

Chapter 1

Introduction

1.1 Introduction

In recent years, programming has emerged as a critical skill in the digital world, but many beginners still find it a bit tricky. The reason may be attributed to one significant challenge with a personal learning strategy that serves unique needs and learning styles for individual beginners. The classical ways in use of learning approaches along with those traditional online coding sites kept along the same approach, so pupils at different levels of competencies—pacing, average, or fast moving-get equal projects. This brings about dissatisfaction, disengagement, and ultimately dwindling interest in programs.

Most of the research done in programming education has been on finding and correcting faults in students' contributions. However, these methods do not include the need for adaptive task sequencing, which is important for promoting ongoing development. The existence of online coding platforms, that are hugely used by institutions and instructors, offers the potential for automating feedback and working allocation based on learners' performance. Despite this promise, however, there is still a need in the proper use of data from these platforms to personalize learning trajectories.

The presented work on this hybrid learning paradigm aims to assess the learner's contribution as well as the performance through machine learning techniques. In that case, the model relies upon algorithms, and the list includes decision tree, KNN, XGBoost, and random forest with which it is trying to gauge the learners' level and suggest assignments accordingly based upon that. This is developed under the AECW by addressing problems of complexity issues and class imbalance problems for the tasks.

The proposed approach attempts to reform programming education by offering a scalable solution for adaptive learning. It enhances learners by offering personalized assignments and equips teachers with actionable knowledge about

individual and group performance. This is in line with the trend of increasing demand for data-informed approaches in education and lays the foundation for future expansion in computer programming instruction.

1.2 Overview of Online Programming Education

The COVID-19 epidemic prompted a paradigm shift in education by accelerating the adoption of online platforms as the primary mode of teaching and learning. With all that comes along with the adaptation to digital tools, no exception is made with programming education. Use of online platforms such as Zoom, Microsoft Teams, Google Meet, and coding websites specific ones such as HackerRank, Coursera and LeetCode spread wide.

Teachers used these platforms to provide live classes, asynchronous video lectures, and virtual coding labs. The ways of assessing students changed and included automated coding evaluation systems, peer reviews, and project-based submissions. The flexibility and scalability of online programming education made It has been an alternative during the pandemic and remains highly visible in the post-COVID era.

1.2.1 Importance of Online Programming Education

Indeed, online programming education has come to be standard as being highly effective in democratizing access to coding information and skills. It is therefore the must have of innovativeness because online educational applications could offer solutions that are easily accessed through scale as an imperative. For the beginners, the online channels provide learning pathways and an interactive toolkit to fill this gap between the theoretical knowledge and practical use. Automated systems of feedback allow learners insight into code correctness as well as performance while providing them with the opportunity to self-correct their skills without reliance on instructors.

Online programming education, from the teacher's point of view, has analytics and performance metrics by which one can track progress in learners. This would provide data-driven and informed teaching strategies, besides enabling

personalized learning through recommendation of tasks and resources that suit a learner's skill level and history of performance, ensuring students get content tailored to their unique needs. Such scalable platforms enable a larger cohort of students to access the educational experience offered without diminishing the quality bar, especially where the supply of resources for the traditional systems is less easily accessible.

Moreover, online programming education creates an environment of collaborative learning through forums, coding challenges, and peer review, which promotes knowledge sharing and teamwork. It also meets industry demand because it focuses on how things are applied in reality as well as giving students projects and challenges that allow them to build portfolios. In discussing this further, the work from this research shows how adaptive the online programming education will become with intelligent task recommendations as well as advanced assessment systems remain deeply embedded in its imperative nature to shape future programmers.

1.2.2 Limitations of Online Programming Education

Although online programming education is revolutionary, there are several limitations at the heart of this research, especially regarding personalized learning paths and the automated assessment systems. The primary reason is that the automated assessment systems cannot accurately calculate the subtleties involved in a learner's programming. Current evaluation tools primarily revolve around binary correctness metrics—whether code passes a set of predefined test cases. They have largely neglected learning, creativity, and the intermediate errors that indicate the quality of understanding and problem-solving in a learner's approach. Such is the case for novice programmers whose advancement relies on constructive feedback and iterative learning.

On top of that, it has reduced face time between tutors and learners thus minimizing the real-time communication on doubts and specific learning guidance. Besides this, another very important softness is relying on datasets and fixed exercise content in which static contents cannot possibly change with

learners' new status and changed skills of time. Current online programming education bases lack the mechanisms to categorize and report the unique issues in each learner, like often syntax errors or conceptual misconceptions. Moreover, problems of academic integrity such as plagiarism or unauthorized collaboration question the reliability of assessments. These limitations underscore the need to develop cultured, adaptive, and context-aware systems, like those investigated in this paper, to move the personalization and efficacy of online programming education much further forward.

1.2.3 Future Direction for Online Programming Education

The future of online programming education relies on the resolution of existing limitations through the utilization of advanced, adaptive technologies as well as pedagogical refinement in support of novices. At the core of this vision is the establishment of intelligent, context-aware systems that adapt to personal and unique needs, errors, and skill levels of the individual learner. This is possible when such models use machine learning capabilities to suggest programming assignments and provide specific feedback regarding aspects where the student tends to struggle. Automated assessments should be advanced beyond a 'correctness' response—using measures such as 'efficiency of code, creativity of code, error analysis, etc. '; this can help determine greater learner understanding while upholding assessment validity.

Gamification and collaborative coding environments can be used to improve engagement in which learners are active participants and collaborators in virtual communities. Immersive technologies, including virtual and augmented reality, can feign real-world coding conditions that students can engage in by doing hands-on exercise in virtual labs that simulate proficient development environments. Data analytics and predictive modeling can also help instructors identify students who are at risk, thereby being intervened and supported at an appropriate time.

Future systems should also highlight ethical dimensions of programming, such as proper coding practices and awareness of problems in cybersecurity.

Combining innovative technology, adaptive procedures, and a learner-centric method will enable a new generation of online programming education to bridge the gaps that exist, offering more equitable learning opportunities and getting learners ready for the very dynamic demands of the world of programming. This research is in line with these directions, aiming to contribute to the design of adaptive systems that empower novice programmers and enhance the overall effectiveness of online education.

1.3 Overview of Machine Learning

Machine Learning is a sub-type of AI that enables computers to learn and improve themselves by gaining experience without being explicitly programmed. Its application lies in algorithms in detecting patterns, making forecasts and inferences from given data; hence, used to generate recommendations, in predictive analytics, and other decisions-making processes. As presented below is a summarized explanation of the general categorizations of the main machine learning algorithm categories with their functionalities.

1.3.1 Supervised Learning Algorithms

Supervised learning algorithms are one of the basic machine learning techniques that depend on labeled datasets to learn the input variables, also known as features, and output variables, also known as labels. These algorithms work by learning a model from historical data where the desired output is clearly given, so the model can predict outcomes for unknown data. The main thought is to minimize the error between the predicted and the actual outputs by iteratively adjusting the parameters of the model.

Among all the simplest supervised learning algorithms, there's Linear Regression, which helps predict output values to be continuous with numerical types while assuming some linear relationship exists between your input features and the target. For example, it's possible to determine house price estimates by such features such as area, location, and the number of rooms. Meanwhile, Logistic Regression is applied when the decision in a problem will be true or false. So, as an example, it was distinguished between spam and non-

spam emails. Despite its name, logistic regression uses the logistic function to model probabilities, making it a solid tool for classification tasks.

Decision Trees are another of the popular supervised learning algorithms which break down the data into segments based on the value of the features. The node in a decision tree act as a decision rule and leaves are output classes or values. Their interpretability, besides the ability to pick up nonlinear relationships, is their main reason for such popularity in applications such as customer segmentation and medical diagnosis. They are prone to overfitting, though especially if the tree gets too deep. To address this issue, ensemble methods such as random forests combine multiple decision trees while training each on a different random subset of data then aggregating their predictions. Thus, it improves the robustness of a model, reduces overfitting, and handles well high-dimensional data.

SVM is effective for classification tasks. It finds a hyperplane in multi-dimensional space that separates the data into distinct classes best. SVM can be used with kernel functions, thus making it suitable for applications like image recognition and bioinformatics, where the boundaries are complex and non-linear. Another proximity-based method is KNN. This algorithm classifies a data point based on the majority class of its nearest neighbors. KNN is simple and intuitive in nature but is computationally expensive for large datasets owing to its dependency on distance calculation.

Ensemble methods like Gradient Boosting (for example, XGBoost) form an advanced category of supervised learning. Such algorithms construct an ensemble of weak learners, typically in the form of decision trees, by sequentially correcting the errors of the previous models. This technique of gradient boosting is highly accurate and is used on most competitive data science projects for the reason that it can deal with missing values and even complex interactions between features.

To Identify Learning Path for Novice Programmer Based on Semi-supervised Dataset
Using Proposed Machine Learning Algorithm

Algorithm	Key Features	Applications
Linear Regression	Predicts continuous numerical outputs, assumes a linear relationship between variables.	Predicting sales, stock prices, house pricing.
Logistic Regression	Binary classification, outputs probabilities for classes.	Spam detection, medical diagnosis.
Decision Trees	Tree-like structure, interpretable, handles categorical and numerical data.	Customer segmentation, fraud detection.
Random Forest	Ensemble of decision trees, reduces overfitting, handles high-dimensional data.	Credit scoring, feature importance analysis.
Support Vector Machines (SVM)	Finds optimal hyperplanes, effective for both linear and non-linear classification.	Image classification, bioinformatics.
K-Nearest Neighbors (KNN)	Proximity-based, simple yet effective for small datasets.	Handwritten digit recognition, recommendation systems.
Gradient Boosting (e.g., XGBoost)	Builds sequential models to correct previous errors, highly accurate.	Predictive analytics, customer churn prediction.

Table 1.1: Supervised Machine Learning Algorithms

Supervised learning algorithms are extremely useful in delivering an accurate prediction if sufficient amounts of labeled data are provided. However, their ability heavily depends on the quality of the data, how feature selection is carried out, and the complexity of the model selected. They find numerous applications in finance (credit scoring), healthcare (prediction of diseases), and educational systems (student performance analysis), thereby rendering them a necessity in the field of machine learning.

1.3.2 Unsupervised Learning Algorithms

Unsupervised learning algorithms are designed to analyze and infer patterns from datasets that lack labeled outcomes or target variables. Unsupervised learning is different from supervised learning, which involves explicit guidance. Instead, unsupervised learning operates in an independent manner to discover latent structures, groupings, or distributions within the data. Such algorithms are very helpful for exploratory data analysis, dimensionality reduction, and clustering, providing insights that may not be immediately apparent.

Algorithm	Key Features	Applications
K-Means Clustering	Groups data into clusters based on similarity, uses centroids.	Market segmentation, document clustering.
Hierarchical Clustering	Builds a hierarchy of clusters, merges or splits recursively.	Gene sequence analysis, social network analysis.
Principal Component Analysis (PCA)	Reduces dimensionality, retains variance, simplifies datasets.	Image compression, data visualization.
Anomaly Detection	Identifies unusual data points or outliers.	Fraud detection, network intrusion detection.

Table 1.2: Unsupervised Machine Learning Algorithms

K-Means Clustering is an often-used, unsupervised learning unguided algorithm that partitions observations in the space of similarities according to a predefined number of clusters. Data points are iteratively assigned to clusters and updated centroids computed. This approach is vastly found in the area of marketing segmentation in grouping of client behaviors to formulate customized marketing policies for the respective groups. Hierarchical Clustering forms a hierarchy of clusters through the merging or splitting of them in iterative steps, based on similarity measures. It has proved useful in the fields

of bioinformatics for the study of hierarchical relationships between gene sequences or protein structures.

For multi-dimensional datasets, Principal Component Analysis is a technique that decreases the number of variables while keeping most of the variance of the dataset. PCA transforms data into a new coordinate system whereby the dimensions are ranked based on the amount of variance that they explain. It has been widely used in image processing and visualization to simplify the complex datasets without losing critical information.

The other important application of unsupervised learning is Anomaly Detection. It identifies data points which significantly stray from the normal and often represent deceitful transactions, equipment catastrophes, or other unusual activities. Clustering and statistical methods are often pragmatic in real-time monitoring systems to flag incongruities.

Unsupervised learning can also be used in more complex models. This includes Autoencoders, which are a special kind of neural network: one that compresses input data into a latent space and then reconstructs the original data. The discrepancies between the original data and the reconstructed data may indicate any pattern or anomaly. Autoencoders are very efficient in compression for images, denoising, and feature extraction.

Though very versatile, unsupervised learning algorithms have the problem of interpreting the results. There is no predefined labeling available to validate discovered patterns. The choice of similarity measures, the nature of the dataset, and appropriate tuning of hyperparameters often determine whether the approach will be a success in using such algorithms.

Unsupervised learning is needed to uncover latent patterns and trends in unlabeled datasets. It is used in diversified industries, including customer categorization in marketing and threat identification in cybersecurity, along with the preprocessing of data in most pipelines of machine learning. The fact

that the algorithms can reveal insights irrespective of explicit labels makes them really very important for data exploratory purposes and knowledge uncoverage.

1.3.3 Semi-Supervised Learning Algorithms

Semi-supervised learning algorithms fill the gap between supervised and unsupervised learning. These algorithms use a mix of labeled and unlabeled data. Such an algorithm is extremely useful when it is costly to obtain labeled data, or gathering such data is time-consuming or even impossible, but there is a huge volume of unlabeled data. The structure and patterns in the unlabeled data combined with the examples labeled improve model performance and reduce dependency on extensive labeling.

Algorithm	Key Features	Applications
Self-Training	Uses labeled data to label and include unlabeled data iteratively.	Text classification, image recognition.
Graph-Based Learning	Leverages graph structures to propagate labels.	Community detection, recommendation systems.
Generative Models (e.g., GANs, VAEs)	Generates new data points, models data distribution.	Medical imaging, data augmentation.

Table 1.3: Semi-supervised Machine Learning Algorithms

The most intuitive semi-supervised learning approach is Self-Training. Here, it starts with the training of a supervised model on the labeled data and then iteratively adds its predictions on the unlabeled data as pseudo-labels. This process is repeated until the model stabilizes. This can be very effective for many tasks, such as text classification and image recognition. However, self-training can propagate errors if the model's initial predictions on unlabeled data are wrong.

The final set is Graph-Based Semi-Supervised Learning that offers representing data in graph structures using data points as the representation and edges to resemble connectivity relationships. Algorithms that do propagate labels from those marked by the user as some key nodes and allow dissemination with propagators through connectivity relationships that function more spectacularly in finding out different social network communities, etc, and recommendation algorithms.

Semi-supervised learning in high-dimensional data can be facilitated through Generative Models such as VAEs and GANs. These models produce new data points but also use the underlying distribution of the data to derive an understanding of the trends present in the data. For instance, for medical imaging, GANs could be used to synthesize realistic images to help improve the accuracy of classifiers through training on augmented datasets.

Semi-supervised learning is very helpful for NLP tasks where there are often fewer labeled data available, such as annotated text. BERT and GPT-like pre-trained models fine-tune on a few thousand labeled samples while still relying on the enormous amount of unannotated text for optimal performance in applications such as sentiment analysis and question answering.

The primary benefit of semi-supervised learning is to minimize the dependence on labeled data with high accuracy. The technique is applied extensively in such domains as medical diagnosis where gathering labeled data needs the intervention of experts, or fraud detection, in which anomalies are rare but important. It determines the success of semi-supervised learning by exploiting assumptions regarding the quality of labeled data and the assumption that the distributions or patterns followed by the labeled and unlabeled data are similar. Heavily degrading the performance will be caused in the case of violation of such assumptions.

High-efficiency algorithms on real-world problems with minimal labeled data apply semi-supervised learning algorithms. These combine the two supervising

and unsupervised approaches to learning tasks. Applications have been shown in many fields, for example, health, cybersecurity, and autonomous systems due to the possibility of fusing the labeled and unlabeled data with the optimization of learning. The bottleneck remains with labeled data yet there is sufficient unlabeled data to tap, hence this will still remain a transformative role in machine learning applications.

1.3.4 Reinforcement Learning Algorithms

Reinforcement learning is a very strong paradigm in machine learning. In this, the agent learns how to make decisions based on interactions with an environment rather than explicit instructions by the environment. It is different from supervised learning because RL focuses on sequential decision-making: here, the agent optimizes a strategy or policy so that cumulative rewards over time are maximized. Q-Learning and Deep Q-Learning allow the key algorithms applied to estimate the value of actions, and Policy Gradient and Actor-Critic methods enable direct optimization of policies in continuous and dynamic environments.

Algorithm	Key Features	Applications
Q-Learning	Learns action values for states, maximizes cumulative rewards.	Game-playing agents, robotic navigation.
Deep Q-Learning (DQN)	Combines neural networks with Q-Learning to handle large state spaces.	Video game agents, autonomous systems.
Policy Gradient Methods	Directly optimizes policies for continuous or discrete actions.	Robotics control, financial trading.
Actor-Critic Methods	Combines policy optimization (actor) and value estimation (critic) for stability.	Autonomous driving, energy optimization.
Proximal Policy Optimization (PPO)	Balances exploration and exploitation, ensures stable updates.	Virtual simulations, industrial automation.

Table 1.4: Reinforcement Learning Algorithms

The use of proximal policy optimization with Monte Carlo algorithms improves higher learning efficiency in complex scenarios. It has been implemented so far in robotics, computer gaming, autonomous navigation systems, and healthcare-related environments for demonstrating the feasibility and validity of RL in unimagined dynamic and uncertain domain environments. However, efficiency has rarely been fully devoid of offset balance by computational cost, exploration-exploitation trade-off, and design in reward-creation mechanisms. Despite these challenges, the RL ability to handle complex tasks in decision-making places it at the heart of the innovation in artificial intelligence.

1.3.5 Neural Networks and Deep Learning

Neural networks and deep learning are highly advanced, sophisticated machine learning techniques inspired from the structure and functioning of the human brain. It includes layers with interconnected nodes that process their input data through weighted connections; this is how they tend to learn complex patterns as well as relationships. Deep learning is a type of neural network characterized by architecture with many hidden layers intended to extract hierarchical features representing intricate, nonlinear relations in large data sets.

Algorithm	Key Features	Applications
Feedforward Neural Networks	Processes input through interconnected layers, handles complex patterns.	Stock price prediction, classification tasks.
Convolutional Neural Networks (CNNs)	Extracts spatial features, ideal for image and video processing.	Object detection, medical imaging.
Recurrent Neural Networks (RNNs)	Processes sequential data, remembers temporal dependencies.	Time-series forecasting, natural language processing.
Transformer Models	Efficient for sequential data, utilizes attention mechanisms.	Language translation, chatbots, sentiment analysis.

Table 1.5: Neural Networks and Deep Learning Algorithms

Other variants, like CNNs, really come into their own in image and video processing by identifying spatial patterns, while RNNs and Transformers handle sequential data, so they are extremely useful in natural language processing and time series analysis. These models run applications such as object detection, speech recognition, language translation, and autonomous driving. Even though these models and neural networks are superior, they are very computational-intensive and require large datasets, so they are criticized for being less interpretable compared to traditional machine learning algorithms. However, the power to address high-dimensional and complex problems has made them tools in modern AI that can propel advancement in all spheres, including healthcare, finance, and technology.

1.3.6 Ensemble Learning

Ensemble learning is one of the most powerful techniques in machine learning that pools the predictions of multiple models together to improve overall performance, accuracy, and robustness. Ensemble methods are particularly successful for complex and noisy datasets by leveraging the strengths of individual models and mitigating their weaknesses. Key approaches include bagging, where models are fitted on different subsets of the data to reduce variance (like in Random Forest), and boosting, which sequentially improves a sequence of weak models by emphasizing what the previous one does wrongly (like in Gradient Boosting and XGBoost). Another method is called stacking, where predictions from multiple models are integrated using a meta-model in order to produce a final output, enhancing generalization.

Applications in finance, healthcare, and recommendation systems are prime domains of ensemble learning where high accuracy is of extreme importance. While these techniques usually perform better than single models, they are computationally expensive and require careful tuning of hyperparameters. Despite these challenges, ensemble learning remains the core foundation of predictive analytics for providing robust solutions in academia as well as in practice.

Algorithm	Key Features	Applications
Bagging (e.g., Random Forest)	Combines predictions from multiple models to reduce variance.	Feature selection, customer analytics.
Boosting (e.g., XGBoost, AdaBoost)	Focuses on correcting errors of weak models, improves generalization.	Loan approval prediction, marketing response models.
Stacking	Combines outputs of different models using a meta-model for final prediction.	Hybrid recommendation systems, predictive analytics.

Table 1.6: Ensemble Learning Algorithms

1.4 Research Motivation

As a programming instructor, the reasons for this research are mainly because novice learners always fail to overcome the persistent challenge of acquiring programming skills in their learning journey. Teaching the programming skill often reveals an uphill task for beginners, as one of the common hurdles is difficulty in understanding syntax and logical problem-solving and then debugging errors. It presents a challenge that students can become discouraged and drop out at alarming rates, lacking the self-assurance to continue studying within the field. Existing structures of education, traditional as well as online, generally do not allow for such individualized support and coaching, which is needed to manage all these different learning requirements in a proper manner.

The COVID-19 pandemic has accelerated the shift to online programming education, further exposing the shortcomings of the traditional teaching methods. Currently, automated assessment tools only test for binary correctness and nothing about the learning process and intermediate steps that reflect a student's approach to solving problems. Due to the lack of personalized feedback and dynamic task allocation, several students end up in unguided struggle without knowing where to improve.

As a teacher, this research is motivated by the desire to empower students with a supportive and adaptive learning environment. It aims to create a system that not only evaluates performance but also suggests appropriate learning paths tailored to individual strengths, weaknesses, and progress. This research is to develop a hybrid model using semi-supervised datasets with the most advanced machine learning algorithms, which can connect the bridge between automated systems and personalized teaching. Such a model revolutionizes the teaching of programming, not only making it more inclusive and effective but also in consonance with the needs of novice learners.

This research is also motivated by the need to reduce the teacher's cognitive load with an intention to focus more on strategic interventions instead of routine evaluations. An intelligent task recommendation system can make the teaching process efficient and result in learning outcomes that will never allow a student to lag behind his or her own programming method. This work will help to fill in all of these critical gaps in the teaching of programming, arm educators, and support students in finding a solid foundation for their futures in technology.

1.5 Problem Statement

Many problems arise through training in programming schooling; concerning the diversity of needs when coding. Many of these students end up frustrated due to their failure in finding other more unique ways to study, leading to the disengagement of classes and eventually dropping off in school. The processes by which things are done from olden classes and new-comer computer-based tools result quite similar in all circumstances in which people find them working. This oneness does not consider how distinct the students of learning and their skills may differ, therefore bringing about some poor performances and outcomes from the resultant learning.

Previous study has mostly looked at ways to find mistakes and give feedback, but it hasn't looked at how learners' needs change and adapt. Machine learning approaches have made known promise in educational data gathering and predicting

success, but they have not yet been used to create personalized learning routes. Also, prediction analytics and task scheduling aren't used together very often to change learning activities based on how sound students are doing right now.

These holes show how important it is to have a strong system that can't only guess what skills students have but also give them flexible, individual learning paths. This kind of system should look at the computer tasks that are turned in, put students into clusters based on how well they did, and then suggest next tasks that are just correct for each student. Taking care of this problem is necessary to boost participation, student achievement, and the general success of computer education.

1.6 Research Objective

Dataset Preparation: The efficacy of any machine learning model in programming instruction is profoundly affected by the quality and relevance of the dataset. This mission aims to develop a comprehensive dataset designed to detect the learning trajectories of rookie programmers. The dataset is derived from authentic code submissions either from online coding platforms or academic coding settings. Attributes like code accuracy, error patterns, execution duration, and test case efficacy are retrieved and preprocessed. Furthermore, derived factors such as job difficulty and submission behavior (e.g., number of tries) are included to guarantee that the dataset reflects the intricate nature of programming skill growth. This study incorporates an analogy of obesity-like features to improve the model, as seen in predictive healthcare studies, ensuring that the dataset is both clean and well-structured, as well as abundant in predictive qualities pertinent to learners' success.

Propose a machine learning algorithm: This aim explores and compares the many algorithms. It lists KNN, decision tree, random forest, and finally XG Boost as major algorithms, as discussed above in order to get best efficient ones for dataset. The investigation of this research provides ensemble approaches, like Adaptive Ensemble Classification Weighting, or known as (AECW). This proposed algorithm gathers together strength from diverse models. This approach helps solve problems such as uneven class distributions, different levels of job complexity, and scalability. The goal is to extant techniques that allow for better task adaptation with

the help of proper classification of learners, hence evaluation in terms of accuracy, precision, recall, and F1-score.

Exploiting the Learning Trajectory of a Novice Programmer: A definition of individualized learning pathways for novice programmers is the foremost objective of this project. The concept of code submittals and performance measurement is applied to a machine-learning-based framework that classifies learners as slow, middle, or fast learners based on an analysis of performance measures such as code quality, mistake repair, job completion rates, and dynamics that suggest further tasks to be offered to the learner. This adaptive method promotes constant development of skills with the aid of engaging students in appropriately challenging activities. The study seeks to explore how predictive modeling might improve instructor capabilities, particularly by offering actionable insights about student development that would permit interventions and refine course designs. This purpose ultimately serves to align the demands of learners with the capabilities of current programming education platforms.

1.7 Research Contribution

This study addresses critical concerns related to programming education through introducing innovative hybrid learning tactics that combine the applications of machine learning to personalize experience tailored for beginner programmers.

Creation of an Adaptive Learning Framework: It is a machine learning architecture that dynamically adjusts according to the needs of every individual learner by analyzing the submitted work for the tasks. This paradigm classifies learners-for example slow, medium, fast and also suggests suitable coding projects tailored to the skill level of the learner, thereby encouraging more participation and higher proficiency development.

Incorporation of Sophisticated Machine Learning Algorithms: The incorporation improves learner classification accuracy as well as its robustness along with task adaptability employing techniques, among them including, are Decision Tree, K-Nearest Neighbors, XGBoost, Random Forest and Adaptive

Ensemble Classification Weighting; in such models, adequate accuracy performance predictions would surely be achieved by addressing complications arising from task complexity or even class imbalances.

Development of an Extensive Dataset: A structured dataset is created with attributes such as accuracy, error patterns, task difficulty, and submission behaviors. This dataset is used to train and test the prediction models so that the framework captures the complexities of programming skill development appropriately.

Dynamic Task Arrangement: The proposed model describes data-driven task sequencing sensitive to real-time changes in learners' learning trajectories. Thus, the method ensures that learners get tasks of appropriate complexity to minimize frustration and enhance learning efficiency.

Assistance for Educators: It does this by providing actionable insights into learner performance to help educators make more informed choices around curriculum design and pedagogical tactics. With automated feedback and performance measurement, the approach reduces teacher workload and makes instruction scalable.

Influence on the Problem Domain: This study links the traditional approach to teaching and the need for individual learning in programming education. Employing predictive analytics and adaptive learning techniques, it makes the programming education more accessible, efficient, and fun to learn for beginner programmers.

1.8 Scope of Study

This study aims to improve the learning experience of novice programmers by designing a strong machine learning-driven framework for learner classification and personalized task recommendation based on their performance. The proposed system is based on the Adaptive Ensemble Classifier with Complexity-Aware Weighting (AECW), which offers a data-driven approach to enhancing

programming education. The scope of this study is therefore defined by the following key areas:

Personalized Learning in Programming Education: Adaptive learning systems are emerging for diverse needs. By creating a set of groups between fast and slow, medium learners get categorized to generate task recommendations according to proficiency and learner progress.

Application of Machine Learning Algorithms: The research assesses the effectiveness of several machine learning algorithms that are implemented to predict the performance of learners and adapt learning paths dynamically by using K-Nearest Neighbors, Decision Tree, Random Forest, XGBoost, and AECW algorithms.

Feature-Rich Dataset Utilization: A semi-supervised dataset with task complexity, trial counts, and error metrics is utilized in this study. This dataset is then used as a foundation for model training and validation purposes for the task recommendation system.

Dynamic Task Recommendation System: one of the main emphases is for the development of a systematic recommendation of tasks of varying levels and complexities based on learners' past performance. The recommendations must be dynamic to avoid overload/underload of the learners on tasks, thereby consistently being engaged and developing skills.

Improve learner classification: The study aims at classifying learners by analyzing submission patterns and performance metrics, hence leading to the accurate identification of learner types and unique needs through personalized learning pathways.

Scalability and Adaptability: The research explores the scalability of the AECW algorithm and its adaptability across diverse educational contexts, ranging from institutional programming courses to online coding platforms. The findings are

expected to benefit both educators and learners by providing a scalable, automated solution to improve learning outcomes.

Future Applications in Education: The proposed methodology can be extended to other domains requiring personalized learning systems, such as mathematics, science, and language education. While the study focuses on programming education, the research lays the groundwork for future interdisciplinary applications of machine learning in education.

This study's scope is limited by its focus on programming education and the dataset that it used for model training and validation. Although this research illustrates the potential of a personalized task endorsement system, more work and implementation in actual educational settings are necessary to establish the full scope of applicability.

1.9 Organization of Thesis

This thesis is presented in six chapters, each of which would be structured to give a comprehensive overview of the research work undertaken, ranging from foundational concepts to experimental results and future directions.

Chapter 1: Introduction

This chapter introduces the research topic and lays down the ground for the study. It gives an overview of the research motivation, problem statement, objectives, and contribution of the work. Additionally, it includes a discussion on the scope of the study and how it is structured.

Chapter 2: Literature Review

This chapter gives a detailed review of the existing literature that is relevant to programming education, machine learning techniques, and adaptive learning systems. It identifies gaps in research and positions the proposed work within the context of current advances in the field.

Chapter 3: Dataset Preparation

This chapter will describe the sources of the dataset, preprocessing techniques, feature engineering, and strategies used in partitioning data in order to ensure that it is adequate for analysis and modeling.

Chapter 4: Existing Methodology

This chapter deals with the methodologies currently in use in programming education and predictive analytics. It provides a comparative analysis of existing approaches, algorithms, and highlights their shortcomings, which the proposed methodology attempts to overcome.

Chapter 5: Development of Proposed Machine Learning Algorithm

This chapter presents the Adaptive Ensemble Classifier with Complexity-Aware Weighting (AECW) algorithm. It describes the design, implementation, and how it will be integrated into the task recommendation framework to achieve personalized learning.

Chapter 6: Results and Discussion

This chapter presents the experimental result obtained after evaluating the proposed algorithm and comparing the performance with some baseline models. Important performance metrics, which include accuracy, precision, recall, and F1-score, are discussed along with insight derived from the analysis. It summarizes the main outcomes and contributions of the research for the conclusion of the thesis and points out the limitations within the current work and directions for further research-including possible improvements and wide applications of the proposed framework.

This structured arrangement guarantees a logical flow of ideas and provides readers with a strong understanding of the research expedition, from identifying the problem to proposing solutions and examining their impact.