



UNIVERSITY of  
RWANDA

*Research and Postgraduate Studies  
(RPGS) Unit*

---

Research thesis title:

**PREDICTING HIGH SCHOOL STUDENTS' PERFORMANCE ON E-LEARNING  
USING EXPERIENCE APPLICATION INTERFACE (xAPI) AND ARTIFICIAL  
INTELLIGENCE**

By

**Marcel NYIRINKINDI**

215026054

Supervisor: Associate Professor, Dr. Ghislain Maurice Norbert ISABWE

Co-supervisor Dr. Frederic NZANYWAYINGOMA

A dissertation submitted in partial fulfillment of the requirements for the degree of Master of  
Science in software engineering

The University of Rwanda/College of Science and Technology

NYARUGENGE CAMPUS KIGALI, RWANDA

2023

## **Declaration**

I, NYIRINKINDI Marcel, declare that this thesis titled PREDICTING HIGH SCHOOL STUDENTS' PERFORMANCE ON E-LERNING USING EXPERIENCE APPLICATION (xAPI) AND ARTIFICIAL INTELLIGENCE, and that work within it is my own and has been generated by me as a result of my original research. I assured that:

1. This work was done in candidacy for a Master of Science degree at the University of Rwanda College of Science and Technology.
2. Wherever I have consulted the printed work of others, this can always be attributed.
3. Wherever I have quoted from the work of others, the supply is often given. Except for such quotations, this thesis is entirely my very own work.
4. I have acknowledged all primary sources of facilitation.
5. Wherever the thesis is predicated on work done alone conjointly with others, I have explicitly shown what was done by others and what I have contributed myself.

## **Acknowledgments**

I am thankful to God for the best fitness and health that have been essential to finish this thesis.

I sincerely thank Dr. Frederic Nzanywayimana, Senior Lecturer and Postgraduate Coordinator at the School of ICT, University of Rwanda College of Science and Technology Nyarugenge Campus, for supplying me with all the important research resources and advice, added for being my Co-supervisor.

I add on record my honest thanks to Associate Prof. Dr. Maurice Ghislain Isabwe from the University of Agder in Norway, on top of being my Main Supervisor, the continuous encouragement, logical orientations, advice, and support on the devices I used mainly in Machine Learning model building and predictions.

I am additionally thankful to Eva Mirembe and Dr. Godfrey Mayende from Makerere University. I am incredibly grateful and indebted to them for sharing their expertise, honest and precious guidance, and encouragement prolonged to me.

I take this possibility to send specific gratitude to all the Department college participants for their assistance and commitment during my journey of studies. I also thank my dad and mom for their encouragement, help, and attention. I am further thankful to my companion, who supported me through this venture.

Furthermore, with a documented location, my gratitude extends to everyone who has directly or indirectly contributed to this endeavor.

## **Abstract**

Motivated by the importance of education to the development of individuals and society, this research investigates using Artificial Intelligence (AI) to predict student performance. This thesis pays special attention to the current problem by evaluating the performance once AI is used as a tool to add value to existing teaching strategies for long-term improvements. Specifically, the thesis explores applications of Artificial Intelligence in education: predicting student performance. The thesis identifies potential challenges, current limitations, and suggestions for further improvement. Rigorous analysis using Artificial Intelligence will lay a solid foundation for further study within the domain. This research discusses the possibility of using xAPI (experience Application Interface) for information collection for Artificial Intelligence dataset input to predict student performance with an e-learning management system. This dissertation pays special attention to the factors influencing academic performance and, viable ways to anticipate and predict students' long-term achievements.

This research provides a way of anticipating student's results using a machine learning-based system, to predict students' marks. Using Jupyter Notebook from Anaconda software the dataset is processed. Using 4 types of models namely Logistic Regression (LR), Support Vector Regression (SVR), K Nearest Neighbours (KNN), and Artificial Neural Network (ANN), the results predicted are in the range between 40 and 85 out of 100, with an accuracy of 95%. From the results also we detected the need of increasing the dataset size, especially increasing quantitative dependent variables for prediction.

During research we also realized that e-learning can't be implemented alone, hence the requirement of a blended teaching system.

Keywords: Alin Education, Machine Learning, Performance-Prediction, e-learning

## List of acronyms

**%:** Percentage

**4G:** Fourth Generation

**ADL:** Advanced Distributed Learning

**AI:** Artificial Intelligence

**ANN:** Artificial Neural Network

**ANOVA:** Analysis of Variance

**DL:** Deep Learning

**DT:** Decision Tree

**EDM:** Education Data Mining

**FT:** Tensor Factorization

**GB:** Giga Bytes

**GHZ:** Giga Hertz or Giga Cycle (GC)

**HMM:** Hidden Markov Model

**IBM:** International Business Machine

**ICT:** Information Communication Technology

**KC:** Knowledge Centred

**KC:** Knowledge Components

**KDD:** Knowledge Discovery and Data Mining Challenge

**KNN:** K Nearest Neighbours

**LMS:** Learning Management System

**LR:** Logistic Regression

**LRP:** Learning Record Provider

**LRS:** Learning Record Store

**MF:** Matrix Factorization

**MI:** Mutual Information

**ML:** Machine Learning

**MOOC:** Massive Open Online Course

**OC:** Opportunity Count

**PISA:** Program for International Student Assessment

**RAM:** Random Access Memory

**RF:** Random Forest

**ROC:** Receiver Operation Characteristic

**SDMS:** School Data Management System

**SPP:** Student Performance Prediction

**SVM:** Support Vector Machine

**SVR:** Support Vector Regression

**TEAM:** Temporal Emotional Aspect Model

**TEL:** Technological Learning Environment

**TVET:** Teacher Vocational Education Training

**UK:** United Kingdom

**xAPI:** Experience Application Interface

## List of figures

Figure 1: The flow of data processing in machine learning.....	23
Figure 2: System architecture design.....	28
Figure 3: xAPI features.....	30
Figure 4: Script for collecting usernames and emails from Storyline .....	32
Figure 5: Script for sending a statement to Veracity .....	33
Figure 6: xAPI parameters configuration in Veracity (Irs.io).....	34
Figure 7: Configuration of Storyline triggers .....	35
Figure 8: Publish contents using Storyline by choosing one of the available options.....	35
Figure 9: Script for collecting activities from the storyline.....	36
Figure 10: Script for collecting open answers .....	36
Figure 11: Script for sending open text answers.....	37
Figure 12: Script for collecting quiz score.....	37
Figure 13: Script for collecting the time spent on the activity.....	38
Figure 14: Script to convert to ISO time.....	39
Figure 15: Send the final time spent in the system .....	39
Figure 16: System flowchart.....	40
Figure 17: List of used libraries .....	42
Figure 18: Statistical exploration of the dataset.....	43
Figure 19: The proportion of students by gender.....	43
Figure 20: The parent's level of education .....	44
Figure 21: Type of lunch often taken.....	44
Figure 22: The number of students who took the test preparation course .....	45
Figure 23: Students grades distribution .....	46
Figure 24: The number of exams a student succeeded .....	47
Figure 25: Performance by gender in math.....	48
Figure 26: Performance by gender in JavaScript.....	49
Figure 27: Performance by gender in VueJS .....	50
Figure 28: The relation between different subjects.....	51
Figure 29: Corelation between different parameters .....	52
Figure 30: Summary of the prediction using logistic regression .....	56

Figure 31: Distribution of predicted scores using logistic regression.....	57
Figure 32: Plot of test RMSE vs K number of neighbours .....	58
Figure 33: Distribution of predicted scores for KNN regression.....	58
Figure 34: Support vector regression RMSE and predictions .....	59
Figure 35: Prediction using the lowest RMSE in combination with grid search.....	59
Figure 36: Distribution of predicted scores for SVR .....	60
Figure 37: Comparison between models.....	61
Figure 38: Mean Absolute Error (MAE) and loss on each epoch.....	62
Figure 39: Model cross-validation score.....	63
Table 1: Advantages and disadvantages of used models.....	54

## Contents

Declaration.....	ii
Acknowledgments.....	iii
Abstract.....	iv
List of acronyms .....	v
List of figures.....	vii
CHAPTER 1 INTRODUCTION .....	2
1.1 Preamble.....	2
1.2 Background and motivation .....	2
1.3 Problem statement .....	3
1.4 Research questions .....	4
1.5 Study objectives .....	4
1.5.1 General objective .....	4
1.5.2 Specific objectives .....	4
1.6 Hypotheses.....	4
1.7 Study scope .....	4
1.8 Significance of the study.....	5
1.9 Organization of the study.....	5
1.10 Contribution .....	7
CHAPTER 2 STATE – OF - THE ART.....	8
CHAPTER 3 METHODOLOGY .....	16
3.0 Introduction.....	16
3.1 Type of research used.....	16
3.2 Data collection method.....	18
3.2.1 Data preparation.....	18
3.2.2 Data cleaning .....	19
3.2.3 Model selection.....	19

3.2.4 Model processing .....	20
source Image by Abid Ali Awan data scientist and blogger.....	23
3.3 Resources, materials and tools .....	23
3.4 Justification behind the analysis.....	24
<b>CHAPTER 4 SYSTEM DESIGN AND STATISTICAL ANALYSIS .....</b>	<b>27</b>
4.1 System architecture.....	28
4.1.1 System configuration .....	28
4.3 xAPI in combination with Articulate Storyline .....	30
4.3.1 Experience Application Interface(xAPI) required statement.....	31
4.3.2 Collecting username and emails of a user from the storyline .....	31
4.3.3 Get and send statements from the storyline .....	32
4.3.4 Configure storyline’s triggers .....	34
4.3.5 Collecting xAPI activities .....	36
4.3.6 Collecting open-text answers from the storyline .....	36
4.3.7 Collecting quiz score information.....	37
4.3.8 Measuring the duration a student spent on a certain activity.....	37
4.3.9 Data visualization in xAPI.....	39
4.4 System Flowchart.....	40
4.5 Statistical analysis.....	40
4.5.1 Importing necessary libraries.....	41
4.5.2 Exploring data.....	42
4.5.2.1 The proportion of students by gender .....	43
4.5.2.2 The Parents’ level of Education.....	43
4.5.2.3 Type of lunch often taken .....	44
4.5.2.4 The number of students who took the test preparation course.....	45
4.5.2.5 Students’ grades distribution.....	45

4.5.2.6 The number of exams a student succeeded.....	46
4.5.2.7 Performance by gender in math .....	47
4.5.2.8 Performance by gender in JavaScript .....	48
4.5.2.9 Performance by gender in VueJS.....	49
4.5.2.10 The relation between different subjects .....	51
4.5.2.11 Correlation between different parameters.....	52
4.5.2.12 Summary of probability .....	53
4.6 Building Artificial intelligence Model.....	53
4.6.1 Linear Regression Model.....	55
4.6.2 K Nearest Neighbours Regression.....	57
4.6.3 Support Vector Regression (SVR).....	59
CHAPTER 5 ANALYSIS OF THE RESULTS .....	61
CHAPTER 6 CONCLUSION AND RECOMMENDATION .....	63
6.1 Conclusions.....	63
6.2 Recommendations:.....	65
References.....	67

# CHAPTER 1 INTRODUCTION

## 1.1 Preamble

The Rwandan Government puts substantial efforts into Internet-related infrastructures, like spreading optic fibers across the country. Nowadays, everything is going digital, and Rwandan institutions for education can't stay behind. Artificial Intelligence has relevance to addressing education-related challenges (like predicting the long-run performance of students) that are rooted in the shortage of standard means of teaching. The present generation and the academic system are becoming more advanced. The use of AI and Machine Learning is becoming a must go for increasing the performance and the quality of education. Especially in particular cases, e-learning is becoming the trusted way to study without the risk of getting contaminated, like with the pandemic Covid-19.

## 1.2 Background and motivation

The field of education is changing its standards. Many decisions relating to students are taken based on their performance, namely advancing from lower level to advanced level, which generally defines students' future well-being. However, (Marks, 2006) found that the background of the student studies, and socioeconomic inequalities including the role of poverty, income and wealth, those factors influence considerably the studying journey of the student. From where the Rwandan government program to feed students at school to mitigate balanced diet problems influences on students' performance.

Through the use of AI , (Pojon, 2017) used a combination of 3 different machine learning methods and engineering features to make comparisons of how much improvement prediction performance. AI is employed to create shared insights from student learning data, allowing for customized learning experiences and aiding in decision-making. It helps improve the complex interaction between students and teachers, enabling teachers to proactively respond to anticipated student outcomes. Many students can make the wrong choice for their choice of studies because of not anticipating their performance earlier. This can affect the future life of the student.

### **1.3 Problem statement**

Students' performance can sometimes be measured on the results of marks obtained at the end of the term. Whether the result is good or bad, it means challenging to track the factors that lead to such a result. Taking adequate measures to improve performance is complex when the cause is unknown.

In high schools (ordinal levels) from S1 to S3 (Senior 1 to Senior 3), students take a mixed course for the purpose of identifying at least 3 courses where they perform better. At the end of the Ordinal Level (OL) they choose a combination to continue in Advanced Level (AL, TVET, or A<sub>2</sub>). The choice of combination should be based on at least 3 courses where the performance is better than others. Likewise when students are completing the advanced level, TVET (Teacher Vocational Education Training) certificates they must have succeeded in at least 2 main courses from their combination or from core and specific courses from TVET schools, in order to continue for the university level.

Early performance detection can help in identifying the student's talent, hence helping in guidance and orientation of students to follow either general courses or Technical and Vocational Education and Training (TVET) program. This will avoid the waste of time on the student side who sometimes get an advanced certificate at the end of Advanced Level (AL) with poor marks and that certificate becomes useless, because it doesn't fit in any requirements for the next level of study or in the labour market.

Not detecting the student's performance has many consequences. To name a few of them we can say: Student bad orientation: Many students follow study combinations just because their parents choose it, or their relative did the same. But, if the student is not involved in studies' decision making, most of time they study against their talent leading to failure. Early problem detection will help in the reorientation then anticipate possible future challenges. A student who missed opportunity due to problem we talked previously, they end up becoming burden to the country while the issue would have been handled if the prediction took place on the right time. Not handling this problem can increase the problem of education standards as students may not be able to compete internationally.

## **1.4 Research questions**

This thesis presents a solution to address learning patterns of students that result in different outcomes. The thesis , addresses the following research questions:

1. What are the causes of low students' performance?
2. What can be done to predict and anticipate better student performance?

## **1.5 Study objectives**

### **1.5.1 General objective**

To use machine learning techniques to predict student performance and then anticipate appropriate decisions to be made based on students' results.

### **1.5.2 Specific objectives**

1. To use different existing predictive models for results prediction namely Logistic Regression. Support Vector Regression, K- Nearest Neighbour, Artificial Neural Network

2. To use xAPI to prepare and publish digital learning resources, then track the students' journey and activities, aiming to identify weaknesses.

3. To assess the accuracy of the employed predictive models and select the one that demonstrates the highest performance for more accurate predictions.

## **1.6 Hypotheses**

Predicting students' performance earlier will increase the chance of students making the right choice of their combination at the end of the Ordinal Level , as well as increase the number of students who complete the requirements to continue to the university level.

## **1.7 Study scope**

This work investigates everyday factors in students' life that affect their learning performance, in order to predict how changes in those factors may affect future performance.

AI – based solutions can help in dealing with educational challenges that are inherent to standard method of teaching this generation and the academic system that is becoming bit by bit complicated. In such complexity, spotting how to facilitate a student and improve performance takes work. In particular, e-learning is creating a large quantity of knowledge that may adapt

Artificial Intelligence to address complex educational challenges and adopt more intelligent educational technology solutions, which might help predict the students' performance. These perspectives indicate that Rwandan educational institutions could adopt Artificial Intelligence to make performance predictions, take necessary measures then keep international standards.

## **1.8 Significance of the study**

The field of education is changing its standards. Many decisions relating to students are taken based on their performance, which defines students' future well-being. AI is applied to generate shared insights from student data, focusing on personalized learning experiences, facilitating decision-making, and improving student-teacher interactions for proactive interventions based on predicted outcomes.

## **1.9 Organization of the study**

Chapter 1: Introduction

It includes the background of this study, the problem statement, the general as well as the specifics objectives, the scope, the significance, the organization of the study and the final conclusion.

Chapter 2: State – of – the art.

This chapter present related works that have been done by other researchers on the field of education mainly about students' performance. Within this chapter we identify the gaps in current research and what is our contribution in this study.

Chapter 3: Research methodology

Within this chapter the methodology used for the study is described the same as the tools that are used.

Chapter 4: System design and statistical analysis

The chapter for analysis and design show how to use xAPI to collect data that are feed into machine learning algorithm to predict the students' learning performance results

Chapter 5: Results and analysis

The chapter for results and analysis shows the outputs from prediction model that has been made in machine learning.

#### Chapter 6: Conclusion and recommendation

The last chapter gives the conclusion based on the output of results from the model of prediction and gives the recommendations to the future researchers

## **1.10 Contribution**

In this thesis our contribution is to assess the accuracy of the employed predictive models and select the one that demonstrates the highest performance for more accurate predictions. Then use the best to predict students' results. Since there is a number of at-risk students, incorporation of student performance prediction model into education system is an urgent need for most educational entities (Zohair, 2019). Based on the evolution of education, and the vision 2050 of Rwanda, the usage of prediction system is becoming a possibility. This thesis demonstrates the importance of involving prediction system in education. This requires proper recording and storing of records of students for long-run which to enable the implementation of prediction system. Traditional recording of students' records is a barrier to this process. The current system of managing school data in a harmonized way School Data Management System ([SDMS](#)) owned by Ministry of education will help in providing data in order to predict students' advancement based on records from past academic years.

## CHAPTER 2 STATE – OF - THE ART

This chapter takes into account the related work that several researchers have done including the way of tracking the student's behaviour to predict performance. Here we are evaluating the current evaluation process and its impact on student performance and arranging the students based on the result they got from different evaluations. Considering the education context, understanding the student's point of view must be considered. Students are the key stakeholders in education, so their point of view, mainly on how they are evaluated, is one of the keys to their best performance. In (Zafer Unal, 2017) researchers found that students prefer multiple-choice questions compared to essay based. Results of the research conducted in a pharmacy school by (Stéphane, 2010) showed that 86% of the students found that a preliminary test was helpful, 84% said that class tests are essential, 81% were in favour of revision, then 63% were for tutorials.

One of the methods in the evaluation of the students is assessing the tasks given to them. Also, (Natriello, 2013) found out that communicating the purpose of assessment is important. The observations indicate a robust correlation between testing methods, objectives, task criteria provided, and student assessment criteria. Assessing students' evaluations involves taking into account how the process aligns with predefined goals and its resultant effects.

The observations reveal a significant connection between testing practices and the roles, tasks, criteria, and standards used in student analysis. The outcome statement of a student's study suggests that the evaluation process is influenced by the predefined objectives.

The analysis conducted on an artificial intelligence system (Shubham, 2021) shows that the systems use formal logic by applying an algorithmic program to see the degree of students' performance from high to lower level. Formal logic is a way of computing the supported degrees of truth rather than using the exact old true or false (1 or zero) Boolean algebra on which normally computers are based.

Results on students' effectiveness, as perceived by their colleagues, have accumulated over multiple semesters. (Stéphane, 2010). The impact is linked to system use, highlighting the value of centralizing colleague evaluations to assess essential skills in higher education using xAPI for tracking learning activities.

Advanced Distributed Learning (ADL) developed software that can records and track various types of learning experiences for learning systems, experience Application , which could overcome the restrictions of single-system recording. It permits the recording of learning processes from various e-learning environments, and thus the format of learning method data contains the “actor,” “verb,” and “object” (Can, 2015). xAPI permits lecturers to accurately understand students’ learning processes and to facilitate the development of learning process-based learning performance assessment and learning designation mechanisms (e-learning, 2015)). Since a colossal amount of learning method information is collected in inquiry primarily based learning environments, and they have multiple attributes, processing technologies like consecutive pattern mining, classification, and a number of them could also be used to explore implicit and vital knowledge and patterns from a colossal amount of learning technique data. The goals are to understand students’ learning behaviours (Jiawei, 2016),to predict and value learning performance, and even to spot the key factors that make sure learning performance by using automatic data analysis methods so that lecturers can modify their teaching ways and develop tons of adaptative teaching modes (Chih-Ming, 2019)Many Artificial Intelligence techniques are used to create precise models for predicting student behaviour and in-class performance. Learning Record Consumers (LRC) are clients who get the right of entry to Learning Record Store (LRS) facts for fact processing: interpretation, evaluation, translation, dissemination, aggregation, etc.

(Jonathan, 2016) highlights the crucial role of analysis tools within the xAPI ecosystem, particularly in the context of Learning Record Stores (LRS) for xAPI, as they play a pivotal role in examining and reporting on trends, gathering metrics, assessing, monitoring, and generating reports to validate the user experience. This assertion can be extended to encompass video-based education evaluation or applied more broadly to video analysis in general. While much of the initial research around xAPI has centered on learning techniques, as evidenced by studies such as those conducted by (Aneesha, 2016) (Alan, 2016) (Thomas, 2017), the versatility of xAPI allows for its application in various fields, including accessibility enhancement, as demonstrated by (Matjaž, 2011). It's essential to remember that xAPI serves as a means to an end, not an end in itself.

Predicting students' performance is a crucial aspect of modern education, facilitated by Artificial Intelligence technologies. This predictive capability enables educational institutions to implement targeted interventions for learners requiring additional support. While assessing student performance, it is essential to consider not only grades but also factors such as learning difficulties. In the context of identifying challenges faced by students in digital design, a study conducted by (Mushtaq, 2019) aimed to predict these issues. To achieve this, the study analyzed data logs from Technological Learning environments using various AI algorithms, including Artificial Neural Networks, Support Vector Machines, Logistic Regression, Decision Trees, and the Naïve Bayes model. Multiple AI models were trained using historical data to forecast student performance in new sessions.

To correlate school characteristics with student performance, AI and mathematical models were utilized in a study by (Chiara, 2018). The study examined student performance data from the 2015 Program for International Student Assessment (PISA) in nine countries. The authors of the study identified the interplay between school-related factors and student outcomes. They employed a mixed versatile tree-based technique to estimate inherent school factors, followed by a boosting model to map these factors to school-level variables. This approach aims to provide a robust platform for timely academic intervention and remedial actions by predicting possible future student performance.

The authors conclude that while learner and faculty characteristics significantly impact student achievements, the extent of these effects varies between countries. This underscores the importance of considering structural differences in international education systems when conducting studies and drawing conclusions.

To predict future student performance effectively, an educational platform should enable timely intervention and remediation. Education Data Mining (EDM) focuses on developing AI models for continuous performance prediction by uncovering hidden insights and patterns. Recent studies, including (Olugbenga, 2018) 's research, conducted experimental investigations comparing various information sources and classification techniques for predicting university students' academic performance.

The research compared algorithms like Decision Trees, Artificial Neural Networks, and Support Vector Machines to achieve higher accuracy in performance prediction. Additionally, (Livieris, 2019), and (Khan, 2019). proposed a semi-supervised combination-based approach to predict student performance early in the academic trimester using a programming model. The J48 algorithm, based on a recursive divide-and-conquer strategy, yielded promising results. They also applied Deep Learning with GritNet, a bidirectional long immediate memory-based model, for long-term student performance prediction, as demonstrated by (Ribeiro, 2018)."

Historically, predicting early course improvement has been challenging, but the authors achieved promising results. Existing literature includes studies analyzing text from students' journals and social network contributions to inform curriculum development (Ribeiro, 2013) (Munezero, 2014) (María Lucía Barrón-Estrada, 2017) (Ravi., 2015). For instance, (Munezero, 2014) analyzed learning journals to predict students' emotions and opinions about their learning experiences. (Kechaou, 2011) focused on extracting user opinions to enhance e-learning systems, using feature selection strategies like Mutual Information, Information Gain (IG), and CHI Statistics. These<sup>1</sup> approaches measure the difference between observed and expected frequencies of outcomes. Additionally, hybrid learning methods leverage Hidden Markov Models and Support Vector Machines.

Settings require tracking students' emotions and learning activities to understand their educational needs. The authors propose the Temporal Emotional Aspect Model (TEAM), which monitors emotions over time and provides two key outputs: (a) specific probabilistic aspect distributions for emotions and (b) their time-based development, revealing significant emotional patterns. The study also noted developmental trends (Mostafa., 2019). The findings indicate that: (i) content-related aspects received the most focus, with a higher likelihood of confusion compared to negative emotions; (ii) emotional fluctuations were more likely at the beginning and end of a semester; (iii) low-performing students showed lower emotional engagement and increased confusion towards the semester's end compared to high and medium-performing students.

---

<sup>1</sup>[https://www.investopedia.com/terms/c/chi-square-statistic.asp#:~:text=A%20chi%2Dsquare%20\(%20CF%872\)%20statistic%20is%20a%20measure.especially%20those%20nominal%20in%20nature.](https://www.investopedia.com/terms/c/chi-square-statistic.asp#:~:text=A%20chi%2Dsquare%20(%20CF%872)%20statistic%20is%20a%20measure.especially%20those%20nominal%20in%20nature.)

Existing research, such as (Zhi-Liu, 2019) analysis, predominantly focuses on small courses and fails to address the overall issue of student underperformance. These studies often predict overall student performance without proposing effective intervention methods. Machine Learning-based studies or those solely assessing student grades tend to fall short of achieving high accuracy or predefined goals.

In related literature, (Minaei-Bidgoli, 2003) explored student behavior in submitting assignments for an online physics course, using clustering and genetic algorithms. (Thai-Nghe, 2010) employed recommendation systems for predictions using publicly available data. Knowledge Centred Service and Online Community Service were used to tackle the problem, showing promise.

(Elbadrawy, 2016) applied linear regression algorithms and matrix solving to predict student outcomes. Moodle, a popular learning management system, has been utilized in various studies. (Konstantinidis, 2013) used Moodle logs for failure prediction and assessment of remedial measures. Some researchers have employed mathematical tools, including analytical devices within Moodle, for rapid data analysis.

(Kotsiantis, 2013) used Moodle to analyze data from 337 college students over three years, identifying variables correlated with final grades, including students' attitudes towards Moodle. (Oreški., 2018) used statistical tools like Analysis of Variance (ANOVA) and linear regression to examine the connection between online student interactions and their grades. (Cristóbal Romero, 2008), employed various machine learning techniques to predict overall student performance based on previous homework. (Abdullah., 2015) used decision trees to predict learning styles of 35 Moodle-based course students. (Pascual-Miguel, 2011) attempted to correlate Moodle results with academic achievement but found no significant correlation due to a limited dataset.

To address intervention challenges, both machine learning strategies and personal tutoring are proposed to achieve better results. Student interaction data is used to manage Moodle activity and assess project grades, primarily focusing on the first half of the course for timely interventions.

(Zhang, 2021) discuss Student Performance Prediction (SPP) as an essential aspect of personalized education, reviewing it from a machine learning perspective. They experimented

with a dataset from their institute, including 1,325 students and 832 courses. SPP benefits students in course selection and study planning, helps teachers adapt materials, identifies at-risk students, and aids education managers in curriculum optimization. The SPP process involves data collection, problem reformulation, model training, performance prediction, and result analysis, considering student and course characteristics.

Data collection includes student and course data, while machine learning models like the pass-fail model, assessment model, and outcome model are applied. Education modes such as offline, online, and blended classes are considered. This revised text provides a concise overview of the key points in the original text.

SPP benefits students in course selection, study planning, and helps teachers adapt materials and identify at-risk students (Ibrahim, 2007) (Bayer, 2012; Kloft, 2014). Education managers can use SPP to review the curriculum and improve educational performance planning (Reeves, 2018). This data-driven approach offers objective insights for the education system. The five steps to solve the problem are data pre-processing, feature selection, problem reformulation, model training, performance prediction, and results analysis.

- (1) The initial phase involves collecting data related to the Student Performance Prediction (SPP) system. This data includes a set of triplets (student, course, grade) that represent student scores. Student attributes encompass age, gender, health, economic status, educational level, and more. Similarly, course attributes involve parameters like frequency, duration, scope, opening season, etc., as referenced by (Elbadrawy, 2014; Kennedy, 2015; Barba, 2016) . Additionally, advanced student features can encompass parental characteristics, classmate group dynamics, and academic records, while course characteristics may include instructor qualities, introductory course details, assistant roles, and others. Three prevalent assessment models are the pass-fail model, assessment model, and outcome model. It's essential to note that this discussion emphasizes the need for clarity. Furthermore, educational contexts can be categorized as offline, online, or blended classes, as outlined by (Rovai, 2004).
- (2) After data preparation, the second step entails redefining problem 1. Broadly speaking, problem 1 can be reframed as involving grouping, classification, and regression tasks. The grouping aspect involves clustering data points into multiple clusters, each cluster

containing instances with high similarity. Numerous studies have segmented data into distinct groups based on SPP students and/or courses, as demonstrated by (Cakmak, 2017). The classification aspect is focused on predicting discrete scores, typically accomplished through machine learning classifiers such as Boolean regression (Elbadrawy, 2014) and Support Vector Machine (Xu, 2016).

The regression component involves forecasting continuous scores, often employing regression models such as linear regression (Alario-Hoyos, 2016) and neural networks (Oladokun, 2008). Furthermore, several studies have extended their findings to various levels of analysis, as highlighted by (Shahiri, 2015).

- (3) During the third stage, the chosen machine learning model is expanded to generate associations relevant to the reformulated problem. Many studies have traditionally employed machine learning techniques such as decision trees (DT) (Al-Radaideh, 2006); (Koprinska, 2015) , neighborhood methods (Meier, 2015) , linear regression (LR) (Anozie, 2006), neural networks (Andrews, 1995) and (Sorour, 2014), as well as kernel-based methods (Boser, 1992).

In the realm of Student Performance Prediction (SPP), novel feature learning approaches have been explored, including lasso regression (Sorour, 2014), (Zhang, 2018b); (Zhang, 2020), matrix factorization (MF) (Slim, 2014), tensor factorization (TF) (Thai-Nghe, 2010a), and deep neural networks (Kim, 2018). Among these methods, matrix factorization (MF) and deep learning have increasingly garnered attention within the SPP domain.

However, it's worth noting that sometimes simpler methods can yield more interpretable results compared to complex learning models, as highlighted by (Merriënboer, 2005).

- (4) In a trained model, the fourth step typically involves evaluating the performance of a new student in a new course, essentially inserting a new instance of {sp, cq} into the model (M) to obtain a prediction (yp, q). During this phase, current research explores various strategies. Some studies, such as those by (Al-Radaideh 2006), (Shovon, 2012), (Ahmed, 2014), (Meier, 2015), and (Al-Barrak, 2016), utilize training data to forecast a participating student's course grade. In contrast, the work by (Ren, 2016) focuses on predicting the assessment for the following semester based on school reports.

Additionally, several studies have examined post-course assessment throughout the learning period, as demonstrated by (Xu , 2017).

However, it's noteworthy that only a limited number of studies concentrate on the original Problem 1 (as defined at the outset of the process), emphasizing the learning model itself rather than specific attributes related to individual students or courses.

- (5) The expectation is that the model's outcomes will uncover interpretable patterns, aiding participants in enhancing their educational endeavors. In essence, Student Performance Prediction (SPP) offers insights into various aspects of teaching and learning, including identifying students at risk of attrition (Quadri, 2010), distinguishing students at varying knowledge levels (Meier, 2015), identifying critical learning factors (Mayilvaganan, 2014), and establishing relationships between courses (Zhang Y. A., 2021a). An effective means of enhancing performance can be the utilization of a grading system that anticipates student progress.

Researchers have employed advanced and sophisticated techniques to address the once seemingly impossible challenge of predicting students' outcomes from earlier weeks. For instance, Ioannis E (Livieris ,2019) utilized the J48 algorithm, achieving promising results by categorizing five classes based on eleven metrics. While these outcomes hold potential, our approach involves the application of fundamental algorithms on more limited metrics with the aim of attaining comparable results.

Furthermore, a deep learning Gridnet algorithm has demonstrated effectiveness when applied to specific datasets. However, it's worth noting that this algorithm typically requires a larger volume of data to yield optimal results. Our study, in contrast, relies on a dataset comprising only 1000 records, potentially introducing bias, as observed in the work conducted by (Ribeiro,2018).

In a related vein, the research conducted by (Zhi Liu ,2019) also centered around a smaller cohort, mirroring our current focus. However, Liu's study did not successfully address the issue of student underperformance. To mitigate these limitations, our research incorporates additional variables, such as dietary habits and parental involvement, with the intention of overcoming these challenges.

## CHAPTER 3 METHODOLOGY

### 3.0 Introduction

Methodology<sup>2</sup>, is the process of systematically solve the research problem. Within we do a study of different steps to carry out while studying the research problem by involving all logics behind those steps. It is also a discourse framework for research, a coherent and logical theme supporting views, beliefs, and values, that guides the various researchers [or alternative users] achieve the fixed goal of research. It includes the theoretical analysis of the body of strategies and principles involving a branch of knowledge such that the methodologies used from differing disciplines vary looking on their historical development.

### 3.1 Type of research used

Many forms of research methodology are often used in the classification of research. In our case, we adopt quantitative research though some important variables like feeding will be qualitative. Machine learning algorithm accuracy is better once the model is trained on significant data. Based on the information found, quantitative analysis is outlined as a style of data analysis using numerical assessments. The strategy will take the shape of performance metrics or data assessments employing a numerical scale, appreciating the rating of a variable on a scale of one to 10. Quantitative research offers exactness for responsive queries that need exact, verifiable responses. In support of our research objectives related to the evaluation of secondary school students' performance and the approach for predicting long-term performance, this form of analysis is particularly beneficial as it relies on quantitative data and statistical metrics. The results of the information that may be used will be of such a structure. To ensure a great outcome from the data used, python libraries like Pandas, SciPy, NumPy, and scikit-learn are used in machine learning prediction models and play essential roles in different aspects of the machine learning workflow. Here's a brief overview of the importance of each library:

Pandas:

Importance: Pandas is used for data manipulation and analysis. It provides data structures like DataFrames and Series, which are essential for storing and manipulating the dataset. You can use

---

<sup>2</sup> <https://southcampus.uok.edu.in/files/link/downloadlink/rm%20u1%20p1.pdf> accessed 28/07/2023 06:17:00

it to read data from various sources, perform data preprocessing, and handle missing values and categorical data.

NumPy:

Importance: NumPy is a fundamental library for numerical computations in Python. It provides support for multidimensional arrays and mathematical functions. It is crucial for performing operations on data, such as array manipulation, linear algebra, and statistical analysis.

Seaborn and Matplotlib.pyplot:

Importance: Seaborn and Matplotlib are data visualization libraries. Seaborn is built on top of Matplotlib and provides a higher-level interface for creating attractive and informative statistical graphics. Data visualization is crucial for understanding data patterns and relationships, which is an important step in any machine learning project.

Scikit-learn:

Importance: Scikit-learn is one of the most widely used machine learning libraries in Python. It offers a comprehensive suite of tools for building, training, and evaluating machine learning models. This library provides a wide range of algorithms for classification, regression, clustering, and more, along with tools for feature selection, model selection, and evaluation.

Dabl:

Importance: Dabl, short for "Data Analysis Baseline Library," is a library that automates many aspects of the data analysis and machine learning process. It helps with data preprocessing, feature selection, model selection, and hyperparameter tuning. Dabl can be especially useful for quick model prototyping and understanding the baseline performance of various machine learning models.

Plotly:

Importance: Plotly is a library for creating interactive, web-based visualizations. It's particularly useful for creating interactive plots and dashboards, which can be valuable for presenting and

sharing the results of your machine learning models in a user-friendly manner. It's often used for creating dynamic and customizable charts and graphs.

In summary, these libraries are essential components of a machine learning prediction model. Pandas and NumPy are used for data manipulation and numerical operations, Seaborn and Matplotlib for data visualization, Scikit-learn for building and evaluating machine learning models, Dabl for automating aspects of the machine learning process, and Plotly for creating interactive visualizations. When used together, these libraries provide a comprehensive toolkit for developing and deploying machine learning models.

### 3.2 Data collection method

Data collection techniques like survey interview and focus group was better choice for this research to gather all required data. The described challenge in next paragraph changed the status then forced me to used data already collected by other researchers [Kaggle](#)<sup>3</sup>

#### 3.2.1 Data preparation

Here we are using the dataset from [Kaggle](#)<sup>4</sup> . The data contain the columns of Gender, Race/ethnicity, Parental level of education, Lunch, test preparation course, reading score, writing score, and Maths score.

	gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score
0	female	group B	bachelor's degree	standard	none	72	72	74
1	female	group C	some college	standard	completed	69	90	88
2	female	group B	master's degree	standard	none	90	95	93
3	male	group A	associate's degree	free/reduced	none	47	57	44
4	male	group C	some college	standard	none	76	78	75

Because the dataset is from Kaggle, based to the research's purpose, some columns, like ethnicity, are not required for our purpose. And some columns have been edited like writing score to JavaScript, then writing score to VueJS. This is to have a proper meaning of the results on Rwandan education. Additionally, we redefined the means of lunch that has the quality and reduced; in our case, we are taking the daily lunch as balanced and regular diet feeding.

---

<sup>3</sup> <https://www.kaggle.com/gauravarena/students-performance-analysis-and-prediction/data> accessed 28/07/2023 06:17:00

<sup>4</sup> <https://www.kaggle.com/gauravarena/students-performance-analysis-and-prediction/data> accessed 05/09/2023 08:51:22

	gender	parent_education	lunch	prep	math	JavaScript	VueJS
0	female	bachelor's degree	standard	none	72	72	74
1	female	some college	standard	completed	69	90	88
2	female	master's degree	standard	none	90	95	93
3	male	associate's degree	free/reduced	none	47	57	44
4	male	some college	standard	none	76	78	75

### 3.2.2 Data cleaning

Data analysis approaches also are vital in methodology. That is why we involve checking missing data in our dataset, invalid records, outliers, like records that by typing error would exceed 100 marks. At the same time, the maximum is fixed at 100, and any other irrelevant information would lead to the analysis's false results.

```
no_of_columns = data.shape[0]
percentage_of_missing_data = data.isnull().sum()/no_of_columns
print(percentage_of_missing_data)
```

```
gender          0.0
race/ethnicity  0.0
parent_education 0.0
lunch           0.0
prep            0.0
math            0.0
JavaScript      0.0
VueJS           0.0
dtype: float64
```

### 3.2.3 Model selection

4 types of models are selected for their accuracy on relatively small dataset like the one used in this research. namely Linear Regression, Support Vector Regression, and K Nearest Neighbor.

#### 3.2.3.1 Typecasting (label encoding)

The used models in this research are not for natural language processing. Because of the independent variable used here are of string type, to facilitate their use in model processing for prediction a typecast is done to convert them into numerical values.

## Label Encoding

```
from sklearn.preprocessing import LabelEncoder

# creating an encoder
le = LabelEncoder()

# Label encoding for test preparation course
data['prep'] = le.fit_transform(data['prep'])

# Label encoding for Lunch
data['lunch'] = le.fit_transform(data['lunch'])

# Label encoding for parental Level of education
data['parent_education'] = le.fit_transform(data['parent_education'])

#Label encoding for gender
data['gender'] = le.fit_transform(data['gender'])
```

### 3.2.4 Model processing

The strategies employed for organizing and comparing diverse results, derived from the analysis of data, which may involve one or multiple columns, have proven instrumental in facilitating informed decision-making. These strategies leverage statistical outcomes generated from various parameters to contribute to the process of drawing conclusions.

OLS Regression Results

Dep. Variable:	average_score	R-squared:	0.210
Model:	OLS	Adj. R-squared:	0.201
Method:	Least Squares	F-statistic:	22.98
Date:	Wed, 14 Jun 2023	Prob (F-statistic):	2.66e-31
Time:	16:26:26	Log-Likelihood:	-2781.9
No. Observations:	700	AIC:	5582.
Df Residuals:	691	BIC:	5623.
Df Model:	8		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	70.5792	1.468	48.094	0.000	67.698	73.461
gender_male	-3.5236	0.984	-3.582	0.000	-5.455	-1.592
parental_bachelor	1.6121	1.758	0.917	0.359	-1.839	5.063
parental_hs	-5.4357	1.593	-3.413	0.001	-8.563	-2.309
parental_masters	4.3183	2.261	1.910	0.057	-0.121	8.757
parental_somecol	-1.5724	1.443	-1.090	0.276	-4.405	1.260
parental_somehs	-4.9937	1.538	-3.246	0.001	-8.014	-1.973
lunch_standard	9.1581	1.027	8.915	0.000	7.141	11.175
test_prep_none	-7.8748	1.027	-7.668	0.000	-9.891	-5.858

Omnibus:	6.193	Durbin-Watson:	2.021
Prob(Omnibus):	0.045	Jarque-Bera (JB):	6.293
Skew:	-0.230	Prob(JB):	0.0430
Kurtosis:	2.929	Cond. No.	9.13

```
plt.plot(final_rmse.k, final_rmse.RMSE)
plt.plot(final_rmse.k[min_index], final_rmse.RMSE[min_index], 'ro')
plt.title('Plot of Test RMSE vs K number of neighbours')
plt.xlabel('K')
plt.ylabel('Test RMSE')
plt.show()
```

K N N

```

model_linear = SVR(kernel="linear")
model_linear.fit(train_X, train_Y)
pred_linear = model_linear.predict(test_X)
linear_rmse = sqrt(mean_squared_error(test_Y,pred_linear))

# kernel = poly
model_poly = SVR(kernel="poly")
model_poly.fit(train_X, train_Y)
pred_poly = model_poly.predict(test_X)
poly_rmse = sqrt(mean_squared_error(test_Y,pred_poly))

# kernel = sigmoid
model_sigmoid = SVR(kernel="sigmoid")
model_sigmoid.fit(train_X, train_Y)
pred_sigmoid = model_sigmoid.predict(test_X)
sigmoid_rmse = sqrt(mean_squared_error(test_Y,pred_sigmoid))

# kernel = rbf
model_rbf = SVR(kernel="rbf")
model_rbf.fit(train_X, train_Y)
pred_rbf = model_rbf.predict(test_X)
rbf_rmse = sqrt(mean_squared_error(test_Y,pred_rbf))

data = {"kernel":pd.Series(["linear","polynomial","sigmoid","rbf"]),
        "Test RMSE":pd.Series([linear_rmse,poly_rmse,sigmoid_rmse,rbf_rmse]),
        "Pred":pd.Series([pred_linear,pred_poly,pred_sigmoid,pred_rbf])}
table_rmse=pd.DataFrame(data)
table_rmse

```

SVR

```

import tensorflow as tf

# Importing necessary models for implementation of ANN
from keras.models import Sequential
from keras.layers import Dense

cont_model = Sequential()
cont_model.add(Dense(100, input_dim=train_X.columns.value_counts().sum(), activation="softmax"))
cont_model.add(Dense(60, activation="relu"))
cont_model.add(Dense(1, kernel_initializer="normal"))
cont_model.compile(loss="mean_squared_error", optimizer = "adam", metrics = ["mse"])

model = cont_model
model.fit(np.array(train_X), np.array(train_Y), epochs=300)

# On Test dataset
pred = model.predict(np.array(test_X))
pred = pd.Series([i[0] for i in pred])

```

ANN

The flow of data processing in machine learning

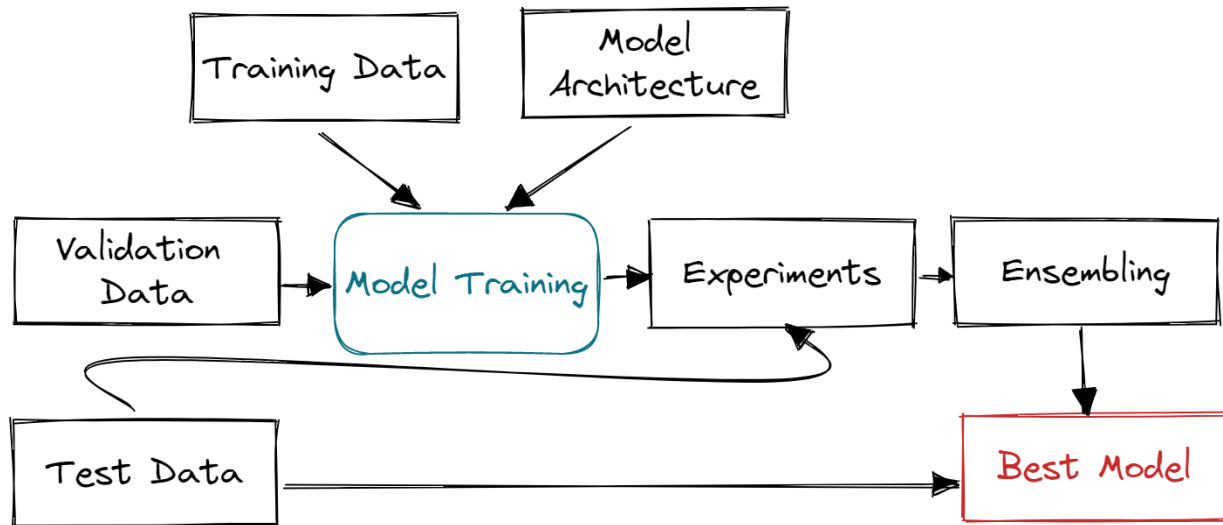


Figure 1: The flow of data processing in machine learning

source<sup>5</sup> Image by Abid Ali Awan data scientist and blogger

### 3.3 Resources, materials and tools

This research wants powerful computing equipment. Particularly, computing processors ought to be a minimum of 3.5 GHz (gigacycle per second), quadcore with a substantial throughput.

A RAM of sixteen GB (Giga Bytes) or higher is desirable due to the amount of data to be processed.

Moreover, the internet high speed bundle is preferable if it is optic fiber. This facilitates data exchange among different components of the system and xAPI. Alternatively, we will use 4G (Fourth Generation) wireless internet. Big amount of computer codes is required to perform this task: Here we used anaconda as one of the software to perform the task. Being a set of various software, it's essential to use it because it contains plenty of important libraries that are used and their importance of use is detailed on page 20& 21. Here a python Jupyter notebook which includes all imported libraries (see figure 17 page 36) is used to run codes and see the output on

<sup>5</sup> <https://www.datacamp.com/blog/machine-learning-lifecycle-explained> accessed 28/07/2023 06:17:00

the interface. Dependent variables used to identify the impact of each of them, are known from the dataset as justified in the chapter 4 about system analysis and design on page 23.

### 3.4 Justification behind the analysis

The importance of this research is the anticipation of students' results. Education is one the foremost vital part of life. Success in life mostly depends on success in studies. The ultimate result of the report does not explain why the student got such results. This research aims to point out one of the ways to anticipate the student's performance and then take useful measures that are coming back as part of the studies supported results from the last part of the term. This prediction is often made double a time or regularly based on how frequently the student is evaluated. The supply of dataset we tend to utilize in this analysis is from Kaggle, as explained from page 17 (data collection).

To answer our research question about how to evaluate the students' performance, different comparisons are done. To evaluate the impact of each parameter in our dataset, the role of gender on various subjects, about the parents' education background, and the regular standard and balanced diet. Finally, a model is trained to predict how the students will perform given the factors influencing their performance.

The model uses essential functions and libraries (see figure 17 page 36): The contents [available](#)<sup>6</sup> ,give clear understanding about libraries that are important in work depending on their functionalities from the next paragraph. For example, Pandas is an open-source library that gives superior information manipulation in Python. It's built on top of the NumPy package (one of the famous packages in machine learning), which implies that NumPy is usually needed for operating the Pandas. The name Pandas comes from the word Panel Data, which means a political economy from three-dimensional data. It's used for data analysis in Python and was developed by Wes McKinney in 2008.

1. **Pandas**<sup>7</sup> is one of python package which is used to provide fast, flexible and expressive data structures. Its design allows to work with relational data in an easy and intuitive

---

<sup>6</sup> <https://www.javatpoint.com/pandas-vs-numpy> accessed 01/09/2023 08:50:00

<sup>7</sup> [https://pandas.pydata.org/docs/getting\\_started/overview.html](https://pandas.pydata.org/docs/getting_started/overview.html) accessed 02/09/2023 10:50:10

way. As purpose of being a high-level building block for the real-world data analysis in python.

2. **NumPy**<sup>8</sup> it stands for Numerical Python. This library is mainly for processing complex mathematical tasks like arrays in our case. It has the ability to process a multidimension arrays and handle its data easily using a relatively low computational capacity. is generally written in C language, it is free accessible thanks to the work done by Travis Oliphant in 2005.
3. **Seaborn**<sup>9</sup> is a plotting library that provides a more accessible interface, smart defaults for plots required for machine learning, most importantly, the plots are highly sophisticated than those in Matplotlib.
4. **Matplotlib.pyplot**<sup>10</sup> is a plotting library used for 2D graphics in a python programming language. It is utilized in python scripts, shells, internet application servers, and different graphical program tool-kits. By using it, we can generate plots, histograms, power spectra, bar charts, error charts, scatterplots, etc with simply a couple lines of code.
5. **Scikit-learn**<sup>11</sup> is a library in Python that gives several unattended and supervised learning algorithms. It's engineered upon a number of the technologies that may be already acquainted with most machine learning users, like NumPy, pandas, and Matplotlib.

Practicality scikit-learn provides the following:

- & Regression, including Linear and Logistic Regression
- Classification, including K-Nearest Neighbours
- Clustering, including K-Means and K-Means++
- Model selection
- Pre-processing, including Min-Max Normalization

information found<sup>12</sup> by Satyam Kumar(August 5<sup>th</sup> 2021) shows the good practice of using scikit-learn library for machine learning

---

<sup>8</sup> [https://www.w3schools.com/python/numpy/numpy\\_intro.asp](https://www.w3schools.com/python/numpy/numpy_intro.asp) accessed 05/08/2023 08:53:05

<sup>9</sup> <https://seaborn.pydata.org/tutorial/introduction> accessed 28/07/2023 06:17:00

<sup>10</sup> <https://www.goeduhub.com/639/what-is-matplotlib-in-data-science> accessed 11/08/2023 11:57:49

<sup>11</sup> <https://www.codecademy.com/articles/scikit-learn> accessed 11/08/2023 13:35:25

<sup>12</sup> <https://towardsdatascience.com/a-python-tool-for-data-processing-analysis-and-ml-automation-in-a-few-lines-of-code-da04b3ba904f> accessed 28/07/2023 06:17:00

6. **Dabl** is a Python tool for knowledge Processing, Analysis, and millilitre Automation in some strains of code. For starters, it is a basic information analysis library that makes it easier and more convenient to model a supervised system for both novices and experienced and does not employ an understanding of science knowledge. Dabl is stirred by suggestions of the Scikit-learn library and tries to adjust system mastering modelling by decreasing the standardized procedures that may be many times without making major changes to the original status by automating the components.
7. **Plotly** is one of the libraries used for higher plotting. Referring to information<sup>13</sup>, plotly is a python open-source graphing library. It's accustomed to create interactive, publication-fine graphs. Thanks, to it we can make line plots, scatter plots, bar charts, error bars, field plots, histograms, heatmaps, subplots, multiple-axes, polar charts, and bubble charts.
8. **Label** writing in coder is one of the elements of the scikit-learn library. It's used to encode the Data Frame of string labels. Suppose that the peak column consists of tall, medium, and short. Then applying label encoding, the peak column is reborn into: 0, that represents tall, 1 which represents medium size 2, which represents short. The machine learning process needs the potency of information. First and foremost, we check the info we are aiming to use. Particularly regarding missing data which can negatively impact the results. Once we describe the data, we've got within the data set, sort of a summary, we have the following results. According to the Oxford wordbook, a dataset is "an assortment of data that is treated as one unit by a computer" .Oxford dictionary<sup>14</sup>

---

<sup>13</sup> <https://plotly.com/python/> accessed 05/09/2023 08:53:44

<sup>14</sup> <https://www.oxfordlearnersdictionaries.com/definition/english/data-set?q=data+set> accessed 04/08/2023 09:23:57

## CHAPTER 4 SYSTEM DESIGN AND STATISTICAL ANALYSIS

This chapter presents the analysis and design of the students' performance prediction system. The system architecture, and system flowchart.

This chapter presents the analysis of the complete system. In this research, we used a combination of Articulate Storyline<sup>15</sup>, xAPI<sup>16</sup> (which makes the essential part of a teacher to prepare contents, assessments, quizzes, tracking the student's behaviour and their journey into the e-learning system), and the Anaconda<sup>17</sup> software mainly Jupyter notebook. The collected data are used as a dataset for the AI model to predict the student's future results, then allow the school administrator to take the appropriate measures in an anticipated way. Prepared data from the storyline and Veracity Learning Record<sup>18</sup> store is published on Moodle for student performance detection.

---

<sup>15</sup><https://www.voices.com/blog/what-is-articulate-storyline/#:~:text=Articulate%20Storyline%20is%20the%20ideal,laptops%2C%20tablets%2C%20and%20smartphones>. Accessed 05/09/2023 08:54:53

<sup>16</sup><https://xapi.com/overview/> accessed 05/09/2023 09:35:07

<sup>17</sup><https://docs.anaconda.com/free/navigator/index.html> accessed

<sup>18</sup><https://lrs.io/> accessed

## 4.1 System architecture

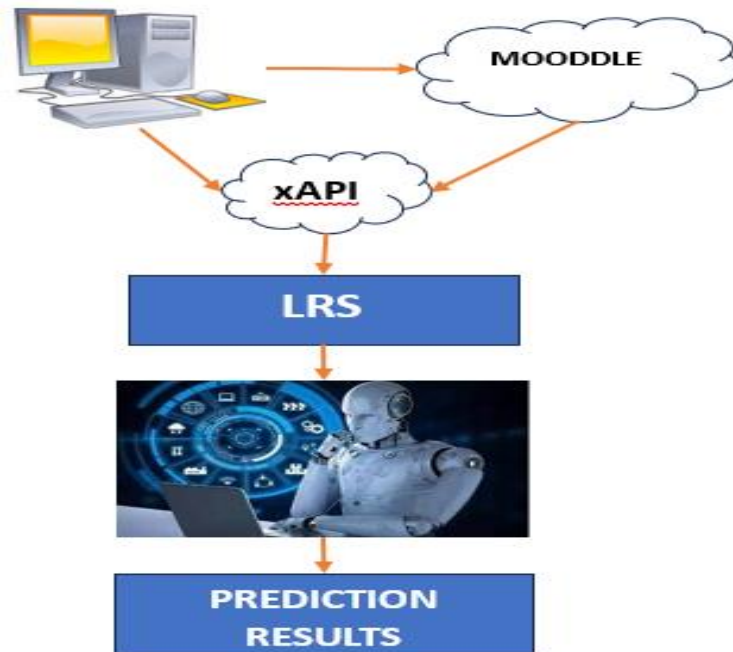


Figure 2: System architecture design

The system is mainly composed by two main parts working complementary to each other: the local part mainly focuses on configuration and observation of the results. The second part work online mainly for data processing including Moddle classes and xAPI tracking activities. LRS (Learning Records Store works as a storage of results from xAPI and Moddle and providing the input dataset to the AI prediction models. From them we have prediction results which we have to analyse for their accuracy then make adjustment if necessary to improve it and rely on the prediction for decision making about how to improve the student performance. The figure 2 shows the system architecture

### 4.1.1 System configuration

This activity is conducted on the computer. Here we prepared the configuration of articulate storyline software which we used to prepared and publish courses to Moddle e-learning system

Here we are using storyline <sup>19</sup> to prepare the contents. The choice made because we found it helpful to prepare and connect it to the xAPI, the same as publishing the contents to Moodle.

Some industries employ articulate storyline software for making interactive courses as they notice it is easy for beginners but powerful enough for experts. And it lets a user create just about any thinkable interaction in minutes. In addition, its presentations have the ability to automatically adapt to the screens of different devices, such as laptops, tablets, and smartphones. This is also helping based on the user who is going to access the content published from the Storyline.

Depending on the user preference also there exist its alternative like iSpring Suite, Adobe Captivate, Brainshark, Easygenerator, Lectora, Lessonly by eismic, <sup>20</sup>Elucidat. or 360Learning.. Contents from Storyline once published to Moodle, was accessed by xAPI which in return processed them to give an overview of the student's journey. Here notice that we have a direct access to the xAPI configuration where we can see and add more parameters for the prediction purposes. From the processed data, the dataset as collected saved to LRS then loaded into an Artificial Intelligence Model.

Student results from Moodle also can be edited and add extra columns if needed for the accuracy of predicted results (example the type of food a student takes). The tracked students' behavior from xAPI can help partially in advising student about his/her activities in the system (example how often he/she visit the system, how much time he/she spend answering the quiz, ...)

## **4.2 Articulate Storyline**

One of the most challenging tasks of daily teaching activities is course preparation, the time a teacher spends creating the course contents. The Storyline<sup>21</sup> by Keaton Robbins (May 11<sup>th</sup> 2022) is one type of software that can be used to prepare interactive contents that can attract students to spend time. Courses qualified in the form of games or tutorials with step-by-step guide to be completed; they generate the students' curiosity and do all possible to be satisfied. During this

---

<sup>19</sup><https://www.voices.com/blog/what-is-articulate-storyline/#:~:text=Articulate%20Storyline%20is%20the%20ideal,laptops%2C%20tablets%2C%20and%20smartphones.> 04/09/2023 11:13:40

<sup>20</sup> <https://www.g2.com/products/articulate-360/competitors/alternatives> 04/09/2023 08:20:23

<sup>21</sup> [What Is Articulate Storyline? | Voices | Voices](#)

quest, students learn and do exercises with enthusiasm. Some industries employ articulate storyline software for making interactive courses as they notice it is easy for beginners but powerful enough for experts. And it lets a user create just about any thinkable interaction in minutes. With Storyline, we can build slide-based lessons that mix instruction, audio, video, and interactions to form participating online courses. Each slide could be a blank canvas for organizing contents into scenes: once users interact with totally different parts on the decline, certain actions trigger responses (see Figure 7 page 34)

### 4.3 xAPI in combination with Articulate Storyline

The experience Application Interface (or xAPI) is a specification for learning technology that produces the potential to gather knowledge regarding the extensive selection of experiences someone has (online and offline). This API captures data consistently from a few people or groups’ activities from several technologies. Different systems are ready to communicate by capturing and sharing this stream of activities using xAPI’s vocabulary.



source: here<sup>22</sup>.

Figure 3: xAPI features (by  
Rahul Aug 27<sup>th</sup>, 2021)y6

A JavaScript file containing instruction codes is required to enable the communication between the storyline, xAPI, and Veracity Learning Record store<sup>23</sup>.

<sup>22</sup> <https://tinyurl.com/2ab2k97x> accessed 11/10/2023 20:46:03

<sup>23</sup> <https://lrs.io/> accessed 05/09/2023 08:44:28

### 4.3.1 Experience Application Interface(xAPI) required statement

XAPI has many reports representing a JSON object, but the following must be present, while others can be added if necessary.

- a. **Actor** as a JSON object, it represents the doer of activity. This object may include other things like the actor's last name, emails, etc. Defining the actor is better to follow the common definition standards provided by adlnet, the owner of the xAPI project available here <sup>24</sup>
- b. **Verb** Represents the action that the actor has done. We can express the action using different verbs but to mean the same thing, and we have to avoid that as much as we can, and that is why we use common standardized verbs to mean actions from the official xapi community website<sup>25</sup> which provides a long list of verbs (but currently, it may be deprecated), alternatively the agreed verb of the activity by the community of XAPI users can be found on. website<sup>26</sup>.

### 4.3.2 Collecting username and emails of a user from the storyline

To identify the user of a particular activity, we set up a form that the user must fill out and then collect his name and email. This last one we use as a unique identifier of a user in the system.

Variables defined in the Storyline are accessed using the built-in JavaScript function from the Storyline, which is called GetPlayer (). After calling it, we can now access other parts like setVar () and getVar(), which allow us to manipulate the value contained by the variable.

---

<sup>24</sup> xAPI website <https://www.xapi.com> accessed 9/5/2023 8:49:40 AM

<sup>25</sup> <https://www.registry.tincanapi.com> accessed 9/5/2023 8:55:06 AM

<sup>26</sup> <http://www.xapi.vocab.pub> accessed 9/5/2023 9:02:27 AM

```
1  const player =GetPlayer();
2  const uNamejs = player.GetVar("uName");
3  const uEmailsjs = player.GetVar("uEmail");
4
5  const statement = {
6    "actor": {
7      "name": uNamejs,
8      "mbox" : "mailto:" + uEmailsjs
9    },
10   "verb": {
11     "id": "http://adlnet.gov/exapi/verbs/completed",
12     "display": { "en-us": "completed" }
13   },
14   "object": {
15     "id": "https://www.example.com/write-xapi-contents",
16     "definition": {
17       "name": { "en-us": "Write xAPI Statements" }
18     }
19   }
20 }
```

**Figure 4: Script for collecting usernames and emails from Storyline**

### 4.3.3 Get and send statements from the storyline

We are sending statements using variables defined in the storyline and connected to the JavaScript file uName variable used to get a user name and uEmail variable to get our user emails. The specified function sendStatement() enclose the whole objects to be used to send data from the storyline.

To specify where we send our data, we need to configure the connectivity surrounding the conf object (see Figure 5 page 32). This is helped by a JavaScript wrapper which is provided by the owner of this project to facilitate the connectivity. JavaScript codes that allow us to send statements to our Learning Record store are depicted in the following image. (see Figure 5)

```
const player =GetPlayer();
const uNamejs = player.GetVar("uName");
const uEmailsjs = player.GetVar("uEmail");

const conf = {
  "endpoint": "https://marcel-thesis.lrs.io/xapi/",
  "auth": "Basic " + toBase64("zamhas:wicpir")
}

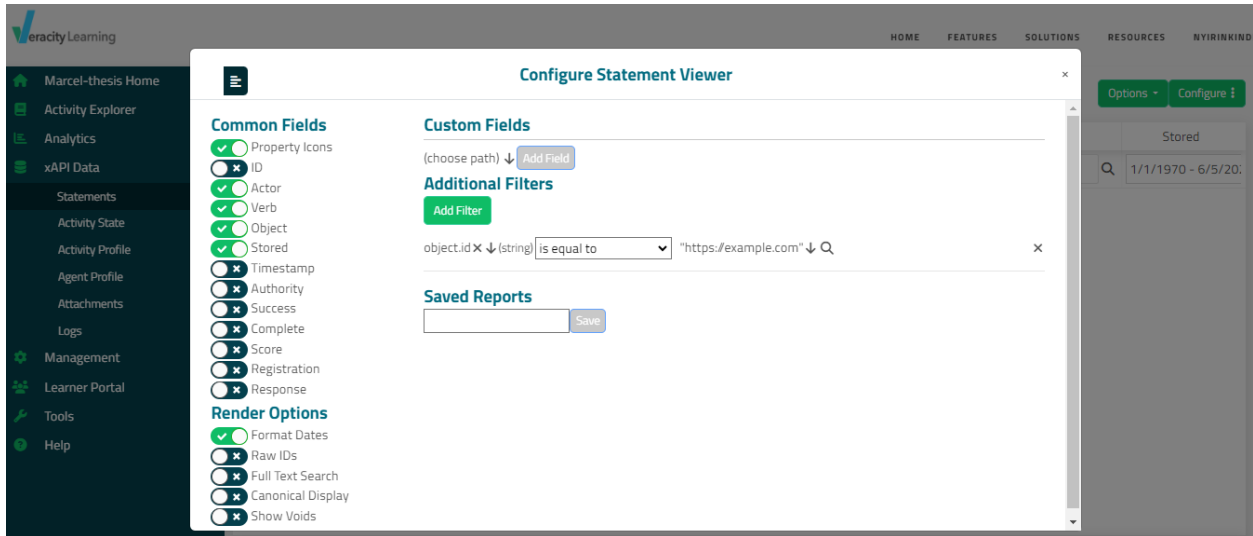
ADL.XAPIWrapper.changeConfig(conf);

const statement = {
  "actor": {
    "name": uNamejs,
    "mbox" : "mailto:" + uEmailsjs
  },
  "verb": {
    "id": verbId,
    "display": {"en-us": verb}
  },
  "object": {
    "id": objectId,
    "definition": {
      "name": {"en-us":object}
    }
  }
}

const result = ADL.XAPIWrapper.sendStatement(statement);
```

**Figure 5: Script for sending a statement to Veracity**

After creating an account on Veracity, we name our LRS the unique name and get our URL as our endpoint. veracity account. We need to be authenticated, create the username and password in the Veracity account, and configure different parameters (see Figure 7) based on what we need to send and receive.



**Figure 6: xAPI parameters configuration in Veracity (lrs.io)<sup>27</sup>**

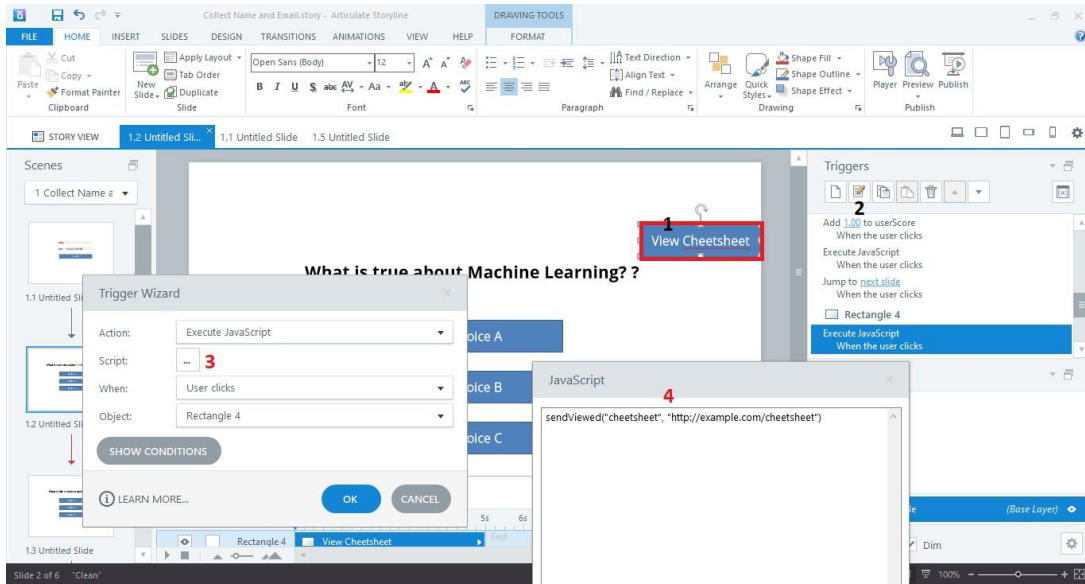
#### 4.3.4 Configure storyline's triggers

The process of configuring what should happen once a slide is open or if a specific button is manipulated is done in the same way depending on what you are preparing to do, either a quiz, a multiple-choice answer, content to be viewed and so on. . . In the following image of the storyline, we are configuring the trigger for viewing the cheatsheet:

1. Click on the View Cheatsheet button
2. go to the triggers tab and click on the edit triggers button
3. On the opened window, click on the script to open the script windows

on that window, we write the function that must be executed once we click on that button of ViewCheatsheet, and that information is directly sent and recorded to our LRS (Learning Record Store).

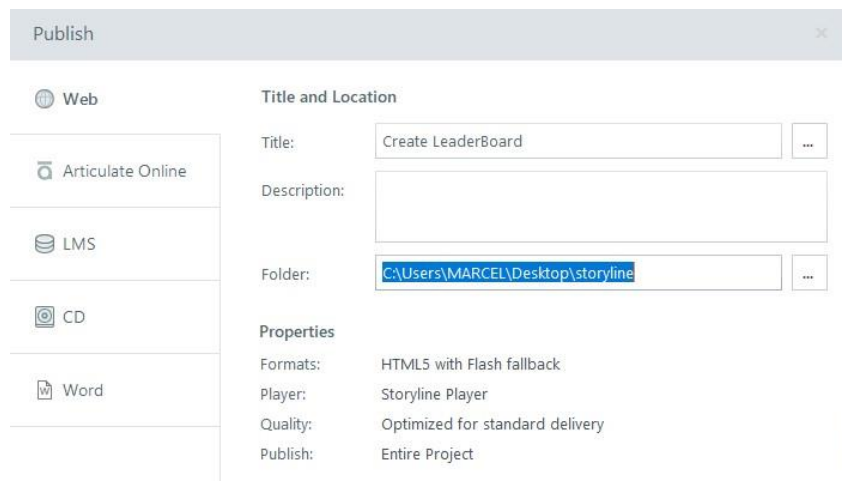
<sup>27</sup><https://tinyurl.com/5cd9kctp> accessed 05/09/2023 08:19:19



**Figure 7: Configuration of Storyline triggers**

Once configured, it's time to publish our work. We do that by selecting one of several options in the storyline, as shown in the image below (Figure 9), either to web if we have a website, to articulate online, to a Learning Management System, to a Compact Disk, or to a word document depending to what we need to do.

Once everything is set up, we can send our statements over the LRS from the storyline's published contents. And now, we can view the user's entire journey in veracity LRS.



**Figure 8: Publish contents using Storyline by choosing one of the available options**

### 4.3.5 Collecting xAPI activities

Figure 10 shows what to be added to the JavaScript code to send the activities done by the student over the Learning Record Store (LRS) which can be monitored for performance.

```
"object": {
  "id": objectId,
  "definition": {
    "name": {"en-us": object},
    "description": {"en-us": objectDescription},
    "type": activityType
  },
  "objectType": "Activity"
},
function sendStatement(verb, verbId, object, objectId, objectDescription, activityType)
```

**Figure 9: Script for collecting activities from the storyline**

### 4.3.6 Collecting open-text answers from the storyline

Among questions that will be asked there can be multiple choices and open-ended questions. In that case, students will be required to answer depending on some logic behind the question. The JavaScript file must be modified to collect open-text answers, as shown in the following images (figure 11 and Figure 12). By the way, the storyline slide configuration must be done accordingly.

**Note:** To avoid conflicts, it's a must to ensure that the variable from the storyline and JavaScript don't have the same names.

```
1 function sendStatement(verb, verbId, object, objectId, objectDescription, activityType, openTextVar {
2
3   const player = GetPlayer();
4   const uNamejs = player.GetVar("uName");
5   const uEmailsis = player.GetVar("uEmail");
6   const userResponse = player.GetVar(openTextVar);
7
```

**Figure 10: Script for collecting open answers**

```
    "objectType": "Activity"
  },
  "result": {
    "response": userResponse
  }
}
```

**Figure 11: Script for sending open text answers**

#### 4.3.7 Collecting quiz score information

Based on the score recorded in our Learning Record Store (LRS) we can do some preliminary overview of the student progress before we load the results into the machine learning model to give precise statistical analysis and predictions. At the end of this step, we can view in our LRS the score obtained from the quiz done by a student, the same as all information we configured till now to be considered within LRS.

```
    "result": {
      "response": userResponse,
      "score": {
        "min": 0, //minimum score
        "max": maxScorejs, // maximum score
        "raw": userScorejs, //maximum number of questions answered correctly
        "scaled": scaledScore //% of score
      },
      "success": success
    }
  }
}
```

**Figure 12: Script for collecting quiz score**

#### 4.3.8 Measuring the duration a student spent on a certain activity

This step aims to track the time a student spent answering each question and the time spent on a slide. This can help to improve students' user experience, as the identification of where they get trapped during their journey in the system can be monitored. Later on, it will be helpful to determine the student's performance, identify the cause then predict the student's future progress. At the top of the JavaScript file, some global variables are added that help to track the timing.

```

1
2  var courseSeconds = 0;
3  var slideSeconds =0;
4
5  var isCourseTimerActive =false;
6  var isSlideTimerActive =false;
7  window.setInterval(()=>{
8      if (isCourseTimerActive === true){
9          courseSeconds += 1
10         }
11         if (isSlideTimerActive === true) {
12             slideSeconds +=1
13         }
14     },1000)
15     const managerTimer = {
16         "course": {
17             "start": () => { isCourseTimerActive = true} ,
18             "stop": () => { isCourseTimerActive = false} ,
19             "reset": () => {courseSeconds = 0}
20         },
21         "slide": {
22             "start": () => {isSlideTimerActive = true} ,
23             "stop": () => {isSlideTimerActive = false} ,
24             "reset": () => {slideSeconds = 0}
25         }
26     }

```

**Figure 13: Script for collecting the time spent on the activity**

Here notice that xAPI has a specific way of recording the time format, ISO 8601, in the form of PT (Period of Time) to indicate the period. For example, if a student spends 3 hours, 20 minutes, and 15 seconds in the system, it will be recorded in the following format PT3H20M15S.

The JavaScript file is accessing the ConvertToIso function using the defined variable finalDuration, which is used to distinguish which time we are accessing for course time or slide time. Then we send the duration over the LRS using the duration statement, where we operate “duration” as a child of the “result” object statement.

```

85 //converting time into ISO 8601
86 function convertToIso(secondsVar) {
87     let seconds = secondsVar;
88     if (seconds > 60) {
89         if (seconds > 3600) {
90             const hours = Math.floor(seconds / 3600);
91             const minutes = Math.floor((seconds % 3600) / 60);
92             seconds = (seconds % 3600) % 60;
93             return `PT${hours}H${minutes}M${seconds}S`;
94         } else {
95             const minutes = Math.floor(seconds / 60);
96             seconds %= 60;
97             return `PT${minutes}M${seconds}S`;
98         }
99     } else {
100         return `PT${seconds}S`;
101     }
102 }

```

**Figure 14: Script to convert to ISO time**

NB images of this script are from this [source](#)<sup>28</sup>

```

38 let finalDuration;
39 if (timer == "course") {
40     finalDuration = convertToIso(courseSeconds)
41 } else if (timer == "slide"){
42     finalDuration = convertToIso(slideSeconds)
43 } else {
44     finalDuration = null
45 }

```

**Figure 15: Send the final time spent in the system**

#### 4.3.9 Data visualization in xAPI

Stored data by LRS (Learning Record Store) collected from xAPI can be visualized via xAPI Dashboards, where we have comprehensive analytics about the student journey into the system. Using this, we can have an overview of the student's journey and know about his weakness and which part of the learning journey we can focus on. From there, we can use Artificial Intelligence models to predict future performance.

<sup>28</sup> <https://www.devlinpeck.com/tags/xapi> accessed 05/09/2023 09:46:46

## 4.4 System Flowchart

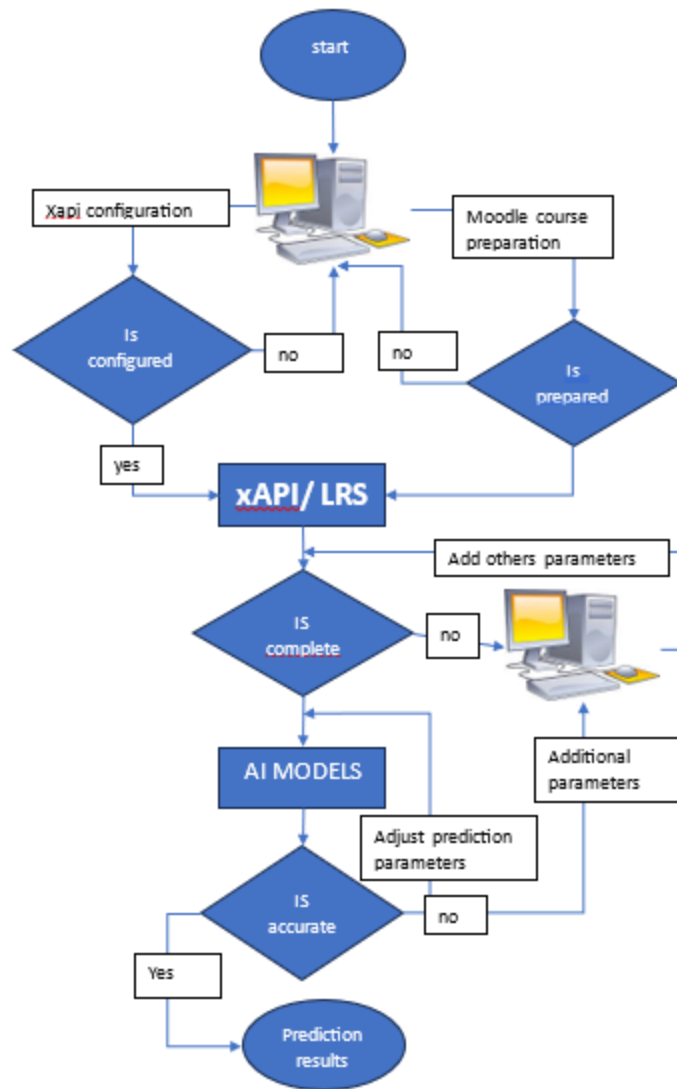


Figure 16: System flowchart

## 4.5 Statistical analysis

To make predictions, we need to have a dataset. Here we are using a dataset from [Kaggle](#). Kaggle is an online platform where data scientists can compete on different machine learning projects. For those purposes, it has a lot of datasets that can be used for data science testing projects for free. The structure of the dataset is explained in chapter 3.2 data collection see page 17.

For this thesis, this dataset has been chosen because it has a lot of similarities to Rwandan education purposes to increase student performance. Here it is used to predict the student's performance. From the results, measures can be taken in an anticipated way to help a student to improve before his/her final results.

In this dataset, there is a population of 1000 students, and all necessary operation to assure a good result have been done (like checking for bias, and invalid data which is done in the process of data preparation of machine learning). From this population, automatically in a randomized way, 20% is used to test and 80% for training the model on the dataset being used. The accuracy of results increases as the population in the dataset increases. Since 20% is used for testing, if the testing sample increases, the error of predicting erroneous results reduces, then the accuracy increases.

To analyse statistically, libraries that are included in Anaconda software (see figure 17 page 42) are used (for example plotly. Matplotlib...).

#### **4.5.1 Importing necessary libraries**

To be able to make an analysis, build a model and make a prediction, there are important libraries that need to be imported (see Chapter 3.4 page 19 and 20). The following figure 17 shows the enumerated list of used libraries

```

import pandas as pd
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
import matplotlib.ticker as mtick
from colorama import Fore, Back, Style

import plotly
import plotly.express as px
import plotly.graph_objects as go
from plotly.offline import init_notebook_mode, iplot
init_notebook_mode(connected=True)

import squarify
from collections import Counter

import tensorflow as tf
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.optimizers import SGD

from sklearn.model_selection import train_test_split
from sklearn.metrics import *

import warnings
warnings.filterwarnings("ignore")

```

**Figure 17: List of used libraries**

#### 4.5.2 Exploring data

As represented in figure 18 statistically the dataset used is composed of: 3 courses Math, JavaScript, VueJS, and the total score is over 100. The maximum value is 100 which shows that we do not have an outlier in the dataset. By counting we find that we have 1000 elements on each of the parameters, which assure that there are no missing values in the dataset to be used. The statistical mean in math is 66.09, in JavaScript is 69.17 in VueJs is 68.05, while the overall mean is 67.77.

The statistical standard deviation in math is 15.16, that of JavaScript is 14.60, in VueJS is 15.20 and that of overall standard deviation is 14.26. The Statistical minimum shows that in math is 0, in JavaScript is 17 marks, in VueJS is 15, while the overall minimum is 9.

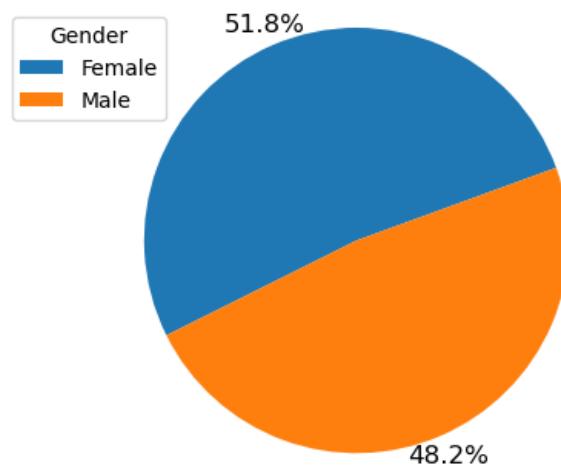
The statistical quarter Median and  $\frac{3}{4}$  in math are 57, 66, and 77 respectively, while in JavaScript is 59.70 and 79. In VueJS, there are 57.75, 69 and 79, and lastly the overall become 58.33, 69.33, and 77.67.

	count	mean	std	min	25%	50%	75%	max
math	1000.000000	66.089000	15.163080	0.000000	57.000000	66.000000	77.000000	100.000000
JavaScript	1000.000000	69.169000	14.600192	17.000000	59.000000	70.000000	79.000000	100.000000
VueJS	1000.000000	68.054000	15.195657	10.000000	57.750000	69.000000	79.000000	100.000000
score	1000.000000	67.770580	14.257311	9.000000	58.330000	68.330000	77.670000	100.000000

**Figure 18: Statistical exploration of the dataset**

#### 4.5.2.1 The proportion of students by gender

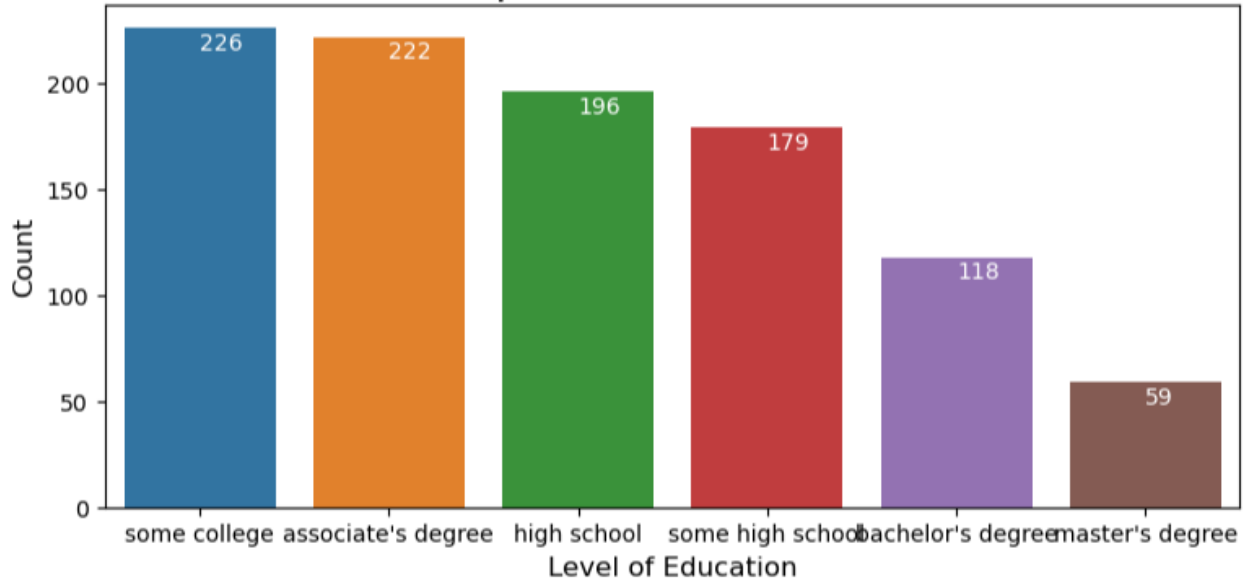
Figure 19 shows that 51.8% of the dataset records are female while 48.2% are male.



**Figure 19: The proportion of students by gender**

#### 4.5.2.2 The Parents' level of Education

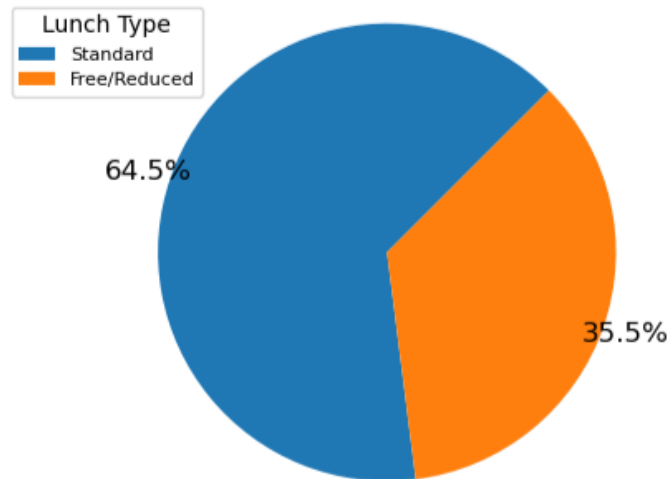
Figure 20 below depict that in the dataset, 226 of the parents attended some college, 222 have an associate degree, 196 with only high school, 179 have some high school, 118 with a bachelor's degree while 59 only have a master's degree.



**Figure 20: The parent’s level of education**

#### 4.5.2.3 Type of lunch often taken

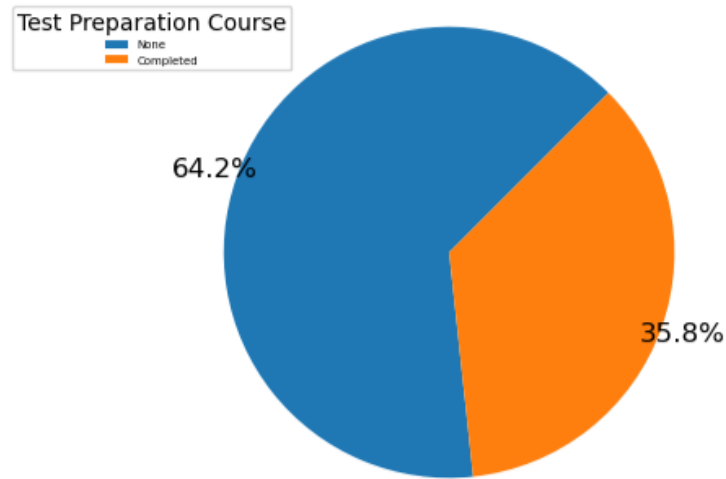
By observing the figure 21, we can see that the food given to the students often is of a balanced diet (standard) 64.5%. in some cases, 35.5% of the food given to the students is of a reduced quality based compared to the balanced diet food requirements.



**Figure 21: Type of lunch often taken**

#### 4.5.2.4 The number of students who took the test preparation course

As shown in figure 22, 64.2% of the students took the test preparation course while a portion 35.8% didn't. The results of this action are shown in the prediction of the performance see page 47



**Figure 22: The number of students who took the test preparation course**

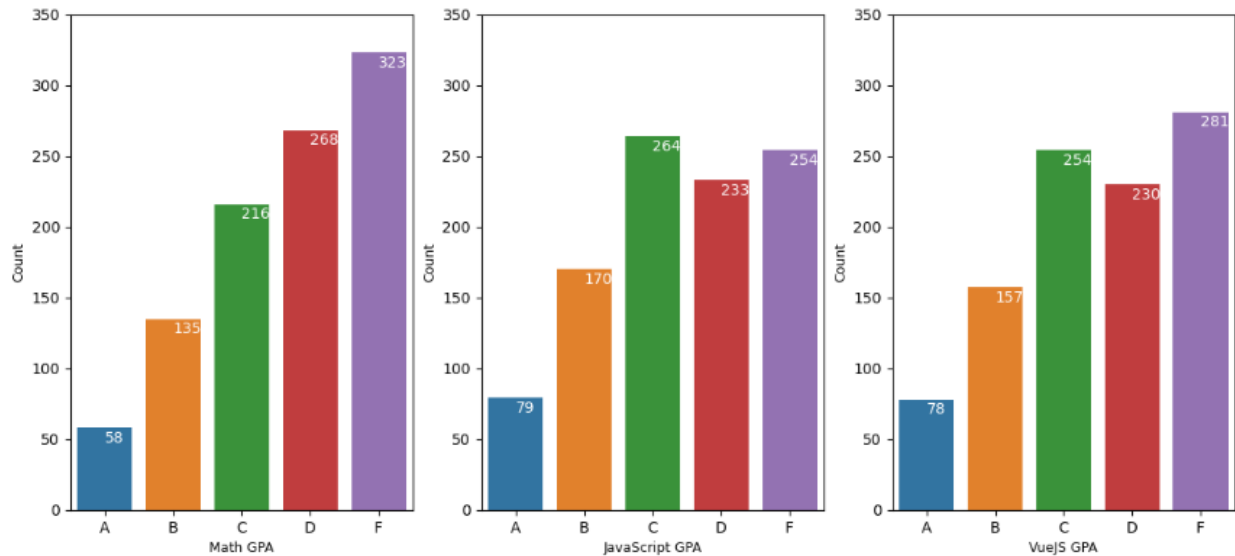
#### 4.5.2.5 Students' grades distribution

The records in the dataset include marks of students. By representing them in GPA (Grade Point Average). This is a number of the level of score a student got which corresponds to an education grading system to determine the academic student performance. To map grades, here marks from 90 and above are marked in grade A, 80 to 89 for grade B, grade C is from 70 to 79, and grade D is between 60 and 69. Marks less than 60 are taken in the grade of F which we can say that the student failed the course.

Figure 23 shows the Math GPA distribution 58 students got a grade A, 135 with a grade of B, 216 have grade of C, 268 have grade of D while 323 students have less than 60 and are classified in grade F.

In JavaScript the grade A is for 79 students, 170 got the grade B, Grade C is for 264 students, and grade D with 233 students, while 254 remaining students have grade F.

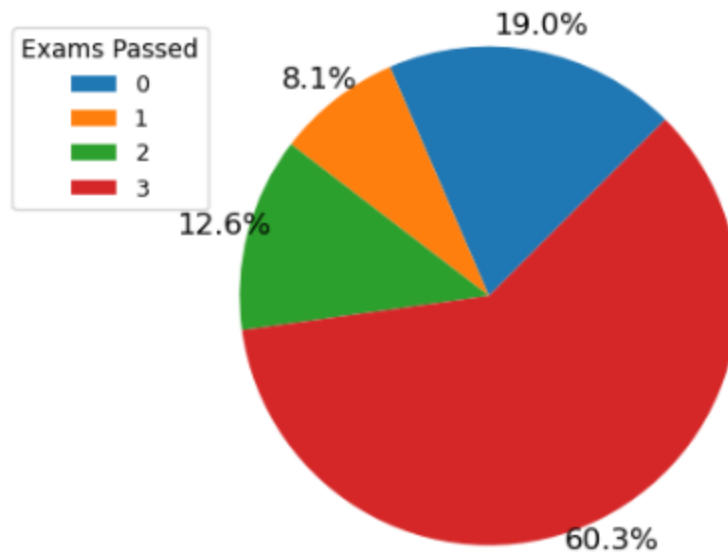
In VueJS, 78 have grade A, 157 with grade B, 254 represent the grade C, 230 got the grade D. while in grade F there is 281. Meaning that the statistical calculations show that there is a need to pay special attention on 323 students in mathematics, 254 in JavaScript, and 281 in VueJS.



**Figure 23: Students grades distribution**

#### 4.5.2.6 The number of exams a student succeeded

The legend of figure 24: 0 represent the % number of students who failed in all exams ( meaning they have less than 60%, 1 represent the % of students who failed in 1 exam, 2 represent the % of student failed in two exams, while 3 represent the % of student who succeeded in all exams. From the same figure, 19.0% Failed in all exams, 8.1% passed in one exam only, 12.6% passed in two exams, while 60.3% passed in all exams.



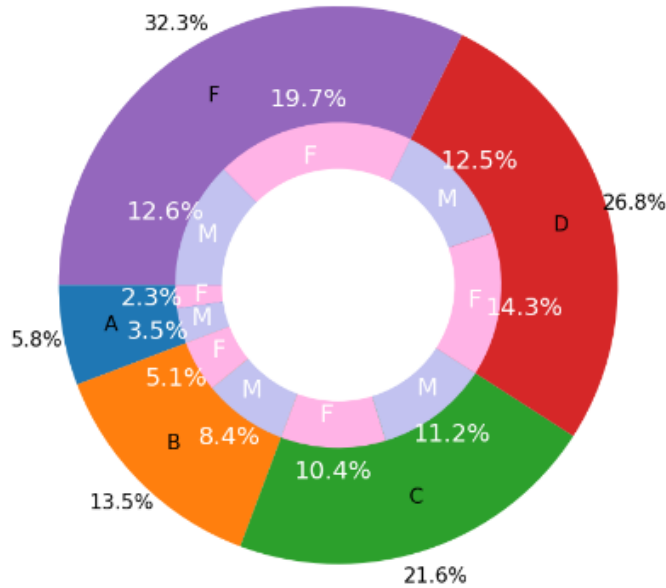
**Figure 24: The number of exams a student succeeded**

#### 4.5.2.7 Performance by gender in math

Males and females are represented by the letter M and F in the following coming figures of maths JavaScript and VueJS respectively.

The percentage depicted in figure 25 below shows that the grade A is 5.8 % where 3.5% are males and 2.3% are females. In grade B, the portion is 13.5% including 8.4% males and 5.1 females. The grade C represent 21.6%, within 11.2% are males and 10.4% being females. For the grade D, its portion is of 26.8% having 12.5% males and 14.3% females. Finally, the grade F its portion is 32.3% counting 12.6% males and 19.7% females.

Based on those statistical numbers, males did better in mathematics as they have a great portion in A, B, and C grades and they have a smaller number of failed students in grade F and D.

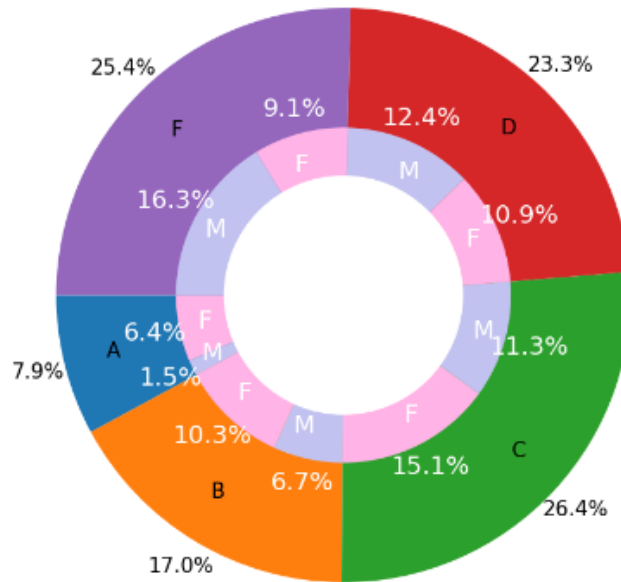


**Figure 25: Performance by gender in math**

#### 4.5.2.8 Performance by gender in JavaScript

In JavaScript by referring to the figure 26 on next page, the students who got the grade A have the portion of 7.9% which include 1.5% males, and 6.4% females. The grade B represent a portion of 17% where 6.7% are males and 10.3% are females. The students who are categorized in the grade C represents the portion of 26.4% having 11.3% males and 15.1% females. The grade category of D has a portion of 23.3% where 12.4% are males and 10.9% are females. The grade category of F counts 25.4% within 16.3% are males, and 9.1% are males.

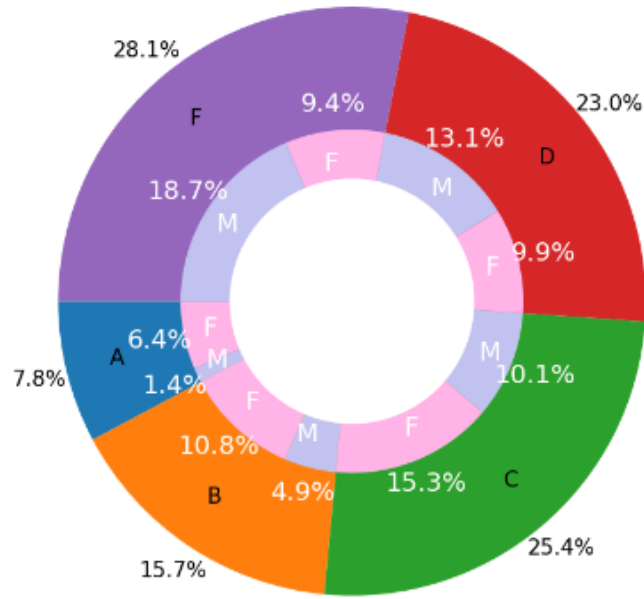
Based on the above numbers, females got better results in the grade A B and C that males who have a great portion in the category of grade F in JavaScript.



**Figure 26: Performance by gender in JavaScript**

#### 4.5.2.9 Performance by gender in VueJS

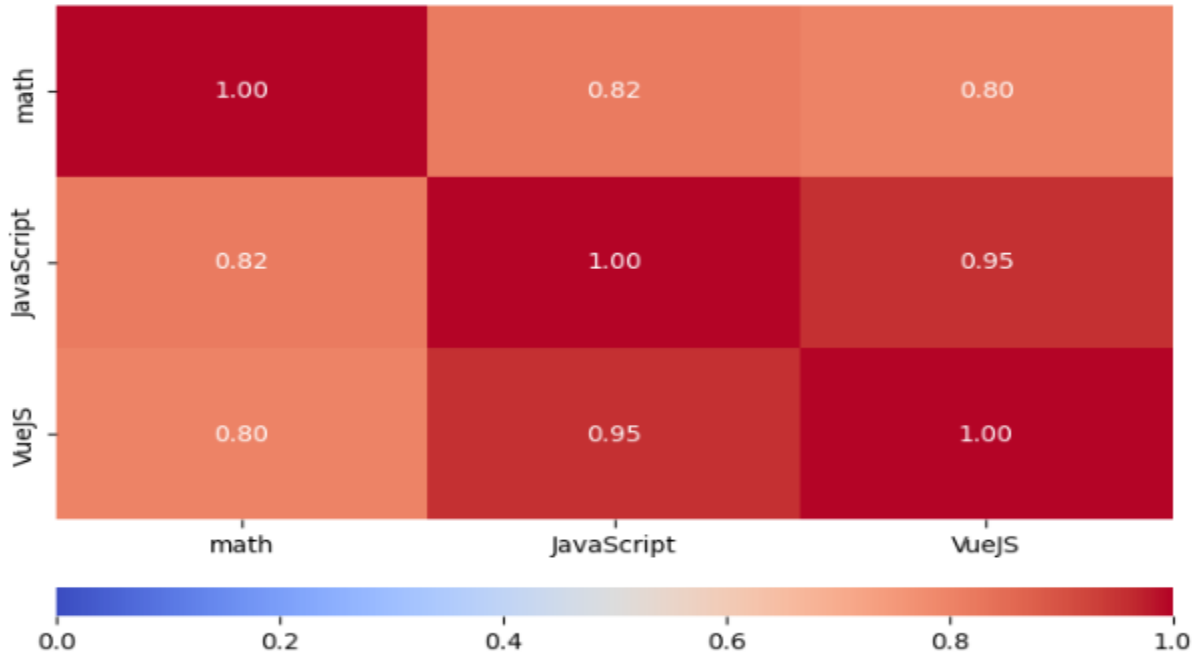
In VueJS, represented by the figure 27 below only 1.4% males again 6.4% females in a total of 7.8% they got the grade A. grade B has a portion of 15.7% where 4.9% are males and 10.8% females. The grade C which is the portion of 25.4%, 10.1% are males and 15.3% females. In category of grade D 23% in total where 13.1% are males and 9.9% females. Lastly the grade category of F covers a portion of 28.1% including 18.7% of males and only 9.4 females.



**Figure 27: Performance by gender in VueJS**

The best performance in VueJS is on the side of females. They got better marks in grade A,B, and C. Even they have small portion of female's students in the category of F and D, which indicate their best performance in VueJS.

#### 4.5.2.10 The relation between different subjects



**Figure 28: The relation between different subjects**

The above figure 28 shows that JavaScript scores and VueJS scores are highly correlated so a student who knows JavaScript finds himself in knowing VueJS easily.

As stated on the website<sup>29</sup>, the formula for the correlation coefficient is expressed as follows:  $r = \frac{(n * \Sigma XY - \Sigma X * \Sigma Y)}{\sqrt{[(n * \Sigma X^2 - (\Sigma X)^2) * (n * \Sigma Y^2 - (\Sigma Y)^2)]}}$ . The elements within this formula are defined as follows: "n" represents the total number of data points, " $\Sigma XY$ " denotes the summation of the product of the x-values and y-values for each data point, " $\Sigma X$ " represents the summation of the x-values in the dataset, " $\Sigma Y$ " represents the summation of the y-values in the dataset, " $\Sigma X^2$ " signifies the summation of the squares of the x-values in the dataset, and " $\Sigma Y^2$ " signifies the summation of the squares of the y-values in the dataset. We calculated correlations between JavaScript, Math, and VueJS as both X and Y values.

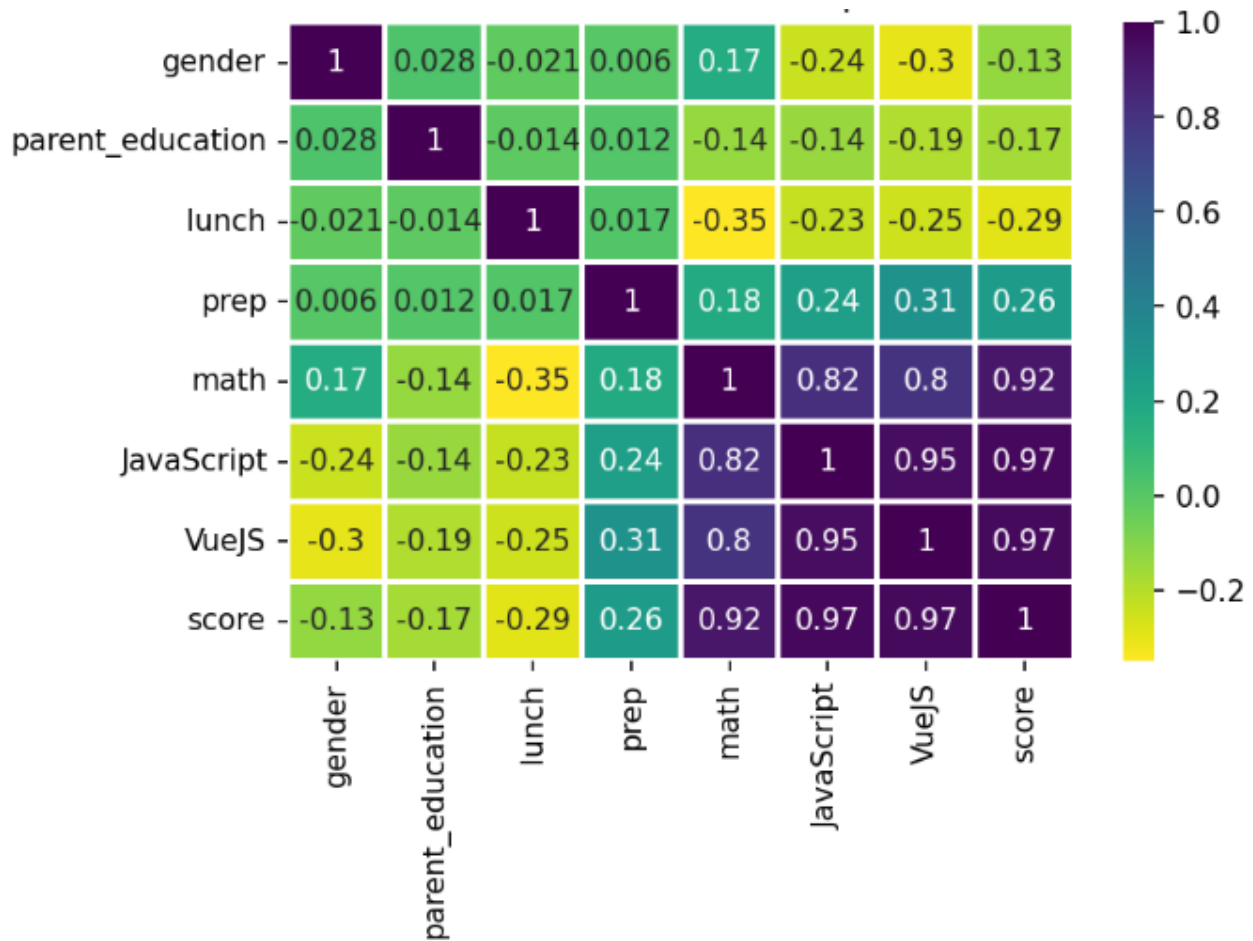
There is a high probability that a student who knows the ways around JavaScript will perform well in VueJS. For example, there is a probability of 95% that a student who succeeded in JavaScript will do also in VueJS. Math scores have lower correlations with other subjects, it is uncertain that if a student is good at math, it implies that he/will be good at other subjects or vice

<sup>29</sup> <https://tinyurl.com/53v557ha> accessed 10/10/2023 10:59:04

versa. Though its probability can't be ignored 80% and 82% are not a small number in probability.

#### 4.5.2.11 Correlation between different parameters

Statistical correlation is the relationship between two variables. It is positive if the value of one variable implies the increase of the dependent variable or one variable decreases the same time with the other, otherwise it is negative. There exists also the linear and non-linear correlation.



**Figure 29: Corelation between different parameters**

The above Figure 29 shows the positive and negative correlation between variables. For example, the relationship between VueJS and course preparation is 0.31 and the math results are negatively correlated to the parent education with a value of -0.14.

#### 4.5.2.12 Summary of probability

A cross-processing between different parameters from the dataset used gives the following probabilities:

1. The Probability of Students Passing in all the Subjects is 79.60 %
2. The Probability of Students Passing in all the Subjects is 2.30 %  
Sample mean for Math Scores: 63.12  
Population mean for Math Scores: 66.089.
3. Sample mean for Reading Scores: 68.5  
Population mean for Reading Scores: 69.169.
4. Sample mean for Writing Scores: 63.12  
Population mean for Writing Scores: 68.054
5. z-critical<sup>30</sup> value: 1.6448536269514722  
Confidence interval<sup>31</sup>: (64.82729483328328, 66.40470516671672)  
True mean: 66.089

The z critical value is a statistical probability value which has a point that cut off under the standard normal distribution. It is a reference of a probability that specific value will have.

Confidence interval is a range of probabilistic values that assure that value of parameter will lie within them.

## 4.6 Building Artificial intelligence Model

In Machine Learning a model is a program that is trained to find some patterns and make decision based on previously unseen dataset<sup>32</sup>. Models are built using different types of algorithms.

---

<sup>30</sup><https://tinyurl.com/p8b6m2e5> accessed 05/08/2023 07:48:06

<sup>31</sup>[https://en.wikipedia.org/wiki/Confidence\\_interval](https://en.wikipedia.org/wiki/Confidence_interval) accessed 14/09/2023 12:08:16

<sup>32</sup><https://tinyurl.com/bddbfbh89> accessed 23/09/2023 14:17:52

Machine learning algorithm is a set of step-by-step instructions given to a computer to execute planned tasks mostly for predicting results from input dataset values. To predict the results of the student, we will make a comparison between different models including: Linear Regression, K Nearest Neighbor Regression, SVR, and Neural Networks ( see pros and cons for them in table below):

**Table 1: Advantages and disadvantages of used models**

Algorithm name	function	advantages	disadvantages
Logistic Regression	$Y = \text{sigmoid}(b_0 + b_1X_1 + b_2X_2 + \dots + b_nX_n)$	<ul style="list-style-type: none"> <li>-Easy to implement and interpret yet efficient in training</li> <li>-The predicted parameters give inferences about the importance of each feature</li> <li>-Performs well on low-dimensional data.</li> <li>Performs well on low-dimensional data.</li> <li>Very efficient when the dataset has features that are linearly separable.</li> <li>Outputs well-calibrated probabilities along with classification results.</li> </ul>	<ul style="list-style-type: none"> <li>Overfits on high dimensional data</li> <li>Nonlinear problems can't be solved with logistic regression since it has a linear decision surface</li> <li>Assumes linearity between dependent and independent variables.</li> <li>Fails to capture complex relationships.</li> <li>Only important and relevant features should be used otherwise model's predictive value will degrade.</li> </ul>
Support Vector Regression	$y = wx + b$	<ul style="list-style-type: none"> <li>It is robust to outliers.</li> <li>Decision model can be easily updated.</li> <li>It has excellent generalization capability, with high prediction accuracy.</li> <li>Its implementation is easy.</li> </ul>	<ul style="list-style-type: none"> <li>Unsuitable to Large Datasets.</li> <li>Large training time.</li> <li>More features, more complexities.</li> <li>Bad performance on high noise.</li> <li>Does not determine Local optima.</li> </ul>
K-NN	Euclidean Distance( $x^*, x$ ) = $\sqrt{\sum_i^m (x_i^* - x_i)^2}$	<ul style="list-style-type: none"> <li>It's easy to understand and simple to implement</li> <li>It can be used for both classification and regression problems</li> <li>It's ideal for non-linear data since there's no assumption about underlying data</li> <li>It can naturally handle multi-class cases</li> <li>It can perform well with enough representative data</li> </ul>	<ul style="list-style-type: none"> <li>Associated computation cost is high as it stores all the training data</li> <li>Requires high memory storage</li> <li>Need to determine the value of K</li> <li>Prediction is slow if the value of N is high</li> <li>Sensitive to irrelevant features</li> </ul>

ANN	$Y=W1X1+W2X2+b$	Uses Large Datasets Non-Linear and Flexible Handles Missing Data Automation and Multitasking	Requires Huge Amounts of Data Issues with Interpretation Adversarial Vulnerabilities High Computational Requirements
-----	-----------------	---	---

**4.6.1 Linear Regression Model**

To validate the model the validation set approach is to split the test portion of the dataset. The training data takes 70% and the test will be 30%. The parameters that are used the average from math JavaScript and VueJS , gender, parent education, type of lunch taken, and the test preparation. The Figure 30 shows that the R-square<sup>33</sup> value is too low 0.223. This can be the results of small number of predictors which are used in the dataset.

---

<sup>33</sup> a statistical measure in a regression model that determines the proportion of variance in the dependent variable that can be explained by the independent variable

Dep. Variable:	average_score	R-squared:	0.223
Model:	OLS	Adj. R-squared:	0.214
Method:	Least Squares	F-statistic:	24.81
Date:	Mon, 12 Jun 2023	Prob (F-statistic):	1.03e-33
Time:	06:35:41	Log-Likelihood:	-2763.9
No. Observations:	700	AIC:	5546.
Df Residuals:	691	BIC:	5587.
Df Model:	8		
Covariance Type:	nonrobust		

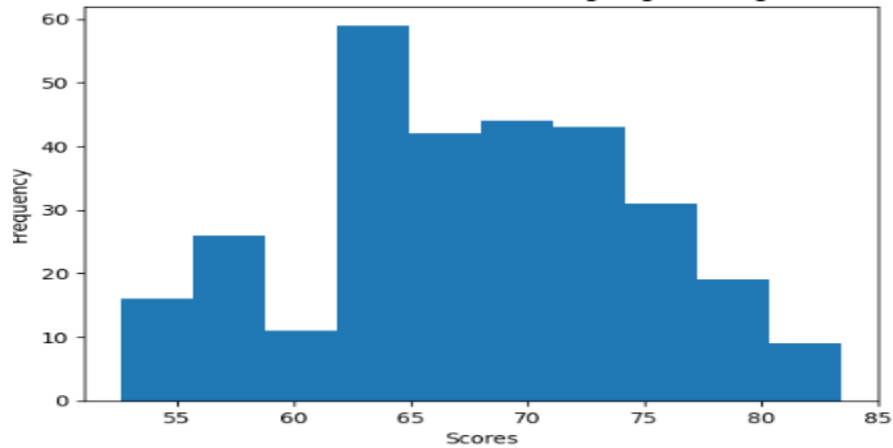
	coef	std err	t	P> t	[0.025	0.975]
Intercept	70.1611	1.478	47.465	0.000	67.259	73.063
gender_male	-4.5280	0.958	-4.729	0.000	-6.408	-2.648
parental_bachelor	1.5988	1.713	0.933	0.351	-1.764	4.962
parental_hs	-5.5560	1.497	-3.711	0.000	-8.495	-2.616
parental_masters	3.5700	2.186	1.633	0.103	-0.721	7.861
parental_somecol	-0.8084	1.473	-0.549	0.583	-3.700	2.084
parental_somehs	-6.0144	1.550	-3.880	0.000	-9.058	-2.971
lunch_standard	9.7456	1.011	9.642	0.000	7.761	11.730
test_prep_none	-7.0187	1.002	-7.003	0.000	-8.987	-5.051

Omnibus:	8.156	Durbin-Watson:	1.951
Prob(Omnibus):	0.017	Jarque-Bera (JB):	8.333
Skew:	-0.265	Prob(JB):	0.0155
Kurtosis:	2.925	Cond. No.	9.47

**Figure 30: Summary of the prediction using logistic regression**

Observing the histogram, it shows that the model is predicting that the results values is in between 40 and around 85 only. Also, we can see the similarity of the data distribution to the average results meaning that the scores are not blindly predicted all the time.



**Figure 31: Distribution of predicted scores using logistic regression**

To facilitate the comparison of models at the end of each type of model used, we calculate the Root Mean Square Error (RMSE<sup>34</sup>)

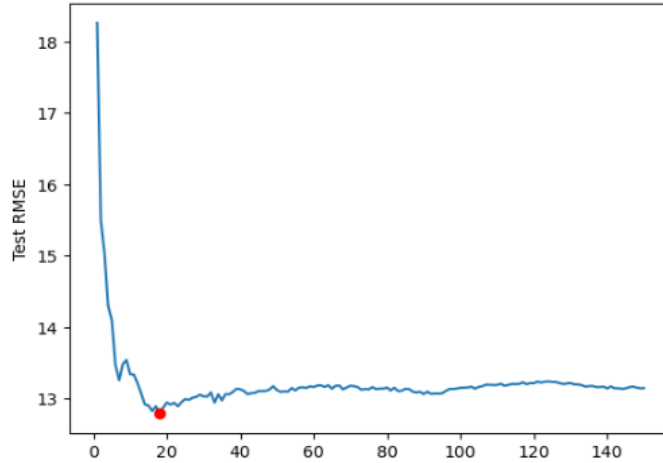
#### **4.6.2 K Nearest Neighbours Regression**

Fitting values of this model is the same as previous (Logistic Regression). But, in this we fit the value of K from 1 to 150 then after we calculate the RMSE as we did early. The K<sup>35</sup> value should not be as large as this to be processed sequentially. Here it is because we are using a small dataset. Otherwise, it should a large amount of computing capacity. The Figure 32 shows the smallest value of K which is 35 where the largest reduction of RMSE(Root Mean Square Error)

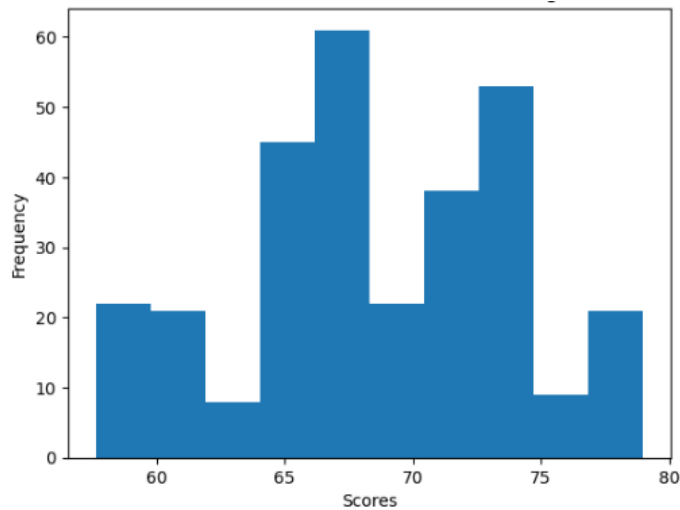
---

<sup>34</sup> It measures the average difference between values predicted by a model and the actual values. It provides an estimation of how well the model is able to predict the target value (accuracy).

<sup>35</sup> Defines how many neighbours will be checked to determine the classification of a specific query point.



**Figure 32: Plot of test RMSE vs K number of neighbours**



**Figure 33: Distribution of predicted scores for KNN regression**

The histogram on figure 33 shows that the predicted results fall in a short range from 55 to around 80. The algorithm of KNN gives the scores averages of the neighbours, then the prediction tends to be closer to the value of the average data. Meaning that the KNN is not giving the very accurate prediction especially when results fall in a lower range.

### 4.6.3 Support Vector Regression (SVR)<sup>36</sup>

The process of providing data is the same as in previous algorithms. In particular, for this model we will use different Kernel and select the best one. Namely: Linear, Polynomial, Sigmoid and Gaussian

	kernel	Test RMSE	Pred
0	linear	12.922286	[60.26050331762064, 64.0011761416976, 56.66086...
1	polynomial	13.092485	[58.89332864251796, 63.89965648539629, 56.0999...
2	sigmoid	12.983928	[64.56950584391251, 64.96699903935362, 61.5995...
3	rbf	13.031163	[60.755662578068645, 65.37235004097501, 58.600...

**Figure 34: Support vector regression RMSE and predictions**

Once the support vector is given enough parameters it gives the best RMSE among other model used here. Next, we are trying to improve the results by using the grid search to get the best parameters. Then the prediction is done for the test data using the best parameters we got from the grid search, and append this test RMSE to the table of the results of the model built for SVR, then it produces the output on the figure 34.

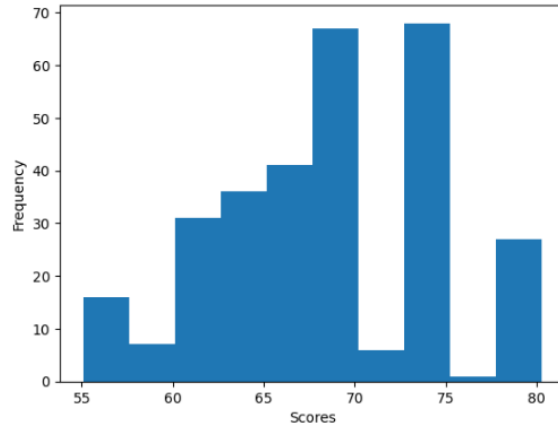
	kernel	Test RMSE	Pred
0	linear	12.922286	[60.26050331762064, 64.0011761416976, 56.66086...
1	polynomial	13.092485	[58.89332864251796, 63.89965648539629, 56.0999...
2	sigmoid	12.983928	[64.56950584391251, 64.96699903935362, 61.5995...
3	rbf	13.031163	[60.755662578068645, 65.37235004097501, 58.600...
0	GS Output	12.910383	[59.80043261256979, 63.7997861472088, 55.80033...
0	GS Output	12.910383	[59.80043261256979, 63.7997861472088, 55.80033...

**Figure 35: Prediction using the lowest RMSE in combination with grid search**

To visualize the distribution of predicted values, we plot the histogram on Figure 36. The figure shows that the predicted scores are the interval from 55 to 80. This result is similar to previous models, meaning that it is not predicting scores on the upper (greater than 85) and lower values (less than 55)

---

<sup>36</sup> a type of machine learning algorithm used for regression analysis. The goal of SVR is to find a function that approximates the relationship between the input variables and a continuous target variable, while minimizing the prediction error.



**Figure 36: Distribution of predicted scores for SVR**

## CHAPTER 5 ANALYSIS OF THE RESULTS

By combining all tests made on the RMSE as depicted in figure 37 in previous page, the final results show that the Linear Regression and Artificial Neural Network provide the best models for predicting students' results.

Anyway, all 4 models used here, predicted the results in a small range between 40 and 85 and all of them are not able to predict the scores which are higher than 85 the same as those lower than 40. This is because the dataset is small and the data picked in those ranges are relatively small such that they may not be picked from while sampling the testing data or training. This can be improved by fitting into the model a larger dataset which can predict and give more precise outputs.

Another influence here is that we have many qualitative parameters like parent level of education, feeding status, test preparation and so on. The models are predicting the results which are numerical. Which means that if the numerical or quantitative data is increased the more the efficiency of the results.

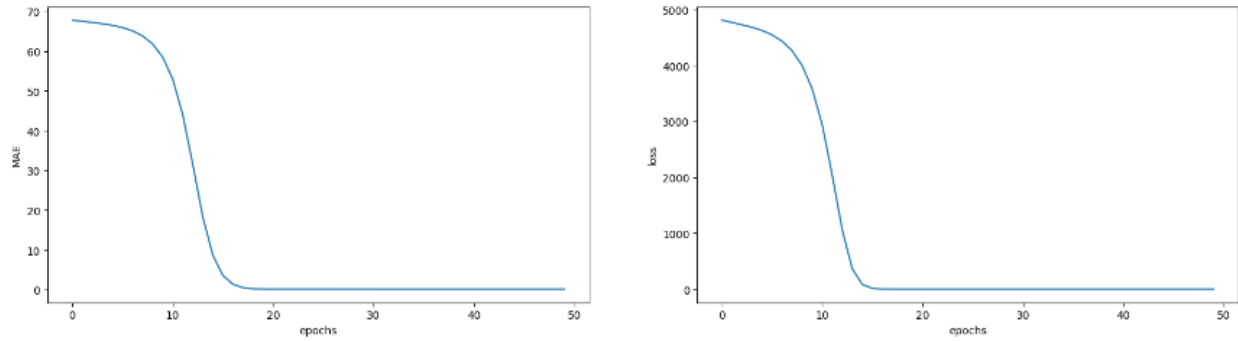
	<b>Model</b>	<b>Test RMSE</b>
<b>0</b>	Linear Regression	11.901006
<b>1</b>	KNN Regression	12.228900
<b>2</b>	SVR	11.997845
<b>3</b>	Neural Network	11.940058

**Figure 37: Comparison between models**

To track the accuracy of the model, we plot the figure 39 produced as a result of MAE<sup>37</sup> against loss

---

<sup>37</sup> measures the average magnitude of the errors in a set of forecasts, without considering their direction. It measures accuracy for continuous variables



**Figure 38: Mean Absolute Error (MAE) and loss<sup>38</sup> on each epoch**

Based on definition from ChatGPT<sup>39</sup>; in machine learning, an epoch refers to one complete iteration through the entire training dataset during the training process of a neural network or other machine learning models. During each epoch, the model processes and learns from all the training examples in the dataset, updating its internal parameters (such as weights and biases) based on the error or loss it computes for the predictions made during that iteration.

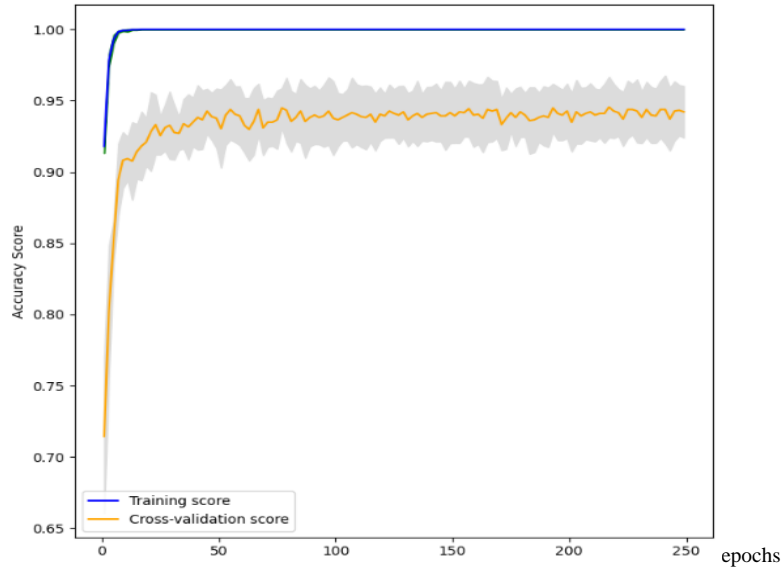
Multiple epochs are typically required in the training process to improve the model's performance gradually. The number of epochs is a hyperparameter that can be adjusted by the practitioner based on factors such as the complexity of the problem, the size of the dataset, and the model's convergence behavior. Increasing the number of epochs can help the model learn more, but it should be balanced to avoid overfitting (learning the training data too well and performing poorly on unseen data) or excessive training time.

Like presented above, in Figure 38 the Mean absolute error is higher when we have few epochs and reduces tending to be 0 as the number of epochs increase, which means that the error in prediction is reducing as the model runs more chunks of epochs.

The same status is on loss values, they reduce too as the epochs number increases.

<sup>38</sup> is a row-level error calculation where the non-negative difference between the prediction and the actual is calculated

<sup>39</sup> <https://chat.openai.com/c/b42e403e-54a1-4026-b338-1b0d1d753459> accessed 10/10/2023 11:16:25



**Figure 39: Model cross-validation score**

The model accuracy in above figure shows that has validated the dataset at a rate of 100% then the results of cross validation score is between 90 and 95%. Meaning that the model can accurately predict the exact value at a rate up to 95%.

The process outlined in Figure 16 involves preparing data from Storyline, configuring xAPI to align with Moodle content, tracking student activities on Moodle with xAPI, storing data in an LRS, and using this dataset for predictive modelling of student outcomes. The results are used as students' predicted performance, then used in decision making.

As suggested in Chapter 6, the results can be improved by increasing prediction defence variables.

## **CHAPTER 6 CONCLUSION AND RECOMMENDATION**

### **6.1 Conclusions**

In the process of prediction, the testing and training datasets are randomly selected from a dataset. In this specific case, we employed a dataset comprising 1000 records, allocating 30% for testing and 70% for training. This means that 300 records were reserved for testing, while 700 were utilized for training. The dataset's size holds significance as it directly influences the model's predictive accuracy, with larger datasets generally leading to improved accuracy.

We have also observed that the models exhibited limitations in predicting scores falling below 40% and exceeding 90%. This limitation arises from the limited representation of data points within these score ranges in the dataset. Given their relatively low frequency, these outliers have a minimal impact on the predictive results. One potential solution to address this issue is to augment the dataset size, particularly focusing on including more data points within these problematic score ranges.

Additionally, it's worth noting that qualitative data had an adverse effect on the model's outputs, even though it holds significant importance in the statistical analysis of result sources. It is essential to recognize that in the context of predicting student performance at the end of a term, quantitative data tends to carry more weight and relevance compared to qualitative data.

Utilizing an online platform for educational purposes offers several advantages when it comes to data recording, streamlining daily teaching-related tasks, and enhancing overall efficiency. However, it is crucial to acknowledge that the educational system cannot exclusively depend on digital tools. There remains a fundamental need for in-person instruction, particularly in areas involving management and tasks that are inherently non-digital.

For instance, within our dataset, the manual recording of a student's daily dietary habits is essential. This manual process is necessitated by the dynamic nature of a student's feeding status, which can change over time.

As demonstrated in the statistical analysis (refer to Figure 21), it becomes evident that dietary habits have a significant influence on students' academic performance, a priority encouraged by the government of Rwanda. Therefore, it is imperative to emphasize the importance of ensuring that students have access to a balanced diet as a key educational strategy.

In summation, this thesis has delved into the critical domain of predicting students' exam performance, with a particular focus on harnessing the capabilities of xAPI (Experience API), Storyline, Moodle, as well as the pivotal roles played by Adaptive AI and Generative AI. This research underscores the paramount importance of this area within the broader educational context.

Through an exhaustive review of existing literature and rigorous data analysis facilitated by xAPI, this study has yielded key insights. These insights highlight the multifaceted nature of factors influencing student outcomes, encompassing prior academic records, study habits, socioeconomic backgrounds, and interactions within the Moodle learning environment.

Moreover, this research has showcased the efficacy of various machine learning models, including Linear Regression, Support Vector Regression, K-Nearest Neighbors (K-NN), and Artificial Neural Networks (ANN), when applied to the task of predicting student performance in exams conducted within the Moodle environment. Furthermore, the integration of xAPI data from Storyline has proven instrumental in enhancing predictive capabilities.

However, it's important to note that the role of Adaptive AI and Generative AI should not be underestimated. These cutting-edge technologies enable personalized learning experiences, adapting content and resources to individual student needs. They also have the potential to generate innovative educational materials and assessments, thereby enhancing the quality of education and facilitating more precise predictive models.

The results underscore the potential of early identification of students at risk of underperformance, facilitated by the utilization of xAPI, Storyline, Moodle, Adaptive AI, and Generative AI. This early recognition empowers educational institutions to proactively implement tailored interventions and support mechanisms, ultimately fostering student success, engagement, and retention in an ever-evolving educational landscape.

## **6.2 Recommendations:**

Building upon the findings of this study and considering the roles of xAPI, Storyline, Moodle, Adaptive AI, and Generative AI, the following recommendations are proposed:

- 1. Enhanced Data Integration:** Educational institutions should prioritize the seamless integration of xAPI, Storyline, and Moodle data to create a comprehensive view of student engagement and performance. This integrated data, coupled with Adaptive AI and Generative AI capabilities, can serve as a foundation for more accurate predictive models.
- 2. Continuous Model Refinement:** Institutions should establish mechanisms for the continuous improvement and refinement of predictive models. These models should

adapt to changes in the learning environment and student demographics while effectively utilizing data from xAPI, Storyline, Moodle, Adaptive AI, and Generative AI.

- 3. Personalized Learning Paths:** Leveraging insights from predictive models and the capabilities of Adaptive AI, institutions can develop highly personalized learning paths for students. These paths can cater to individual needs and promote better outcomes, with content generated by Generative AI.
- 4. Ethical Data Usage:** While utilizing these technologies, institutions should adhere to ethical standards, particularly concerning data privacy and fairness in algorithmic decision-making. Transparent and accountable data practices remain essential.
- 5. Inter-Institutional Collaboration:** Collaborative initiatives among educational institutions, technology providers, data scientists, and AI specialists can facilitate the sharing of best practices and the development of standardized data interoperability frameworks, further enhancing the utility of these technologies.
- 6. Future Research Directions:** Given the dynamic nature of educational technologies, further research should continue exploring advanced analytics techniques, the impact of mental health data, and the continuous evolution of Adaptive AI and Generative AI in education.

By implementing these recommendations and harnessing the capabilities of xAPI, Storyline, Moodle, Adaptive AI, and Generative AI, educational institutions can fully realize the potential of predictive analytics to enhance student outcomes, engage learners more effectively, and elevate the overall quality of education.

## References

- Abdullah., M. A., 2015. Learning style classification based on student's behavior in moodle learning management system. *Transactions on Machine Learning and Artificial Intelligence*, Volume 3(1), p. 28.
- Ahmed, A. B. E. D. a. E. I. S., 2014. Data mining: a prediction for student's performance using classification method.. *World J. Comput. Appl. Technol.* 2,, p. 43–47.
- Alan Berg, M. S. H. D. S. T. M. S., 2016. Dutch cooking with xapi recipes: The good, the bad, and the consistent.. Austin, TX, USA, IEEE, p. 234–236.
- Alario-Hoyos, C. M.-M. P. J. P.-S. M. D. K. C. a. P. G. H. A., 2016. Who are the top contributors in a mooc? relating participants' performance and contributions.. *J. Comput. Assist. Learn*, p. 232–243.
- Al-Barrak, M. A. a. A.-R. M., 2016. Predicting students final gpa using decision trees: a case study.. *Int. J. Inform. Educ. Technol.* 6, p. 528.
- Al-Radaideh, Q. A. A.-S. E. M. a. A.-N. M. I., 2006. International Arab Conference on Information Technology (ACIT'2006), Yarmouk University, Jordan.
- Al-Radaideh, Q. A. A.-S. E. M. a. A.-N. M. I., 2006. s.l., s.n.
- Al-Radaideh, Q. A. A.-S. E. M. a. A.-N. M. I., 2006. Mining student data using decision trees., s.l., s.n.
- Al-Radaideh, Q. A. A.-S. E. M. a. A.-N. M. I., 2006. Mining student data using decision trees., s.l., s.n.
- Andreas Konstantinidis, C. G., 2013. Using excel macros to analyse moodle logs..
- Andrews, R. D. J. a. T. A. B., 1995. Survey and critique of techniques for extracting rules from trained artificial neural networks.. *Knowl. Based Syst.* , p. 373–389.
- Aneesha Bakharia, K. K. A. P. D. G. S. D., 2016. Recipe for success: lessons learnt from using xapi within the connected learning analytics toolkit. *Proceedings of the sixth international conference on learning analytics & knowledge*, p. 378–382.

Anozie, N. a. J. B. W., 2006. Predicting end-of-year accountability assessment scores from monthly student records in an online tutoring system. s.l., Papers from the AAAI Workshop Menlo Park, CA: AAAI Press.

Asmaa Elbadrawy, A. P. Z. R. M. S. G. K. a. H. R., 2016. Predicting student performance using personalized analytics.. *Computer*, Volume 49(4), p. 61–69.

Bayer, J. B. H. G. J. O. T. a. P. L., 2012. Predicting drop-out from social behaviour of students. *International Conference on Educational Data Mining*, p. 103–109.

Behrouz Minaei-Bidgoli, D. A. K. G. K. W. F. P., 2003. Predicting student performance: an application of data mining methods with an educational web-based system.. *33rd Annual Frontiers in Education*, Volume 1, p. T2A–13.

Boser, B. E. G. I. M. a. V. V. N., 1992. A training algorithm for optimal margin classifiers. *Proceedings of the fifth Annual Workshop on Computational Learning Theory (Pittsburgh, PA: ACM)*, p. 144–152.

Cakmak, A., 2017. Predicting student success in courses via collaborative filtering. *Int. J. Intell. Syst. Appl. Eng.* , pp. 10-17.

Can, A. T., 2015. What is the tin can api.. *What is the tin can API*, Volume 31, p. 2015.

Chiara Masci, G. J. a. T. A., 2018. Student and school performance across countries: A machine learning approach.. *European Journal of Operational Research*, Volume 269(3), p. 1072–1085.

Chih-Ming Chen, M.-C. C., 2019. Mobile formative assessment tool based on data mining. *Computers & Education*, Volume 57(2), p. 256–273.

Cristóbal Romero, S. V. a. E. G., 2008. Data mining in course management systems: Moodle case study and tutorial.. *Computers & Education*, Volume 51(1), p. 368–384.

De Barba, P. K. G. E. a. A. M., 2016. The role of students' motivation and participation in predicting performance in a mooc.. *J. Comput. Assist. Learn.*, pp. 218-231.

Elbadrawy, A. S. S. a. K. G., 2014. Personalized multi-regression models for predicting students performance in course activities. *UMN CS*.

- e-learning, C. s. o. x. a. t., 2015. Case studies of xapi applications to e-learning.. s.l., s.n., pp. 1-3.
- Félix Pascual-Miguel, J. C.-P. Á. H.-G. a. S. I.-P., 2011. A characterisation of passive and active interactions and their influence on students' achievement using moodle lms logs. *International Journal of Technology Enhanced Learning*, Volume 3, p. 403–414.
- Gary N. Marks, J. C. & J. A., 2006. Explaining socioeconomic inequalities in student achievement. The role of home and school factors, *Educational Research and Evaluation*, Issue 12:2, pp. 105-128.
- Ibrahim, Z. a. R., 2007. Predicting students academic performance: comparing artificial neural network, decision tree and linear regression. 21st Annual SAS Malaysia Forum, 5th September (Kuala Lumpur)..
- Ijaz Khan, A. A. S. A. R. A. N. J., 2019. Tracking student performance in introductory programming by means of machine learning.. 2019 4th mec international conference on big data and smart city (icbdsc), p. 1–6.
- Ioannis E Livieris, K. D. V. T. T. T. A. M. a., 2019. Predicting secondary school students' performance utilizing a semi-supervised learning approach. *Journal of educational computing research*, Volume 57(2), p. 448–470.
- Jiawei Han, M. K., 2016. *Data Mining: Concepts and Techniques*. [Online] Available at: [www.cs.uiuc.edu/~hanj](http://www.cs.uiuc.edu/~hanj)
- Jonathan M Kevan, P. R. R., 2016. Experience api: Flexible, decentralized, and activity-centric data collection.. *Technology, knowledge and learning*, Volume 21(1), pp. 143–149, .
- Kennedy, G. C. C. D. B. P. a. C. L., 2015. Predicting success: how learners' prior knowledge, skills and activities predict mooc performance. *Proceedings of the Fifth International Conference on Learning Analytics and Knowledge* (New York, NY: ACM), pp. 136-140.
- Kim, B.-H. V. E. a. G. V., 2018. Gritnet: Student performance prediction with deep learning. arXiv preprint arXiv.

Kloft, M. S. F. Z. Z. a. P. N., 2014. Predicting mooc dropout over weeks using machine learning methods. in Proceedings of the EMNLP 2014 Workshop on Analysis of Large Scale Social Interaction in MOOCs (Doha), pp. 60-65.

Koprinska, I. S. J. a. Y. K., 2015. Predicting student performance from multiple data sources., Madrid, Springer, p. 678–681.

María Lucía Barrón-Estrada, R. Z.-C. R. O.-B. F. G.-H., 2017. Sentiment analysis in an affective intelligent tutoring system.. 2017 IEEE 17th international conference on advanced learning technologies (ICALT), p. 394–397.

Matjaž Debevc, P. K. A. H., 2011. Improving multimodal web accessibility for deaf people: sign language interpreter module.. Multimedia Tools and Applications, Volume 54(1), p. 181–199.

Mayilvaganan, M. a. K. D., 2014. Comparison of classification techniques for predicting the performance of students academic environment., 2014 International Conference on Communication and Network Technologies (Sivakasi: IEEE), pp. 113-118.

Meier, Y. X. J. A. O. a. V. d. S. M., 2015. Comparison of classification techniques for predicting the performance of students academic environment. s.l., s.n., p. 113–118.

Mostafa., L., 2019. Student sentiment analysis using gamification for education context.. International Conference on Advanced Intelligent Systems and Informatics, p. 329–339.

Mushtaq Hussain, W. Z. W. Z. S. M. R. A. a. S. A., 2019. Using machine learning to predict student difficulties from learning session data.. Artificial Intelligence Review, Volume 52(1), pp. 381-407.

Myriam Munezero, C. S. M. E. S. a. J. P., 2014. Are they different? affect, feeling, emotion, sentiment, and opinion detection in text.. IEEE transactions on affective computing, Volume 5(2), pp. 101-111.

Natriello, G., 2013. The impact of evaluation processes on students.. School and classroom organization, Issue Routledge, pp. 227-246.

- Nguyen Thai-Nghe, L. D. A. K.-G. a. L. S.-T., 2010. Recommender system for predicting student performance. *Procedia Computer Science*, Volume 1(2), p. 2811–2819.
- Oladokun, V. O. A. A. T. a. C.-O. O. E., 2008. Predicting Students Academic Performance Using Artificial Neural Network: A Case Study of an Engineering Course[J]. Hilo, HI: Akamai University,, p. 72–79.
- Olugbenga Wilson Adejo, T. C., 2018. Predicting student academic performance using multi-model heterogeneous ensemble approach.. *Journal of Applied Research in Higher Education*.,
- Oreški., N. K. a. D., 2018. Analysis of student behavior and success based on logs in moodle.. 2018 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO), p. 0654–0659.
- Pojon, M., 2017. Using machine learning to predict student performance (Master's thesis).
- Quadri, M. M. a. K. N., 2010. Drop out feature of student data for academic performance using decision tree techniques.. *Global J. Comput. Sci. Technol.*.
- Ravi., K. R. a. V., 2015. A survey on opinion mining and sentiment analysis: tasks, approaches and applications.. *Knowledge-based systems*, Volume 89, p. 14–46.
- Reeves, B., 2018. Development of rubrics to support teacher judgement of student proficiency in ethical Decision-Making.
- Ren, Z. R. H. a. J. A., 2016. Predicting performance on mooc assessments using multi-regression models. *arXiv preprint arXiv*.
- Ribeiro., F. B. a. R., 2013. *Procesamiento del lenguaje natural*, Volume 50, p. 77–84.
- Ribeiro., F. B. a. R., 2018. *arXiv preprint arXiv*, Issue 1804.07405.
- Robert A Sottolare, R. A. L. B. S. G., 2017. Enhancing the Experience Application Program Interface (xAPI) to Improve Domain Competency Modeling for Adaptive Instruction. Massachusetts, Cambridge, USA, s.n., p. 265–268.

Rovai, A. P. a. J. H., 2004. Blended learning and sense of community: a comparative analysis with traditional and fully online graduate courses.. nt. Rev. Res. Open Distribut. Learn., pp. 1-13.

Shahiri, A. M. a. H. W., 2015. A review on predicting student's performance using data mining techniques.. Proc. Comput. Sci. 72, p. 414–422.

Shovon, M. I. H. a. H. M., 2012. An approach of improving students academic performance by using k means clustering algorithm and decision tree.. arXiv preprint arXiv:.

Shubham Joshi, R. K. R. P. C., 2021. Evaluating artificial intelligence in education for next generation.. Journal of Physics: Conference Series, Volume 1714, p. 012039.

Slim, A. H. G. L. K. J. a. A. C. T., 2014. Employing markov networks on curriculum graphs to predict student performance,. 2014 13th International Conference on Machine Learning and Applications (Detroit, MI: IEEE), p. 415–418..

Sorour, S. E. M. T. G. K. a. H. S., 2014. Predicting students' grades based on free style comments data by artificial neural network. 2014 IEEE Frontiers in Education Conference (FIE) Proceedings, pp. 1-9.

Sotiris Kotsiantis, N. T. A. F. a. V. K., 2013. Using learning analytics to identify successful learners in a blended learning course.. International Journal of, Volume 5(2), p. 133–150.

Stéphane Brutus, M. B. D., 2010. Improving the effectiveness of students in groups with a centralized peer evaluation system.. Academy of Management Learning & Education, Volume 9(4), pp. 652-662.

Thai-Nghe, N. D. L. K.-G. A. a. S.-T. L., 2010a. Recommender system for predicting student performance.. Proc. Comput. Sci. 1, p. 2811–2819.

Thomas Rabelo, M. L. J. C. V. a. R. A., 2017. Comparative study of xAPI validation tools. Indianapolis, IN, USA, s.n.

Van Merriënboer, J. J. a. S. J., 2005. Cognitive load theory and complex learning: recent developments and future directions.. duc. Psychol. Rev. 17, p. 147–177.

- Xu, B. a. Y. D., 2016. Motivation classification and grade prediction for moocs learners.. Comput. Intell. Neurosci..
- Xu, J. M. K. H. a. V. D. S. M., 2017. A machine learning approach for tracking and predicting student performance in degree programs. IEEE J. Sel. Top. Signal Process. 11., pp. 742-753.
- Zafer Unal, A. U., 2017. Comparison of student performance, student perception, and teacher. International Journal of Instruction, p. 10(4):145.
- Zhang, Y. a. L. S., 2020. Integrated sparse coding with graph learning for robust data representation. IEEE Access 8, p. 161245–161260..
- Zhang, Y. A. R. C. J. a. S. X., 2021a. Undergraduate grade prediction in chinese higher education using convolutional neural networks. LAK21: 11th International Learning Analytics and Knowledge Conference., pp. 462-468.
- Zhang, Y. X. M. a. Y. B., 2018b. Hierarchical sparse coding from a bayesian perspective.. Neurocomputing 272, p. 279–293.
- Zhang, Y. Y. Y. A. R. C. J. D. H. & S. X., 2021. Educational data mining techniques for student performance prediction: method review and comparison analysis.. Frontiers in psychology,.
- Zhi Liu, C. Y. S. R. S. L. L. Z. a. T. W., 2019. Temporal emotion-aspect modeling for discovering what students are concerned about in online course forums. Interactive Learning Environments, Volume 27(5-6), p. 598–627.
- Zied Kechaou, M. B. A. a. A. M. A., 2011. Improving e-learning with sentiment analysis of users' opinions.. 2011 IEEE global engineering education conference (EDUCON), p. 1032–1038..
- Zohair, L. M. A., 2019. Prediction of Student's performance by modelling small dataset size. International Journal of Educational Technology in Higher Education volume.