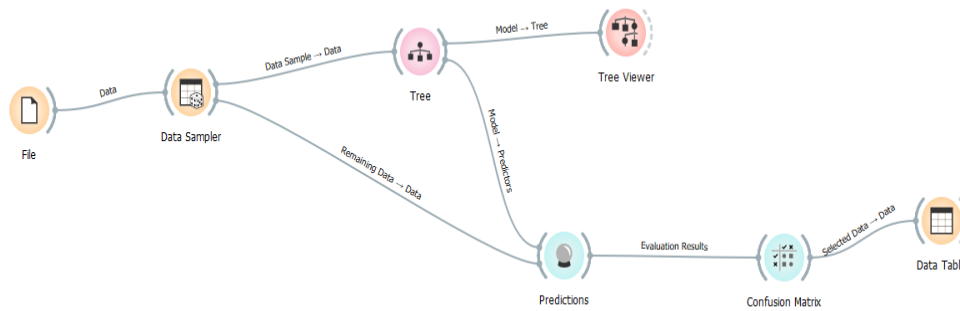


Figure 4. Workflow for the Prediction process

To measure the effectiveness of our model confusion matrix is used (figure 5). The Precision, Recall, and F1 of the model are 0.987 which means the high accuracy of the model. F1-score combines two competing metrics: precision and recall. The best results in prediction are achieved for mobile OS, such as Android, iPad, and iPhone.

The next important measurement is AUC (Area Under the Curve) ROC (Receiver Operating Characteristics) curve which is an important evaluation metric for checking any classification model's performance. The AUC value in our case is 0.992 which is considered outstanding and means better performance of the model.

		Predicted									
		Android	Linux	Mac OS X	Ubuntu	Windows 7	Windows 8.1	Windows 10	iPad	iPhone	Σ
Actual	Android	205	0	0	0	0	0	0	0	0	205
	Linux	0	492	1	0	3	0	23	0	0	519
	Mac OS X	0	3	617	1	0	0	20	0	0	641
	Ubuntu	0	0	0	160	0	0	0	0	0	160
	Windows 7	0	0	0	0	101	0	4	0	0	105
	Windows 8.1	0	0	0	0	0	25	7	0	0	32
	Windows 10	0	2	4	1	0	0	3724	0	0	3731
	iPad	0	0	1	0	0	0	0	6	0	7
	iPhone	0	0	0	0	0	0	0	0	41	41
	Σ	205	497	623	162	104	25	3778	6	41	5441

Figure 5. Confusion matrix

## 6. CONCLUSION

The APTITUDE project has the main role in learning analytics for adaptation and recommendation of learning contents and activities. Four components are the core of the platform which deliver the main processes: Analysis, Prediction, Recommendation, and Adaptation.

The current case study focuses on the validation of the Analysis and Prediction components. The data from the peer learning system is investigated to create a model for prediction. The results are promising with the high values of AUC and F1.

The limitation of the work is that the collected parameters (variables) for every log record are too limited to provide a more meaningful conclusion. According to Winne (2020) it is important to decide what trace data should be gathered and how trace data can contribute to recommendations for improving learning. Better prediction can be dependent on popularity, according to the subject, type of essay (scientific, overview, with more practical examples, or literacy), and level of essay (easy, intermediate, advanced). To produce such results, it needs external labeling of the essays.

Therefore, our future work will be to enrich the number of parameters and make colorations between them.

## ACKNOWLEDGEMENT

The research reported here was funded under a project entitled “An innovative software platform for big data learning and gaming analytics for a user-centric adaptation of technology enhanced learning (APTITUDE)” by the Bulgarian National Science Fund with contract №: KP-06OPR03/1 from 13.12.2018.

## REFERENCES

- Chenchev, I., (2018). MySQL Scripts Evaluation with Moodle and Virtual Programming Lab. *Journal “Cax technologies*, issue 6, ISSN 1314-9628, pp.83-87
- Dankov, Y. et al, (2021). Designing software instruments for analysis and visualization of data relevant to playing educational video games. International Conference on Human Interaction and Emerging Technologies. Springer, Cham.
- Getova, I., et al. (2021). "Results and trends in e-learning of students, conducted during a lockdown. *ICERI2021 Proceedings. RATED*
- Hovde, Karen. (2000). Check the citation: library instruction and student paper bibliographies. *Research Strategies* 17.1 3-9
- Kaur, P. et al, (2015). Classification and prediction-based data mining algorithms to predict slow learners in education sector. *Procedia Computer Science*, 57, 500-508.
- Mai, Tai Tan, et al, (2022). Learning behaviours data in programming education: Community analysis and outcome prediction with cleaned data. *Future Generation Computer Systems* 127 (2022): 42-55
- Pallathadka, Harikumar, et al, (2021). Classification and prediction of student performance data using various machine learning algorithms. *Materials today: proceedings*.
- Petrov, M. Y. (2022). Design and implementation of a configurable software architecture for peer learning. Paper presented at the 2022 57th International Scientific Conference on Information, Communication and Energy Systems and Technologies, *ICEST 2022*, doi:10.1109/ICEST55168.2022.9828640
- Vasileva-Stojanovska, Tatjana, et al. 2015. An ANFIS model of quality of experience prediction in education. *Applied Soft Computing* 34 (2015): 129-138
- Winne, Philip H. (2020) Construct and consequential validity for learning analytics based on trace data. *Computers in Human Behavior* 112 106457
- Xenos, Michalis, (2004). Prediction and assessment of student behaviour in open and distance education in computers using Bayesian networks. *Computers & Education* 43.4 (2004): 345-359.