

# AI Based Student Performance Prediction In Online Education

Harshal Nimbraj Patil  
MCA(Computer Science)  
JSPM University Wagholi, Pune

Dr. Anita Pisote  
Assistant Professor  
JSPM University Wagholi, Pune

## Abstract

This research looks at using artificial intelligence (AI) to predict student performance in online learning environments. It's important for the digital education space. Most platforms have high attrition and academic interventions can go unnoticed for a long time on platforms that allow aztech Technology. So we built a predictive framework using Machine Learning (ML) algorithms. Random forest and long short term memory networks It helps identify at risk students in the early phase of a semester early enough to act.

Traditional models usually only use past grades. Our approach adds multi-factorial data features. At the clickstream level we include duration of video watches and participation in forum discussions. We also use demographic data and a temporal axis, to name a few. Model performance was excellent with strong predictive accuracy. And by applying Explainable AI techniques such as SHAP values we were able to extract transparent insights into potential behavioral triggers of academic success or failure. Using SHAP we pulled out clear, interpretable signals tied to specific behaviors and patterns.

**Keywords:** AI in Education, Student Performance Prediction, Machine Learning, Learning Analytics, Online Education, Explainable AI.

## I. INTRODUCTION

Education is going digital fast. Around the world the traditional classroom has turned into a busy, data-filled virtual space. MOOCs and institutional Learning Management Systems, LMS, are now the main way higher education happens. They generate huge volumes of digital footprints. But this shift

brought a stubborn problem. Dropout rates are high. Engagement often drops compared with face-to-face classes. So predicting student performance with Artificial Intelligence, AI, is not just a tech project. It's become a pedagogical necessity to help students succeed and to keep institutions efficient.

Schools used to judge performance with summative tests like mid-terms and final exams. Those after-the-fact checks usually come too late to help. When a student fails a major exam the chance for corrective support is often gone. AI changes that. It makes learning analytics proactive. AI looks at real-time data. Things like login frequency and video watch time and forum activity. Models can spot subtle disengagement patterns teachers miss. An Early Warning system can flag at risk student in the first few weeks so educators can give tailored support.

Student behavior online is messy. Complex you need advanced Machine Learning, ML, techniques. Simple statistical models can spot correlations. Deep learning architecture like recurrent neural network, and long short term memory, units handle the sequential nature of learning. They treat education like a journey. A sudden dip in activity in week five can be more telling than a student's average performance over a year. Explainable AI, XAI, tackles the black box problem. It makes predictions transparent and actionable for instructors.

This research develops and validates an AI-based framework to predict student outcomes in online education. It compares demographic factors with behavioral engagement metrics to see which matter more. The aim is a model that's accurate and ethically sound. The study ties raw data to educational psychology. It wants to give institutions tools to build a more inclusive and responsive digital learning ecosystem. With AI we move closer to a future where every student's unique learning path.

## I. Literature Review

Student performance prediction is a big research area now in Artificial Intelligence and machine learning. It's also a focus in Educational Data Mining (EDM). Online education grew fast. Learning Management Systems and Massive Open Online Courses spread widely. Schools and universities use predictive analytics to spot at-risk students and try to boost learning outcomes. That raises academic performance.

### A. Student Performance Prediction

A lot of research has tried to predict student academic performance using data-driven methods. Early work used traditional statistical techniques. Lately it's shifted toward machine learning and AI. Romero and Ventura [1] say educational data mining helps analyse student data and find hidden patterns you can use to predict academic success or failure. They examine static features. Things like demographic data and prior academic records, with socio-economic background often included. They also look at dynamic features. Attendance and assignment submissions plus online activity and engagement metrics. Recent work leans more on dynamic behavioral data, especially in online courses. Hussain [2] showed that adding student interaction data Login frequency and time spent on learning platforms, along with participation in online discussions significantly improves prediction accuracy. Kotsiantis [3] pointed out that early prediction of student performance lets educators act sooner. They can provide tutoring tailored to students and targeted interventions to cut dropout rates and improve outcomes.

### B. AI in Education

Artificial Intelligence has changed education a lot. It lets systems adapt to each learner. AI-driven tools in education cover a wide range. From personalized learning to recommendation engines and tools that predict outcomes. Holmes et al. [4] showed AI can boost teaching and learning. It gives adaptive content and automated feedback. And it monitors performance in real time. They are

used a lot online because courses create huge amounts of student interaction data. Baker and Inventado [5] noted AI is widely used in EDM to

model student behaviour and predict outcomes. AI models find patterns that traditional methods miss. That helps with decisions in education. Online, AI looks at engagement and predicts performance from behaviour. That lets institutions set up early warning systems to catch students who might underperform or drop out.

### C. Models Used in Previous Work

Researchers have proposed many machine learning and deep learning model to predict student performance. They vary a lot. Some are simple and easy to interpret. Others are complex and often more accurate.

#### 1. Traditional Machine Learning Models

Many studies used classic machine learning algorithms for prediction. Linear regression is often used for continuous outcomes. Logistic regression is used for categorical ones. Decision trees are common. Random forests too. They handle non-linear relationships and feature interactions. Kotsiantis [3] said decision trees give interpretable results. So they're a good fit for education, where you want to see why the model decides. Random forest is an ensemble method. It has shown better accuracy and more robustness than single classifiers. Support vector machines and K-Nearest neighbours have also been used for classification. They perform satisfactorily on different datasets [6].

#### 2. Neural Networks and Deep Learning

As computers got faster, deep learning took off in recent years. Artificial Neural Networks can pick up complex patterns in big datasets. More advanced architectures like long short-term Memory networks work well for sequential data. Think student activity logs over time. Multiple studies show deep learning beats traditional machine learning on large-scale and time-series data in many cases [7]. But they usually need big datasets and lots of compute and training time. That can limit their use in some educational settings, though.

#### 3. Ensemble and Hybrid Models

Researchers have tried ensemble and hybrid approaches to improve prediction performance. Ensemble methods mix approaches like bagging and boosting. Or they stack models to combine predictions and get higher accuracy. Zhang et al. [8] showed hybrid models that pair decision trees with neural networks can significantly improve prediction accuracy over single models. They use

each algorithm's strengths and cut down on their weaknesses.

#### D. Research Gaps

There has been a lot of research in this area. But gaps still exist.

##### 1. Limited Use of Real-Time Data

Most studies use historical datasets. They rarely tap real-time streams from online platforms. So immediate feedback and timely interventions are hard.

##### 2. Lack of Generalization Across Datasets

Models get trained on specific datasets. They often fail to generalize to different institutions or learning environments. That cuts down on practical use.

##### 3. Underutilization of Behavioral Data

Online learning platforms generate lots of interaction data. Many studies don't make full use of that information for predictions.

##### 4. Data Quality and Imbalance Issues

Educational datasets often have missing values, noise and class imbalance. That can hurt model performance.

##### 5. Model Interpretability

Deep learning models can be accurate. But they act like black boxes. Educators find it hard to interpret results or trust the predictions.

##### 6. Integration with Learning Management Systems (LMS)

Few predictive models are actually implemented inside real-time LMS environments. So their usability in real settings is limited.

##### 7. Privacy and Ethical Concerns

Using student data raises privacy, security and ethical worries. Many studies don't address these issues adequately.

This section describes the overall approach used to develop the AI-based student performance prediction system. The methodology covers data collection and preprocessing, how we select features, model development and evaluation..

## I. Methodology

### A. Data Collection

We pulled data from online learning platforms and public education datasets, like the UCI Machine Learning Repository by Cortez and Silva [9]. This dataset includes all sorts of student info: demographics, grades, and how they behave online. So, you'll see stuff like age, gender, how often they show up, assignment and quiz scores, how often they log in, time spent on the platform, and how much they join in on discussion forums.

#### 1. Data collection from online education platforms

#### 2. Data pre-processing and cleaning

#### 3. Feature selection and transformation

#### 4. Model training using machine learning algorithms

#### 5. Performance evaluation and comparison

### B. Data Pre-processing

First, we needed to clean things up before the data could work with any machine learning algorithms. For missing values, we filled them in with the mean or mode, depending on the type. Duplicates and noisy records gone. We took care of categorical stuff like gender or course type by translating them into numbers using label and one-hot encoding. And for all those number features, we normalized them—so everything played nice on the same scale. You'll find approaches like these in work by Jiawei Han, Micheline Kamber, and Jian Pei [10].

### C. Feature Selection

Next, we had to figure out which factors mattered most for how students perform. We ran correlation analysis and used feature importance rankings from the Random Forest algorithm (thanks to Leo Breiman [11]) to drop irrelevant stuff and shrink things down. We kept important indicators—like attendance, assignment completion, quiz scores, and engagement metrics—so our models would have the essentials.

### D. Model Development

We tried out several machine learning and deep learning models. For traditional approaches, we used Logistic Regression, Decision Tree, Random Forest, and Support Vector Machine—all through Scikit-learn (developed by Fabian Pedregosa and team [13]). On the deep learning side, we built Artificial Neural Networks and LSTM models (after the work of Sepp Hochreiter and Jürgen Schmidhuber [12]) in Tensor Flow, hoping to catch deeper, more complex patterns in student activity.

#### E. Model Training and Evaluation

The dataset split was simple: 80% for training, 20% for testing. We trained the models, then checked their accuracy, precision, recall, F1-score, and looked at confusion matrices to see where things went right—or not. We also used cross-validation to guard against over fitting and make the models more reliable.

#### F. Explainable AI Integration

To make our predictions less of a black box, we added Explainable AI tools specifically SHAP, developed by Scott M. Lundberg and Su-In Lee [14]. SHAP values break down how much each feature pushes the outcome one way or another, so we could actually explain what the model was “thinking.” That way, we got clearer insights into what drives student learning and performance.

#### G. Tools and Technologies

The proposed system was developed with these tools

Programming Language - Python

Libraries -

NumPy and Pandas for data manipulation

Scikit -Learn for ML models

Matplotlib and Seaborn for visualization

Development environment - Jupyter Notebook or Google Colab

#### I. Implementation

The proposed AI-based student performance prediction system is implemented in Python because it's simple, flexible and has strong library support for data science and machine learning. Development happens in an interactive

environment like Jupyter Notebook or Google Colab. That makes running, testing and visualizing results easy. For data handling and pre-processing we use NumPy and Pandas. NumPy handles numerical computations. Pandas gives data structures for organizing and analysing datasets. These tools help clean the dataset, handle missing values and turn raw data into a format suitable for models training. we build machine learning models with Scikit-learn. It provides efficient implementations of classification and regression algorithms. Model like Logistic Regression and Decision tree are trained and tested. Also Random Forest and Support Vector Machine. We also use Tensor Flow to implement deep learning models such as Artificial Neural Networks. These can capture more complex relationships in the data.

#### II. Conclusion

We built an AI system to predict student performance in online courses. We tried traditional machine learning model. Logistic Regression, Decision Tree, Random Forest. Support Vector Machine. We also used deep learning, with Artificial Neural Network. The system analysed student data. It looked at academic records plus behavioural info and signals of engagement. The methodology was hands-on first data pre-processing. Then feature selection and model training. We evaluated the models with standard performance metrics to see how reliable the predictions were Random Forest came out on top. It scored highest on accuracy, precision, recall and F1-score. It handled complex nonlinear relationships and reduced overfitting because it uses ensemble learning. So it's a good fit for predicting student outcomes in a changing online learning setting. Deep learning showed promise for spotting intricate patterns. But Random Forest was easier to interpret and stayed strong across different datasets. What this means for educators and institutions. The predictive system can be added to online platforms or Learning Management Systems (LMS). It can flag students at risk of underperforming. Then instructors can step in with timely help. Offer personalized guidance or set up targeted tutoring. It can also support adaptive learning strategies. The model can help with curriculum planning, resource allocation and guide policy decisions. That can boost retention and engagement while improving academic outcomes.

This work shows how AI-driven approaches could change online education. They provide actionable insights that support both instructors and learners.

## VI. References

- [1] C. Romero and S. Ventura, —Educational Data Mining: A Review of the State of the Art,|| IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews), vol. 40, no. 6, pp. 601–618, Nov. 2010.
- [2] M. Hussain, W. Zhu, W. Zhang, and S. M. R. Abidi, —Student Engagement Predictions in an e-Learning System and Their Impact on Student Course Assessment Scores,|| Computational Intelligence and Neuroscience, vol. 2018, pp. 1–21, 2018.
- [3] S. B. Kotsiantis, —Use of Machine Learning Techniques for Educational Proposes: A Decision Support System for Forecasting Students' Grades,|| Artificial Intelligence Review, vol. 37, no. 4, pp. 331–344, May 2012.
- [4] W. Holmes, M. Bialik, and C. Fadel, Artificial Intelligence in Education: Promises and Implications for Teaching and Learning. Boston, MA, USA: Center for Curriculum Redesign, 2019.
- [5] R. S. Baker and P. S. Inventado, Educational Data Mining and Learning Analytics,|| in Learning Analytics, J. A. Larusson and B. White, Eds. New York, NY, USA: Springer, 2014, pp. 61–75.
- [6] V. Vapnik, The Nature of Statistical Learning Theory. New York, NY, USA: Springer, 1995.
- [7] Y. LeCun, Y. Bengio, and G. Hinton, Deep Learning,|| Nature, vol. 521, no. 7553, pp. 436–444, May 2015.
- [8] Y. Zhang, S. Oussena, T. Clark, and H. Kim, Use Data Mining to Improve Student Retention in Higher Education – A Case Study,|| in Proc. 12th International Conference on Enterprise Information Systems, Funchal, Portugal, 2010, pp. 190–197.
- [9] P. Cortez and A. Silva, —Using Data Mining to Predict Secondary School Student Performance,|| in Proc. 5th Annual Future Business Technology Conference, Porto, Portugal, 2008, pp. 5–12.
- [10] J. Han, M. Kamber, and J. Pei, Data Mining: Concepts and Techniques, 3rd ed. Waltham, MA, USA: Morgan Kaufmann, 2011.
- [11] L. Breiman, —Random Forests,|| Machine Learning, vol. 45, no. 1, pp. 5–32, Oct. 2001.
- [12] S. Hochreiter and J. Schmidhuber, —Long Short-Term Memory,|| Neural Computation, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [13] F. Pedregosa et al., —Scikit-learn: Machine Learning in Python,|| Journal of Machine Learning Research, vol. 12, pp. 2825–2830, 2011.
- [14] S. M. Lundberg and S.-I. Lee, —A Unified Approach to Interpreting Model Predictions,|| in Advances in Neural Information Processing Systems (NeurIPS), Long Beach, CA, USA, 2017, pp. 4765–4774.

