

# Predicting learning outcomes in an online learning platform

Ying Bai, Michael Bosu, Diab Abuaiadah

## Abstract

This paper indicates a comprehensive analysis of a dataset collected from the Open University Online Learning Platform, aiming to understand the relationships between various factors and learners' final results. Data preprocessing and analysis show that factors such as age, gender, different course modules, educational background, IMD band, and total clicks positively correlate with learners' performance. K-Means clustering is utilised to identify distinct learning behaviours among learners, forming learners into three groups. Furthermore, the Random Forest algorithm was used to build machine learning models based on the identified learning behaviours, achieving good prediction accuracy. The findings emphasise the importance of targeted interventions and support tailored to the specific needs of different learner groups. The contribution of this paper is that it is the first paper to use of k-means clustering in dividing the data into groups prior to using the Random forest algorithm to predicting the final result of learners at the Open University. Secondly, this is the first paper to apply Random forest algorithm to the Open University Online Learning Platform dataset with commendable outcome in the prediction of learners final results.

## Keywords

E-learning and online training, machine learning, personalized content for online learners, individualized instruction, online machine learning

## Introduction

In recent years, there has been a significant shift towards online education, particularly during the global COVID-19 pandemic where it became the predominant access to education for many learners globally. Maatuk et al. (2021) conducted a study to investigate the perspectives of students and faculty tutors regarding the challenges and issues associated with remote learning during that period. Maatuk et al. (2021) study found that that students perceive online learning platforms as beneficial for their educational journey. This is primarily because e-learning platforms offer the advantage of accessing courses and classrooms at any time and from anywhere. Consequently, it becomes crucial for providers of online training platforms to offer personalized learning pathways to accommodate learners' cognitive abilities, knowledge structures, and study paces (Maatuk et al., 2021).

Online education platforms are web-based systems that provide users with access to educational materials and enable interaction with instructors and fellow learners. These platforms have gained popularity over the past decade due to their flexibility and convenience in acquiring new skills and knowledge. Notably, various online learning platforms, such as Massive Open Online Courses (MOOCs), Learning Management Systems (LMS), Adaptive Learning Platforms, Language Learning Platforms, and Skills Training Platforms, have effectively leveraged data mining techniques to develop automatic grading and recommendation systems. By utilising intelligent algorithms, these platforms gather valuable user information, including the frequency of platform usage, accuracy in answering questions, and time spent engaging with learning materials (Despotović-Zrakić et al., 2012). Aher et al. (2013) proposed that machine learning (ML) methods, including collaborative filtering, content-based filtering, decision trees, and artificial neural networks, can be employed to process and analyse the acquired information on e-learning platforms. Machine learning, a subset of artificial intelligence, involves learning from data, identifying patterns, and making predictions. The rise of data availability and advancements in cloud computing have contributed to the rapid growth of machine learning, enabling efficient analysis of complex data and mitigation of unforeseen risks (Aher et al., 2013). Although online education offers flexibility and affordability compared to traditional on-campus learning, it does present challenges due to reduced interaction between learners and instructors (Eom & Ashill, 2016). To address this, long-term log data from online platforms can be leveraged for learner and course assessment. Machine learning algorithms can assist in predicting students at risk and estimating dropout rates through analysis of preprocessed log data (Essalmi et al., 2015).

The Felder and Silverman Learning Styles Model (FSLSM) is a well-recognized framework for identifying an individual's learning style. It assesses learners across four dimensions: information processing, input, understanding, and perception (Nafea et al., 2019). In the context of online learning, Moodle, an online learning platform, utilises the FSLSM to offer adaptive course delivery (Essalmi et al., 2015). Despotović-Zrakić et al. (2012) conducted a study focusing on data mining techniques to classify students into clusters based on the FSLSM learning styles model. Although the research primarily focused on learning styles, the results have been applied in Moodle learning management systems to enhance personalized strategies. Machine learning techniques like K-Means clustering and the Random forest algorithm can be employed to develop a model that leverages learners' profiles and online operation logs to refine personalisation strategies on a specific online learning platform. It is important to acknowledge that personalization strategies can also be influenced by other factors, such as prior knowledge and expectations (Despotović-Zrakić et al., 2012).

The objectives of this paper are to:

- identify distinct learning behaviours or learning styles among learners in an Online Learning Platform using K-means clustering.
- divide the learners into groups based on the identified learning behaviours.
- predict learner outcomes using the Random Forest algorithm.

The rest of the paper is as follows. Section 2 is the review of the literature. Section 3 discusses the data used for this research. Section 4 presents the K-Means clustering algorithm and the Random Forest algorithm that have been used to build machine learning models for this study. Section 5 presents the findings and conclusion and section 6 is the future work.

## Literature Review

Various performance metrics are used to evaluate the effectiveness of online education platforms. These metrics include course completion rate, student engagement, satisfaction, and learning outcomes (Essalmi et al., 2015). Machine learning models can leverage these metrics and learner data to make predictions and provide personalized content and feedback to learners (Eom & Ashill, 2016). Eom & Ashill (2016) conducted an empirical investigation and identified six independent variables that influence student outcomes, including course structure, instructor feedback, self-motivation, learning style, interaction, and instructor facilitation. Among these variables, instructor feedback and learning style were found to significantly impact learning outcomes, while user satisfaction could predict outcomes (Eom & Ashill, 2016). Instructor feedback encompasses cognitive, diagnostic, and prescriptive feedback delivered through various channels (Eom & Ashill, 2016).

Learning style refers to an individual's preferred method of acquiring and processing information, influenced by physiological dimensions, cognitive dimensions, and personality characteristics. Machine learning algorithms can automatically detect learning styles based on learners' sequences extracted from e-learning system log files, enabling the provision of personalized content (El Aissaoui et al., 2019). Essalmi et al. (2015) found that a single personalisation strategy may not fit all courses and teachers, leading to the investigation of a generalized approach that considers learners' profiles and appropriate personalisation parameters (Essalmi et al., 2015). Additionally, Aher et al. (2013) demonstrated the practicality of a combined algorithm, using Simple K-means clustering and Apriori-association rule algorithm, for building a course recommendation system. Khanal et al. (2020) reviewed 101 recommendation system-related papers and found that multiple algorithms are often used in a single system, making classification of machine learning algorithms challenging (Khanal et al., 2020).

Hyper-parameter optimization is crucial for improving the performance of trained machine learning models. Different techniques exist for hyper-parameter tuning, and their selection depends on specific scenarios and models (Yang & Shami, 2020). Identifying key hyper-parameters requires data mining and analysis expertise, and choosing the appropriate tuning technique for a given model or dataset is important (Yang & Shami, 2020).

Automating hyper-parameter tuning through optimal configuration improves reproducibility, model performance, and reduces human effort (Yang & Shami, 2020). The implementation of a hyper-parameter tuning approach for the machine learning models built in this research will be discussed in the future.

## Dataset Description and Analysis

The dataset used for this paper was sourced from the Open University (OU) Online Learning Platform. The Open University is a public British University that have students who are mainly off-campus and study online from everywhere in the United Kingdom (Kuzilek et al., 2017). Figure 1 shows the structure of the dataset, including course content in the courses table, assignment marks in the assessment table, students' demographics (such as location, age group, disability, education level, and gender) in the studentInfo table, forum discussion in the vle table, students' interactions with the Virtual Learning Environment (VLE) in the studentVle table and students registration information in the studentRegistration table.

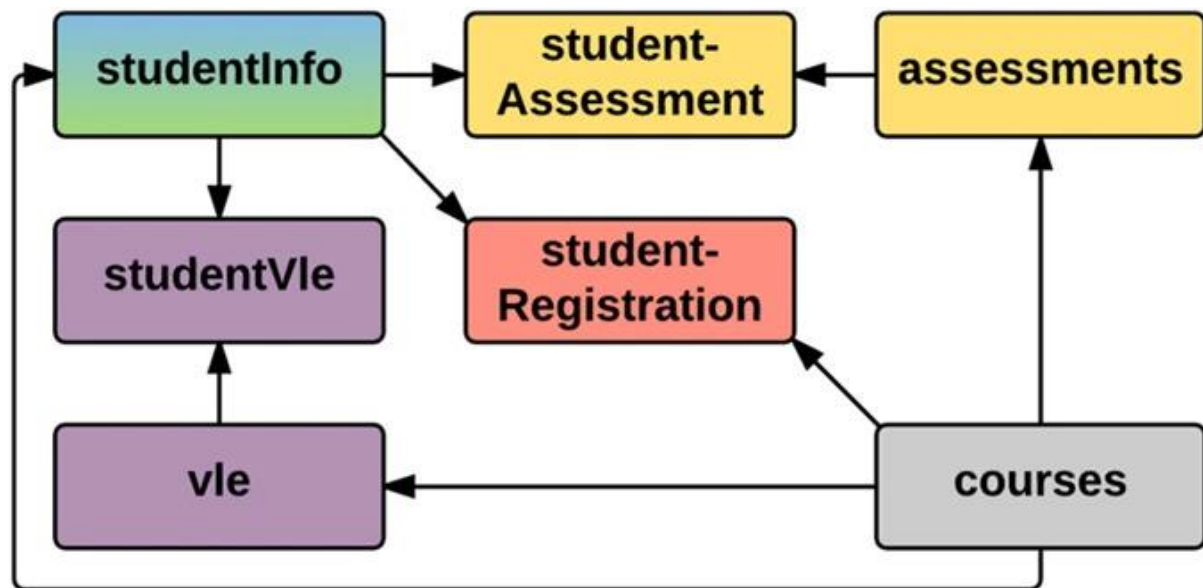


Figure 1. Data structure of Open University Online Learning Platform dataset  
Source: Kuzilek, J., Hlosta, M., & Zdrahal, Z. (2017). Open university learning analytics dataset. Scientific data, 4(1), 1-8.

Figure 2 presents the details of each field in the tables. The final data used for building the machine learning models in Section 4 is obtained by combining data from studentInfo table, studentRegistration table, studentAssessment table, assessments table, courses table and studentVle table. The final data used in building the models comprises of 24845 records of 22488 students and twenty-two (22) course modules.

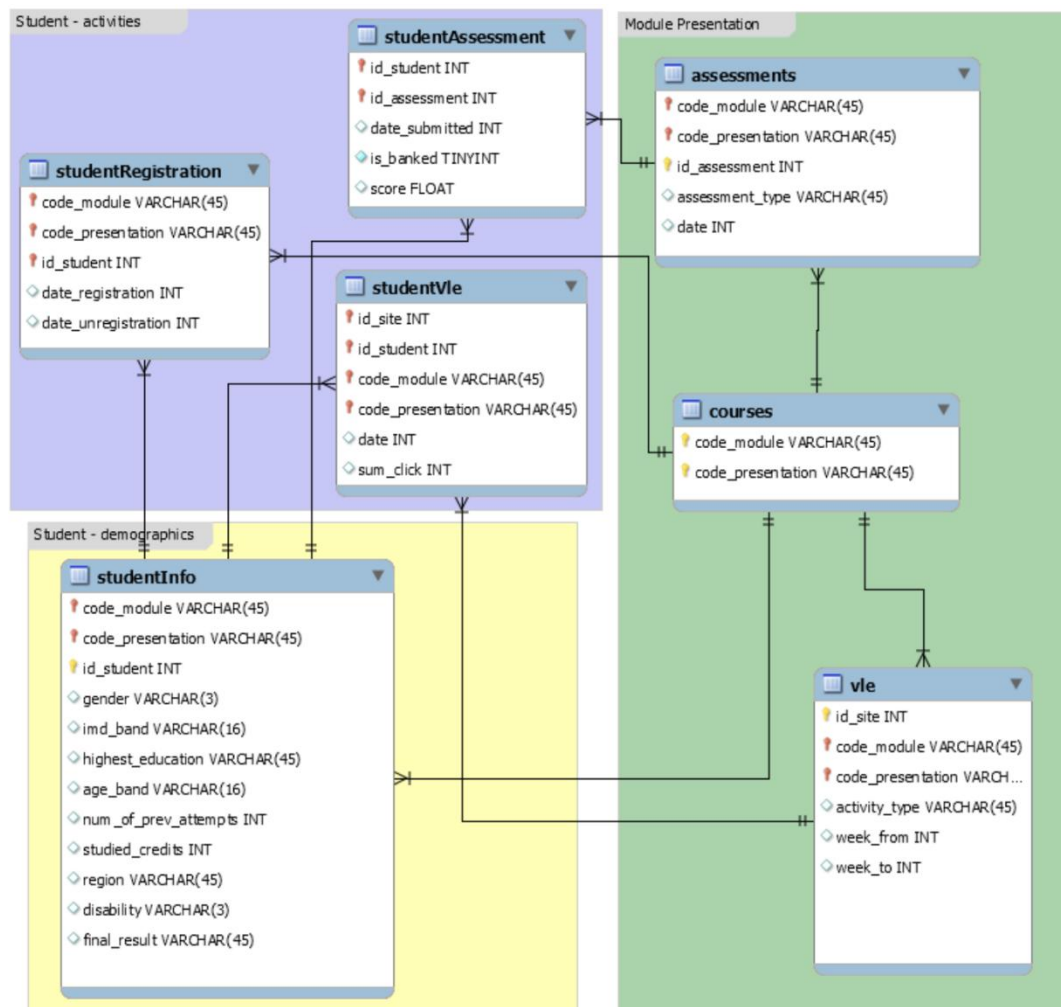


Figure 2 Data schema Source: Doijode, V., & Singh, N (n.d).

Figure 3 presents the relationship between parameters and students' final results. Parameters include age band, education background, the index of multiple deprivations (IMD) band and total clicks which refer to the cumulative number of clicks made by a student within the Open University (OU) online platform. The Index of Multiple Deprivation (IMD) is a measure used in the United Kingdom to assess relative deprivation at the small-area level. It combines various indicators across multiple domains, such as income, employment, education, health, crime, and living environment, to provide an overall measure of deprivation (Noble et al., 2006).

Preliminary analysis of the data as demonstrated by Figure 3 indicate that older learners tend to perform better than younger learners. In fact, around 40% of learners older than 55 years old received a distinction in the final results, as shown in Figure 3(a). Additionally, Figure 3(b) illustrate that learners with higher education qualifications have a higher chance to pass course and achieve distinction performance when compared to learners with lower qualifications. Figure 3(c). indicate that learners from high IMD areas have a higher chance of achieving better performance than learners from low IMD areas. Lastly, Figure 3(d) indicates that frequent interaction with the OU online platform has a positive impact on final results. Learners with a total of over 1000 clicks tend to pass courses and perform better overall.

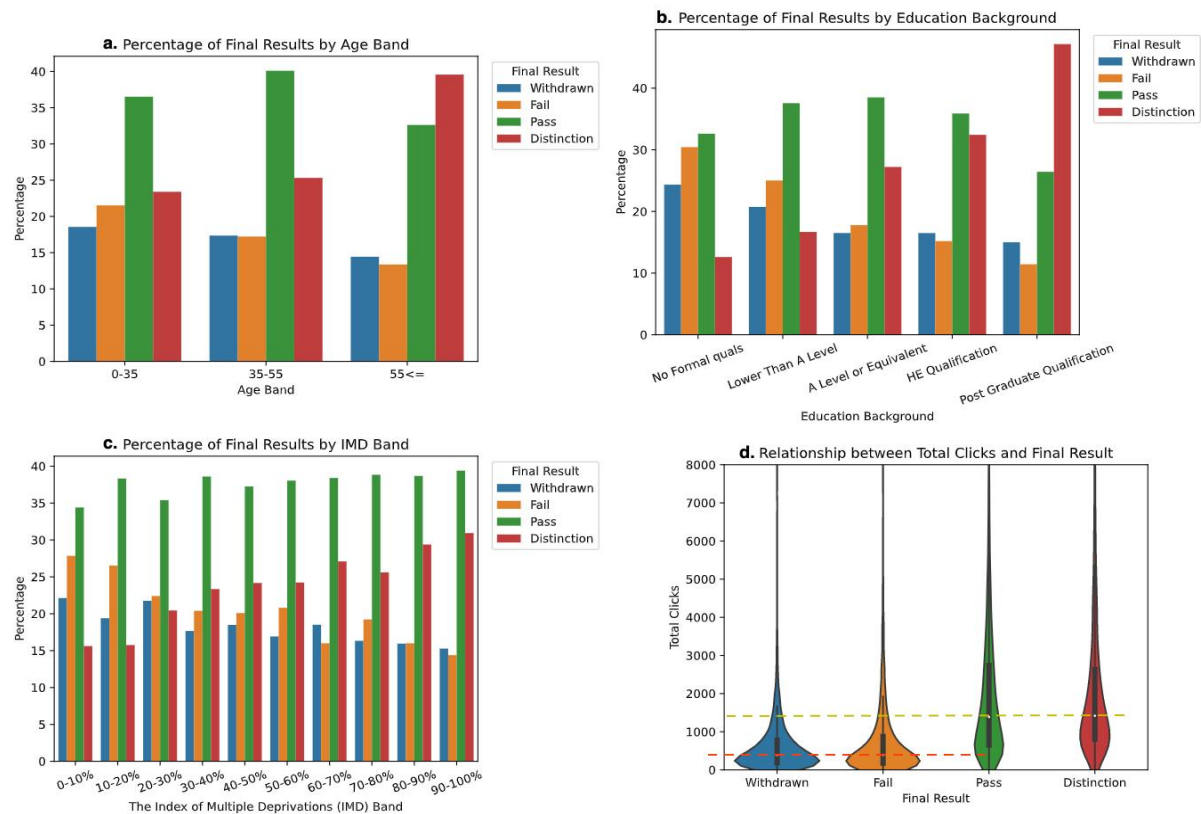


Figure 3. Preliminary results of data analysis

Figure 4 shows a positive correlation between six factors and the final results because their values in the following heat-map are more significant than 0. These factors include IMD band (0.14), age band (0.04), education background (0.17), course modules (0.022), gender (0.049), and total clicks (0.34). The parameter with the strongest positive correlation with final results is total clicks, which has a value of 0.34, the highest among the other parameters.

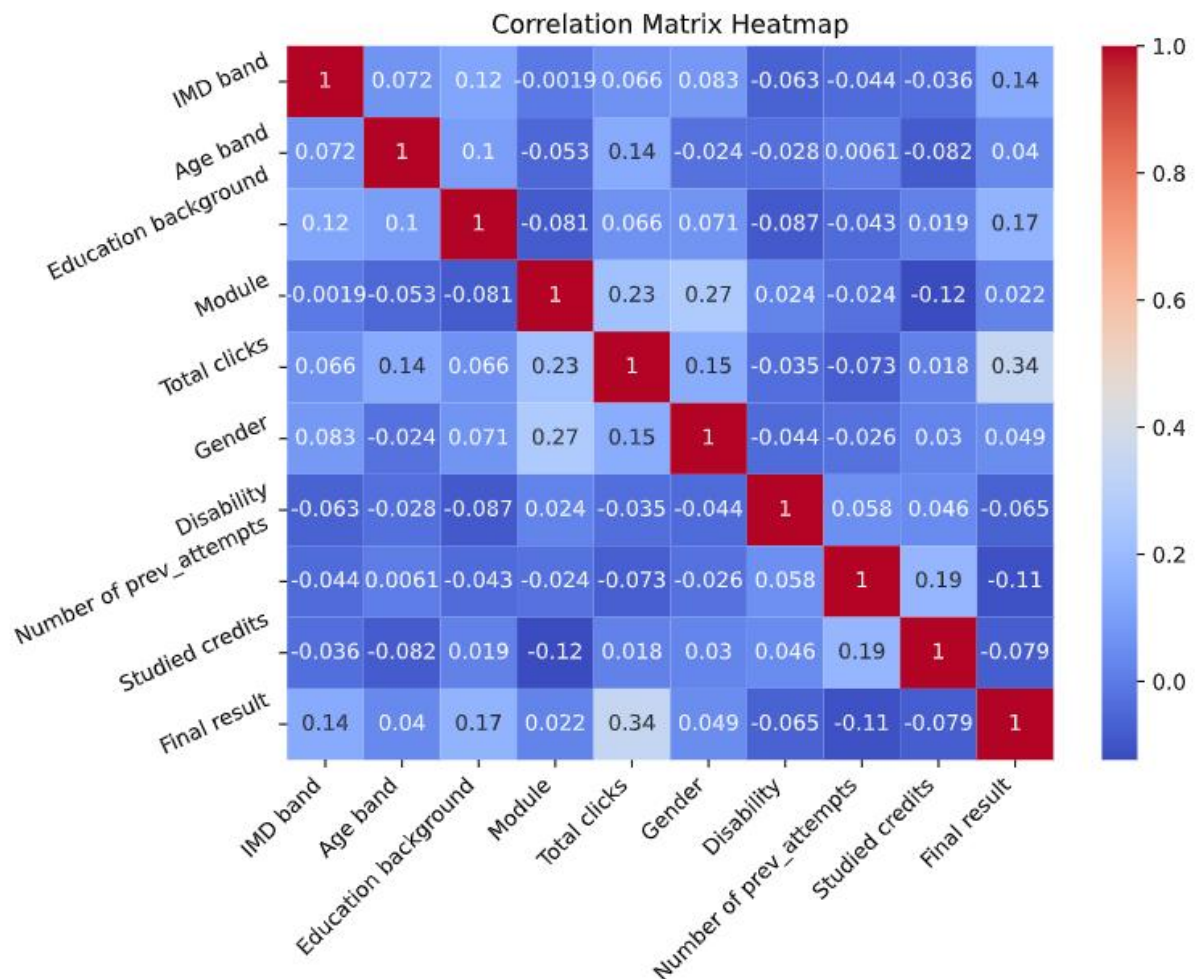


Figure 4. Correlation of age, education background, and IMD

The dataset analysis found that six factors are closely related to the learners' performance in the final results. These factors include age, educational background, IMD in the location, course, gender, and frequency of interaction with the Open University online learning platform.

## Model Development

Based on the dataset analysis in Section 3, age (it refers to "age\_band" column in the combined dataset), educational background (it refers to "highest\_education" column in the combined dataset), IMD in the location (it refers to "imd\_band" column in the combined dataset), course modules (it refers to "code\_presentation" column in the combined dataset), a learner's gender (it refers to "gender" column in the combined dataset), and frequency of interaction with the OU online learning platform (it refers to "total\_click" column in the combined dataset) are going to be used as parameters to build machine learning models to predict learners final results.

There are two fundamental types of machine learning algorithms: supervised and unsupervised. K-Means clustering is an unsupervised algorithm to group learners based on shared characteristics or patterns. By identifying clusters of learners with similar attributes, K-Means clustering helps reveal underlying structures within the data (Mahesh, 2020).

This section uses K-Means clustering to identify distinct groups in the combined dataset. Learners with similar characteristics are grouped based on factors, including age, educational background, IMD in the location, course modules, and frequency of interaction with the OU online learning platform. Then, the Random forest algorithm is utilised to build a machine learning model. The random forest offers a flexible and robust ensemble of decision trees, making it suitable for analysing complex datasets and capturing intricate relationships between variables (Mahesh, 2020). By leveraging the power of Random Forest, this model aims to provide predictions of learners' final results while considering a wider range of factors.

#### 4.1. Utilising the K-Means Clustering Methodology to classify learners

To determine the appropriate number of clusters when utilizing K-Means clustering, it is essential to address the optimal value denoted as "k." The Elbow Technique is a commonly employed method for identifying this optimal k value (Marutho et al., 2018). The Elbow Technique is based on the observation that the within-cluster sum of squares (WCSS) tends to decrease as the number of clusters increases. The WCSS calculates the sum of squared distances between each data point and the centroid of its assigned cluster (Marutho et al., 2018). The quality evaluation of a clustering algorithm, particularly in K-Means clustering, employs the Sum of Squared Error (SSE) metric. The SSE quantifies the variation or dispersion within the clusters. By minimizing the SSE, the algorithm aims to create compact and well-separated clusters (Bholowalia & Kumar, 2014).

By employing the Elbow Technique to determine the optimal k value, a plot is generated with the number of clusters (k) on the x-axis and the corresponding SSE values on the y-axis. Typically, the plot exhibits a downward trend, where the SSE decreases as k increases. However, there is a certain point beyond which increasing k does not significantly reduce the SSE (Bholowalia & Kumar, 2014). The "elbow" in the plot represents the point where the SSE starts to level off, forming a bend similar to an elbow. This bend indicates a diminishing return in SSE reduction beyond that point. The optimal k value is frequently selected at the elbow as it signifies a suitable trade-off between capturing meaningful patterns within the data and avoiding excessive complexity. Figure 5 illustrates SSE values on the y-axis and the number of clusters on the x-axis. It is evident that the optimal k value is 3, as indicated by the elbow point in the plot.

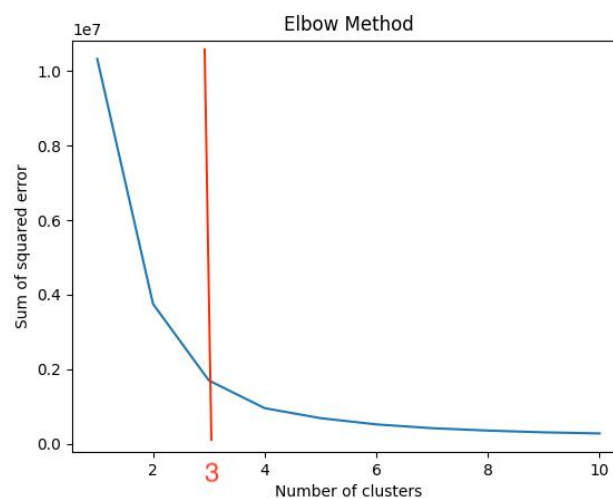


Figure 5. Using the Elbow technique to the find optimal k value



Based on the implementation of the Elbow technique, it was determined that the optimal number of clusters (k) is 3. As a result, the learners in the combined dataset have been divided into three distinct groups, as shown in Figure 7. In Figure 6, for k=3, in Group 1, the group size is 5409, and the overall performance is "Pass". The size of Group 2 is 1066 and the overall performance is also "Pass" whilst the size of Group 3 is 18370 and the performance comprises of "Fail" or "Withdrawn" or "Pass". Also, the group sizes and the overall performances for group numbers 2 and 3 are 1066, 18370 and "Pass", and "Fail or Withdrawn or Pass".

a. Learner's performance index		b. K=3		
Total Score	Learner's performance	Group #	Group size	Overall learner's performance
below 4000	Fail or Withdrawn	1	5409	Pass
4000-8000	Pass	2	1066	Pass
8000 and above	Distinction	3	18370	Fail or Withdrawn or Pass

Figure 6. Group sizes and types

Figure 7 provides insights into the performance of learners in three different groups. It reveals that learners in Group 1 and Group 2 learners have a higher likelihood of successfully passing their courses, with Group 2 showing promising results. However, a significant proportion of learners in Group 3 are at a higher risk of withdrawing or failing their courses.

These findings highlight the existence of distinct learner profiles within the dataset, indicating the need for targeted interventions and support tailored to each group's specific needs and challenges. For example, to address the higher withdrawn and failure rate observed in Group 3, personalised support should be provided by staff of OU online learning platform to assist and motivate learners in this group. By recognising and addressing the unique requirements of each group, educational institutions can enhance learner outcomes and improve overall success rates.

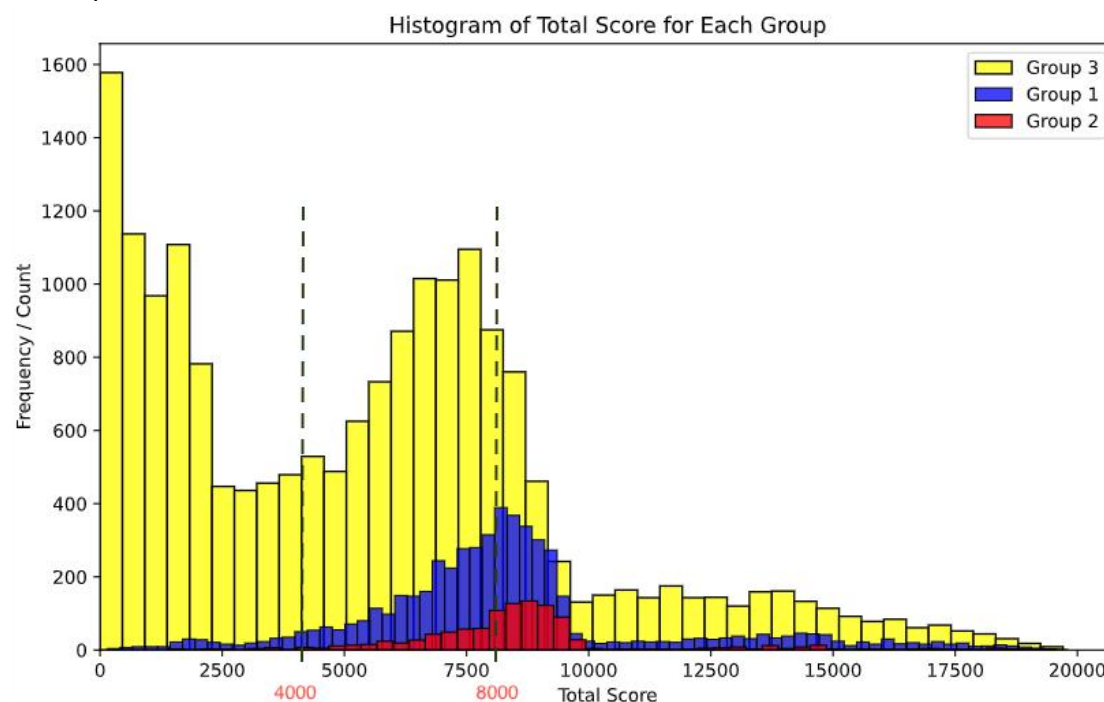


Figure 7 Dividing learners into three groups

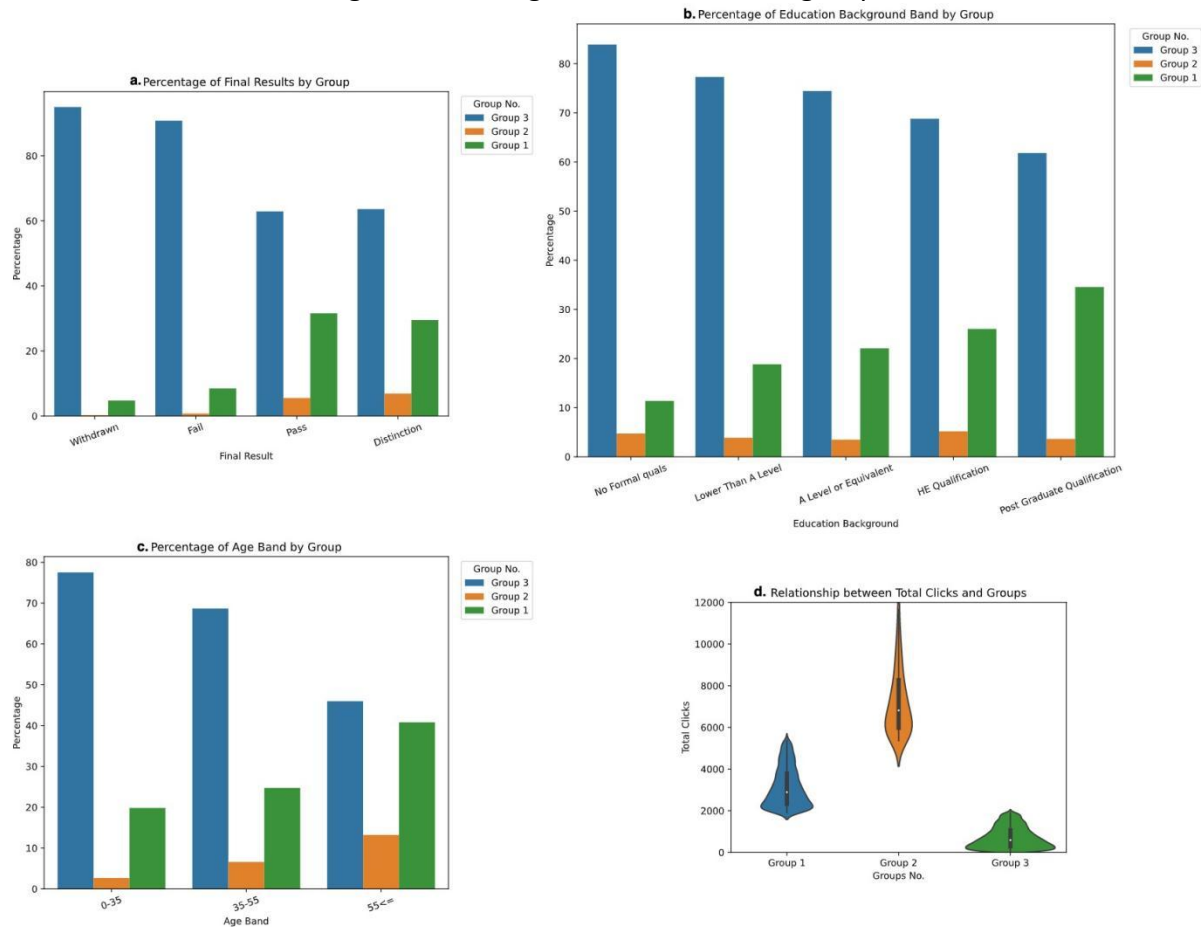


Figure 8. Comparison of three groups

Figure 8 illustrates that almost all learners in Group 2 pass courses, and the learners in Group 1 have relatively high possibility to pass courses compared with learners in Group 3. Figures 8b, 8c and 8d show that most learners in Group 1 have A level or higher education qualification, are usually more than 55 years old, and have moderate interaction frequency with the online learning platform. In contrast, most learners in Group 3 have lower than A level qualifications, are less than 35 years old, and have the lowest interaction frequency with the online learning platform.

## 4.2 Utilising Random Forest to Build Machine Learning (ML) Model

In subsection 4.1, learners were divided into three groups using K-means clustering, which means learners in the combined dataset have distinct characteristics, and by applying K-means clustering, learners with similar characteristics are grouped together. In this subsection, the Random forest algorithm is utilised to build a machine learning model to predict learners' outcomes (final results) based on their learning characteristics identified through K-means clustering in subsection 4.1.

The Random forest algorithm is a classification and regression algorithm to build a forest containing many individual decision trees. It generates reasonable predictions across various parameters while requiring little configuration. Moreover, when processing extensive data, several missing data will not affect its accuracy significantly (Boulesteix et al., 2012).

The parameters used in Random forest model training include two main parts. The first part are the relevant factors found in section 3; age, educational background, IMD in the location, course, gender, and frequency of interaction with the OU online learning platform.

Moreover, the second part is the factors that correlate less with learners' outcomes, including disability, number of previous attempts of a course, and learners studied credits according to the heatmap in Figure 4.

The table of courses in the datasets contains seven different programmes. The domains of these programmes are social science and STEM (Science, Technology, Engineering, and Mathematics). The assessment table shows three assessment types: tutor-marked assessment, computer-marked assessment, and final examination (Kuzilek et al., 2017).

Figure 9 illustrates the visualisation of a decision tree generated by the Random Forest algorithm in the context of creating the machine learning model. The decision tree begins with a topmost node, which is split based on the most significant feature. In this particular tree, the topmost node represents the number of interactions a learner has had with the OU online learning platform (`total_click`). The internal nodes of the tree correspond to binary decisions made based on specific features. These decisions guide the data flow down the tree, ultimately leading to the terminal nodes. The terminal nodes represent the predicted outcomes or classes, such as "Pass" or "Fail," based on the given features and their associated paths within the decision tree. Implementation of the K-means clustering and the Random Forest was conducted in Jupyternotebook in Python 3. X programming environment, and the decision tree in Figure 9 is virtualised in "Graphviz Online" (<https://dreampuf.github.io/GraphvizOnline>).

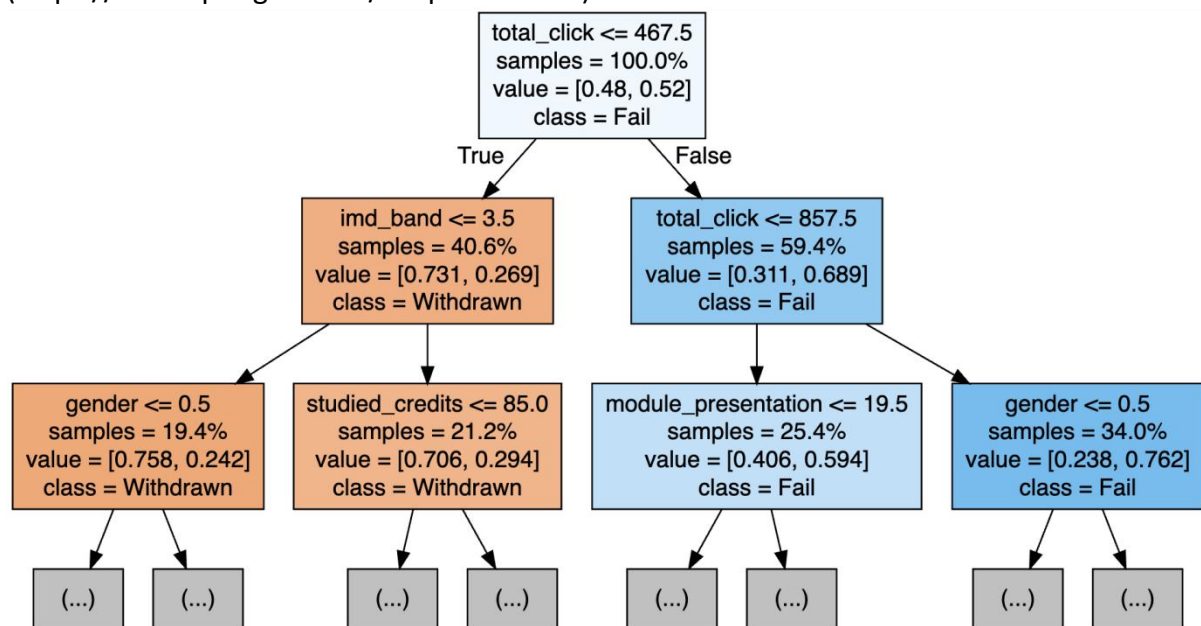


Figure 9 An example of decision trees

Based on identifying distinct learners' characteristics through K-means clustering in subsection 4.1, the original combined dataset was divided into three individual datasets. Using the Random Forest algorithm, these datasets are used to build three separate ML models (Model 1, Model 2, Model 3). Among the three models, Model 2, which is based on learners in Group 2, achieved the highest prediction accuracy of 93%, indicating that this group exhibits more consistent and distinct learning characteristics.

On the other hand, Model 3, which is based on learners in Group 3, demonstrated lower prediction accuracy than the other two models. The findings suggest that learners in Group 3 exhibit a more comprehensive range of learner characteristics, making it more challenging to predict their outcomes accurately.

Furthermore, Model 4 was built using the original combined dataset without applying K-means clustering. It achieved a prediction accuracy of 79%, indicating that the Random Forest algorithm can effectively classify datasets even without the explicit grouping provided by K-means clustering.

Model #	Data size (rows)	Train data size (rows) 80%	Test data size (rows) 20%	Prediction accuracy
1	5409	4327	1082	89%
2	1066	852	214	93%
3	18370	14696	3674	76.3%
4	24845	19876	4969	79%

Figure 10 Results of ML models built by Random forest algorithm

K-means clustering and Random Forest can be used for dataset classification. However, K-means clustering is an unsupervised algorithm that does not require predefined labels or target variables. It aims to identify natural groupings or clusters within the data based on similarities in features or variables. While K-means clustering can be useful for exploratory analysis and understanding data patterns, it does not directly predict a target variable (Mahesh, 2020). On the other hand, Random Forest is a supervised algorithm that requires predefined labels and a designated target variable. It constructs an ensemble of decision trees, where each tree is trained on a subset of the data with random feature selection (Boulesteix et al., 2012). Random Forest can predict the target variable based on the learned patterns and relationships in the training data, as shown in Figure 10.

Therefore, in this specific scenario, K-means clustering is more suitable for discovering underlying patterns and grouping similar learners. At the same time, Random Forest is a supervised algorithm that enables the prediction and classification of target variables (final results). Combining K-mean clustering with the Random Forest algorithm can improve prediction accuracy. In Figure 10, the average accuracy of Models 1, 2 and 3 is 86.1% which is higher than Model 4, 79%, built by the original combined data without applying K-means clustering. It indicates that combining K-means clustering with the Random forest can efficiently improve the ML models' accuracy.

## Conclusion

In conclusion, this paper focuses on collecting, preprocessing, and analysing Open University Online Learning Platform dataset. Data analysis and preprocessing showed that age, gender, educational background, IMD band, and frequency of interaction with the OU online learning platform (total clicks) positively correlated with learners' final results. K-Means clustering was then employed to identify distinct learning behaviours among learners, forming three groups. Furthermore, the Random Forest algorithm was used to build machine learning models based on learner groups with the identified learning behaviours.

Comparisons between the models built with and without K-Means clustering showed the effectiveness of clustering for improved classification. According to Figure 10, the accuracy of the model without applying K-means clustering is 79%, whereas the accuracy of the model with applying K-means clustering is 86.1%. To the best of our knowledge the contribution of this paper is that it is the first paper to use of k-means clustering in dividing the data into groups prior to using the Random forest algorithm to predicting the final result of learners at the Open University. Secondly, this is the first paper to apply Random forest machine learning algorithm to this specific datasets with superior outcome in the prediction of learners final results.

## Future Work

The next stage is finding an approach to set hyperparameters automatically, learning from data instead of manually setting them and executing hyper-parameter optimization techniques, such as Bayesian optimization hyperband (BOHB), genetic algorithms (GA), and particle swarm optimization (PSO), to improve the performance of the machine learning models built above. Therefore, future work could explore additional machine learning techniques, consider more features in the analysis, and evaluate the effectiveness of personalized interventions in improving learner outcomes.

## References

- Aher, S. B., & Lobo, L. M. R. J. (2013). Combination of machine learning algorithms for recommendation of courses in E-Learning System based on historical data. *Knowledge Based Systems*, 51, 1–14. <https://doi.org/10.1016/j.knosys.2013.04.015>
- Angeline, D. M. D. (2013). Association Rule Generation for Student Performance Analysis using Apriori Algorithm. *The SIJ Transactions on Computer Science Engineering & Its Applications (CSEA)*, 01(01), 16–20. <https://doi.org/10.9756/sijcsea/v1i1/01010252>
- Bholowalia, P., & Kumar, A. (2014). EBK-Means: A Clustering Technique based on Elbow Method and K-Means in WSN. *International Journal of Computer Applications*, 105(9), 975–8887. <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=5771aa21b2e151f3d93ba0a5f12d023a0bfcf28b>
- Boulesteix, A.-L., Janitza, S., Kruppa, J., & König, I. R. (2012). Overview of random forest methodology and practical guidance with emphasis on computational biology and bioinformatics. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(6), 493–507. <https://doi.org/10.1002/widm.1072>
- Despotović-Zrakić, M., Jeremic, V., Bogdanović, Z., Barać, D., & Krco, S. (2012). Providing Adaptivity in Moodle LMS Courses. *Educational Technology & Society*, 15(1), 326–338. [http://www.ifets.info/journals/15\\_1/28.pdf](http://www.ifets.info/journals/15_1/28.pdf)
- Doijode, V., & Singh, N. (n.d.). Predicting student success based on interaction with virtual learning environment. Retrieved May 26, 2023, from [https://www.lexjansen.com/sesug/2016/EPO-271\\_Final\\_PDF.pdf](https://www.lexjansen.com/sesug/2016/EPO-271_Final_PDF.pdf)
- El Aissaoui, O., El Madani El Alami, Y., Oughdir, L., & El Alloui, Y. (2019). A hybrid machine learning approach to predict learning styles in adaptive E-learning system. In *Advanced Intelligent Systems for Sustainable Development*

- (AI2SD'2018) Volume 5: Advanced Intelligent Systems for Computing Sciences (pp. 772-786). Springer International Publishing.
- Eom, S. B., & Ashill, N. (2016). The Determinants of Students' Perceived Learning Outcomes and Satisfaction in University Online Education: An Update\*. *Decision Sciences Journal of Innovative Education*, 14(2), 185–215.  
<https://doi.org/10.1111/dsji.12097>
- Essalmi, F., Ayed, L. J. B., Jemni, M., Graf, S., & K. (2015). Generalized metrics for the analysis of E-learning personalization strategies. *Computers in Human Behavior*, 48, 310–322. <https://doi.org/10.1016/j.chb.2014.12.050>
- Khanal, S. S., Prasad, P. W. C., Alsadoon, A., & Maag, A. (2020). A systematic review: machine learning based recommendation systems for e-learning. *Education and Information Technologies*, 25(4), 2635–2664.  
<https://doi.org/10.1007/s10639-019-10063-9>
- Kuzilek, J., Hlosta, M., & Zdrahal, Z. (2017). Open University Learning Analytics dataset. *Scientific Data*, 4, 170171. <https://doi.org/10.1038/sdata.2017.171>
- Maatuk, A. M., Elberkawi, E. K., Aljawarneh, S., Rashaideh, H., & Alharbi, H. (2021). The COVID-19 pandemic and E-learning: challenges and opportunities from the perspective of students and instructors. *Journal of Computing in Higher Education*, 34(1). <https://doi.org/10.1007/s12528-021-09274-2>
- Mahesh, B. (2020). Machine learning algorithms-a review. *International Journal of Science and Research (IJSR)*. [Internet], 9, 381-386.
- Marutho, D., Handaka, S. H., & Wijaya, E. (2018). The determination of cluster number at k-mean using elbow method and purity evaluation on headline news. In *2018 international seminar on application for technology of information and communication*. 533-538. IEEE.
- Nafea, S. M., Siewe, F., & He, Y. (2019). On Recommendation of Learning Objects using Felder-Silverman Learning Style Model. *IEEE Access*, 1–1.  
<https://doi.org/10.1109/access.2019.2935417>
- Noble, M., Wright, G., Smith, G., & Dibben, C. (2006). Measuring multiple deprivation at the small-area level. *Environment and planning A*, 38(1), 169-185.

## Authors

Ying Bai is a student in the Centre for Information Technology at Waikato Institute of Technology, Te Pūkenga. [yinbai10@student.wintec.ac.nz](mailto:yinbai10@student.wintec.ac.nz)

Michael Franklin Bosu is a Lecturer in the Centre for Information Technology at Waikato Institute of Technology, Te Pūkenga. [michael.bosu@wintec.ac.nz](mailto:michael.bosu@wintec.ac.nz)

Diab Abuaiadah is a Lecturer in the Centre for Information Technology at Waikato Institute of Technology, Te Pūkenga. [diab.abuaiadah@wintec.ac.nz](mailto:diab.abuaiadah@wintec.ac.nz)