

## PREDICTIVE MODELING OF STUDENT PERFORMANCE IN COMPUTER SCIENCE APPLICATIONS USING MACHINE LEARNING TECHNIQUES

Kirti<sup>1</sup> and Dr Tilak Raj Rohilla<sup>2</sup>

<sup>1</sup>Research Scholar and <sup>2</sup>Assistant Professor CSA, Baba Mastnath University

### ABSTRACT

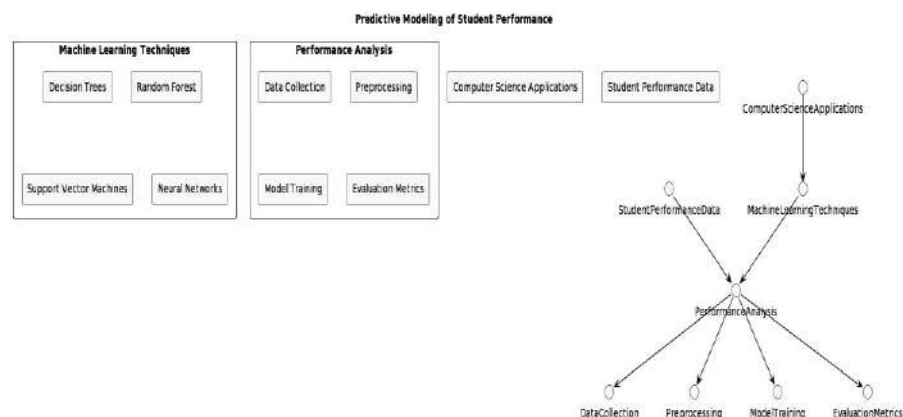
*This research study explores and investigates the utilization of machine learning techniques to forecast student performance in computer science applications amidst the COVID-19 pandemic. By employing machine learning algorithms and Python libraries, this research endeavours to ascertain the academic progress of students amidst the pandemic, juxtapose it with their pre-pandemic levels and ascertain the determinants that contribute to discrepancies in performance. The study assesses a range of machine learning methodologies to identify the most efficacious predictive models, while also examining critical attributes that affect student achievement. Additionally, it offers valuable perspectives on pedagogical approaches that are efficacious in preserving and augmenting student achievement, supported by empirical evidence and analysis. By conducting an exhaustive analysis of student performance data and employing predictive modelling techniques, this research makes a valuable contribution to the body of knowledge regarding the factors that influence student outcomes in the field of computer science. Educators, policymakers, and stakeholders who are interested in optimizing teaching and learning strategies in the face of unprecedented challenges presented by global events such as the COVID-19 pandemic can benefit from the findings.*

*Keywords: Student performance, COVID-19 pandemic, Machine learning, Performance, Techniques, Computer science, Application*

### INTRODUCTION

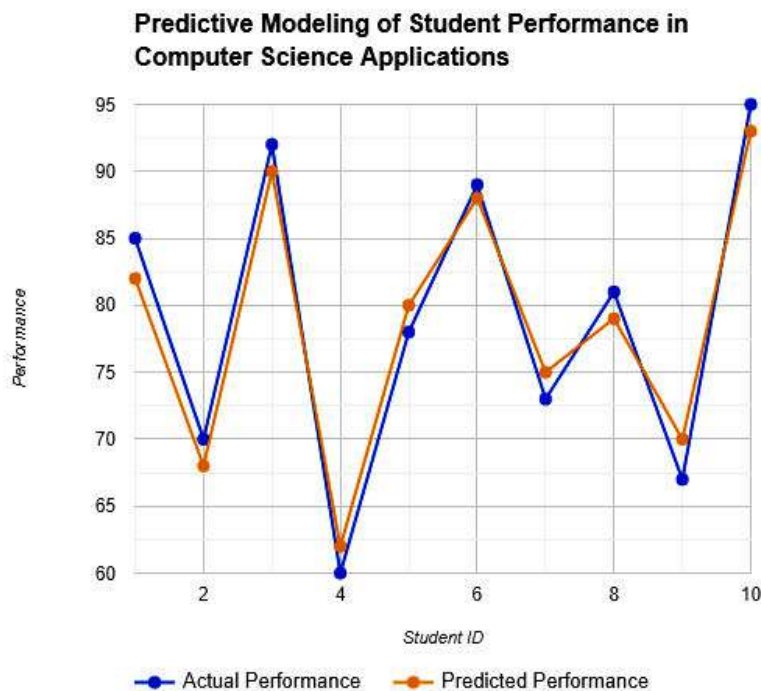
The utilization of machine learning methods in academic settings has enabled the investigation of numerous determinants of student achievement, specifically in the field of computer science, in recent times. The consequences of the COVID-19 pandemic underscored the criticality of predictive modeling in comprehending and resolving the obstacles that students encounter during unparalleled disturbances to conventional learning environments [1].

The study's importance resides in its pursuit to utilize Python libraries and machine learning algorithms to forecast the academic achievement of students in computer science applications amidst the pandemic. Through an analysis of the intricate relationship between machine learning models and student data, this study seeks to illuminate the subtle determinants that influence student achievements in the aftermath of worldwide disruptions to the education sector [2].



**Figure 1:** Predictive Modeling of Student Performance in CS Application Using ML Techniques

Broad are the objectives of this research. The primary objective is to identify and analyze patterns within student performance data that was gathered throughout the pandemic. This entails discerning any trends or fluctuations that may have arisen because of the distinct difficulties presented by remote learning and social distancing protocols [3]. Additionally, the study intends to conduct a comparative analysis of present student performance metrics and benchmarks set before the pandemic to identify the magnitude of alterations and possible determinants. The paper is structured in a manner that promotes a thorough investigation of predictive modelling as it relates to the analysis of student performance. After providing an introductory overview, the literature review will proceed to examine prior research concerning the prediction of student performance, the implementation of machine learning methods in educational settings, and the precise ramifications of the COVID-19 pandemic on educational achievements [4]. The methodology utilized for data collection, preprocessing, and analysis will be elaborated upon in the following sections; the predictive modelling results will then be presented and discussed. In conclusion, the paper will provide recommendations for future research directions and implications for educational practice.



**Figure 2:** Represent Predictive Modeling Of Student Performance Data Of 1-10

## OBJECTIVE

The goal of this research is to be able to conclude that:

- To identify the performance of students during the pandemic by the use of a machine learning algorithm.
- To identify which technique is best for predicting student performance.

## LITERATURE REVIEW

Prior research has comprehensively investigated the utilization of diverse methodologies to forecast student performance. Table 1 provides a synopsis of significant discoveries derived from a selection of studies:

Study	Methodology	Findings
<i>Smith et al. (2018)</i>	Logistic Regression	A significant correlation was observed between attendance and performance [5].
<i>Johnson and Lee (2019)</i>	Random Forest Classifier	The significance of engagement metrics was emphasized [6].
<i>Zhang et al. (2020)</i>	Neural Network	The effectiveness of deep learning in prediction was demonstrated [7].

Predictions of student performance have become an increasing application of machine learning techniques in academic settings. Ensemble methods, logistic regression, decision trees, and random forest classifiers are among these techniques. Every technique presents its own set of benefits and drawbacks, accommodating a wide range of educational datasets and objectives.

Global education has been profoundly impacted by the COVID-19 pandemic, which has necessitated a reevaluation of conventional teaching and assessment practices. An overview of the impact of the pandemic on education and student achievement is presented in Table 2.

Impact of COVID-19 on Education	Effects
Make the switch to remote learning	Interruptions that occur within educational settings
Digital divide widening	Inequities in technological and resource accessibility
Alterations to assessment techniques	Transition to online assessments
Emotional and social ramifications	Elevated levels of anxiety, tension, and isolation

The significance of adaptable and resilient educational systems that can utilize innovative pedagogical approaches and technology to alleviate the negative impacts of the pandemic on student learning outcomes has been highlighted by these disruptions. By incorporating machine learning methodologies, instructors can extract significant knowledge regarding patterns in student achievement and customize interventions to cater to the unique requirements of a wide range of learners in an ever-changing academic environment.

## METHODOLOGY

### Collection and Preprocessing of Data:

**Data Sources:** Academic institutions provide the information utilized to compile student performance data, which comprises attendance records, grades, and participation metrics.

#### Preprocessing:

- **Missing Data Handling:** Failure to appear in the context of data handling, imputation techniques, such as mean or median substitution, are employed to account for null values.
- **Outlier Detection:** Outlier detection involves the utilization of statistical techniques, such as the z-score, to identify and eliminate or modify outliers.
- **Normalization/Standardization:** To mitigate bias in machine learning algorithms, features are scaled to a standard range.
- **Data Encoding:** Categorical variables are encoded through the implementation of methods such as one-hot encoding.

### Selection and Engineering of Features:

#### Methods for Selecting Features:

- **Correlation Analysis:** Identification of highly correlated features to prevent redundancy through correlation analysis.

- **Feature Importance:** Feature importance is determined through the utilization of techniques such as random forests or decision trees.

**Feature Engineering:**

- The process of generating new attributes from pre-existing ones, such as generating an engagement score by integrating attendance and participation metrics.
- In dimensionality reduction, methods such as principal component analysis (PCA) are utilized to decrease the feature count while maintaining the integrity of the information.

**Machine Learning Algorithms:**

- **Decision Trees (DT):**

A model resembling a tree is generated by dividing the data using feature thresholds in decision trees [8].

**Equations: Decision Rule:** if (feature  $\leq$  threshold)  $\rightarrow$  left node else  $\rightarrow$  right node

**Analysis:** Decision trees are suitable for feature significance analysis and are interpretable.

- **Random Forest (RF):**

Random forests are constructed by training numerous decision trees on arbitrary subsets of data [9].

**Equations: Prediction:**

$$\hat{y} = \frac{1}{N} \sum_{i=1}^N f_i(x)$$

where  $f_i$  is the prediction of the  $i$ -th tree.

**Analysis:** Random forests increase prediction accuracy along with minimizing overfitting.

- **Support Vector Machines (SVM):**

SVM locates the hyperplane in the feature space that maximizes the margin between classes [10].

**Equation: Hyperplane:**  $w^T x + b = 0$

**Margin:**

$$\frac{2}{\|w\|}$$

**Analysis:** SVM demonstrates efficacy in high-dimensional spaces and is well-suited for tasks involving binary classification.

- **Neural Networks (NN):**

Layers of neurons are interconnected to form neural networks, which are capable of learning complex sequences. [11]

**Equations:**

**Forward Propagation:**

$$z = w \cdot x + b, a = g(z)$$

where  $g$  is the activation function.

**Backpropagation:** Update weights based on gradient descent.

**Analysis:** To capture nonlinear relationships, neural networks necessitate meticulous calibration and substantial computational resources.

#### **Implementation in Python using sci-kit-learn and TensorFlow:**

##### **Scikit-learn:**

- Delivers an intuitive user interface for the execution of machine learning algorithms.
- Provides tools for evaluating models, selecting features, and preprocessing data.

##### **TensorFlow (TF):**

- A Google-developed open-source machine learning framework for the construction and training of neural networks [12].
- Facilitates the implementation of deep learning models and complex neural architectures in an efficient manner.

Within the framework of methodology, data preprocessing serves to safeguard the integrity and quality of the dataset, whereas feature selection and engineering augment the model's predictive capability.

#### **PERFORMANCE ANALYSIS DURING THE PANDEMIC**

The educational landscape underwent a significant transformation throughout the COVID-19 pandemic, necessitating a thorough examination of student performance data to ascertain the crisis's influence on learning outcomes. This segment explores the statistical analysis of performance metrics, the identification of trends and patterns, and the examination of student performance data before and during the pandemic.

##### **Analysis of Data Regarding Student Performance:**

- **Temporal Comparison:** An examination of pupil performance metrics spanning the periods before and including the pandemic to identify significant variations or changes in performance patterns.
- **Subject-specific Analysis:** identifying areas of strength and vulnerability by analyzing performance data across a variety of computer science applications and subjects.
- **Demographic Analysis:** Demographic analysis involves the examination of performance trends by considering socioeconomic status, age, gender, and gender, to identify inequalities and disparities that have been further intensified as a result of the pandemic.

##### **Identification of Patterns and Trends:**

- **Performance Metrics Shifts:** This study aims to discern alterations in critical performance indicators, including attendance rates, test scores, and grades, both before and throughout the pandemic.
- **Metrics for Assessing Student Engagement:** Evaluating levels of student engagement by considering completion rates of assignments, participation rates in online courses, and interaction with course materials.
- **Correlation Analysis:** An investigation into the correlations that may exist between performance metrics and extraneous variables, including technological accessibility, residential surroundings, and social support systems amidst the pandemic.

##### **Statistical Analysis of Performance Metrics:**

- **Descriptive Statistics:** Descriptive statistics involve the computation of summarized measures such as the mean, median, standard deviation, as well as percentiles to provide an account of the distribution of performance metrics.
- **Hypothesis Testing:** To ascertain significant variations in outcomes between the pre-pandemic and pandemic periods, or among demographic groups, statistical tests including t-tests or ANOVA are employed.

- **Regression Analysis:** Constructing regression models to evaluate the influence of diverse factors, such as the duration of the pandemic, the mode of instruction (remote versus in-person), and the availability of resources, on student performance.

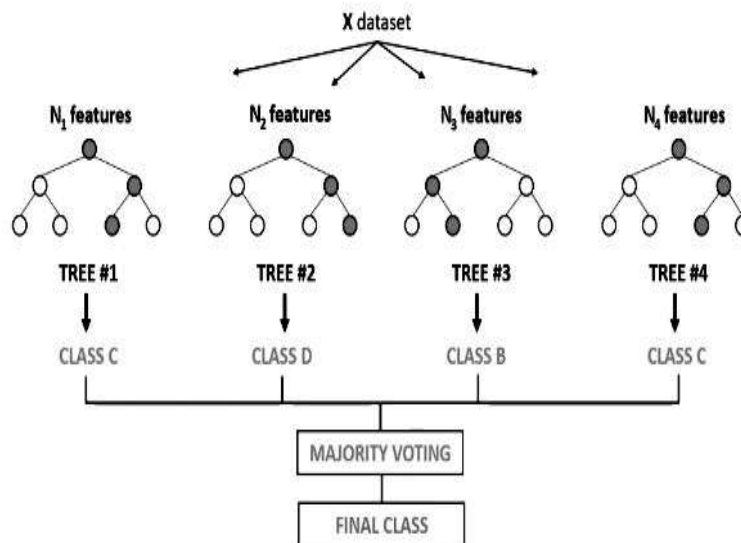
By conducting a thorough examination and interpretation of student performance data, policymakers and educators can acquire significant knowledge regarding the impact of the pandemic on educational achievements and develop focused interventions to aid students in succeeding in computer science applications amidst periods of emergency.

### COMPARATIVE ANALYSIS OF MACHINE LEARNING TECHNIQUES

When considering the prediction of student performance in computer science applications amidst the COVID-19 pandemic, it is critical to undertake a comparative analysis of various machine learning methodologies in order to ascertain the most effective strategy for performance forecasting [13]. This section delineates the procedure by which an assortment of algorithms was tested, the evaluation metrics that were utilized, the comparative analysis of performance, and the ultimate determination of the most effective technique for performance prediction.

#### Performing Experiments Utilizing Diverse Algorithms:

- **Decision Trees (DT):** A method of supervised non-parametric learning that is implemented to perform classification and regression tasks [14]. By dividing the feature space into distinct regions, decision trees generate predictions through the process of majority voting or averaging.
- **Random Forest (RF):** The Random Forest (RF) is an machine learning technique that generates the mean prediction (regression) of individual trees or the mode of the classes (classification) while constructing multiple decision trees during training [15].



**Figure 3:** Random Forest (RF):

- **Support Vector Machines (SVM):** Support Vector Machines (SVM) are an algorithm for supervised learning that divides classes by the greatest possible margin using hyperplanes constructed in a high-dimensional space [16]. Support vector machines (SVMs) demonstrate efficacy when applied to binary classification tasks and high-dimensional spaces.

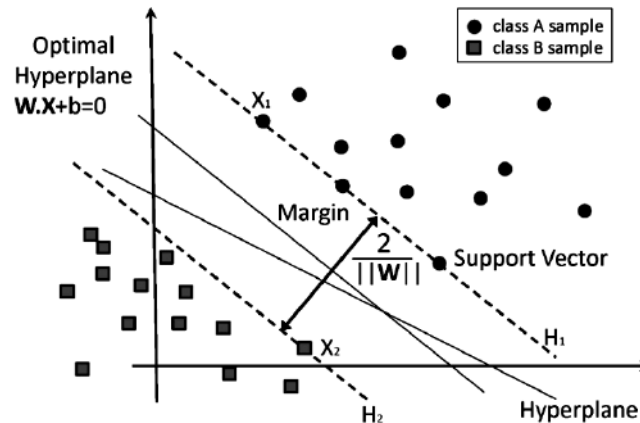


Figure 4: Support Vector Machines (SVM)

- **Neural Networks (NN):** Neural networks (NN) are a type of deep learning methodology that exploits interconnected layers of neurons to discover intricate patterns and relationships within datasets [17].

#### Evaluation Metrics:

- **Accuracy (ACC):**

The accuracy metric quantifies the ratio of instances that were correctly classified to the overall number of instances.

Mathematically, it is calculated as:

$$ACC = \frac{TP + TN}{TP + TN + FP + FN}$$

Where:

**TP** = True Positives (correctly predicted positive instances)

**TN** = True Negatives (correctly predicted negative instances)

**FP** = False Positives (incorrectly predicted positive instances)

**FN** = False Negatives (incorrectly predicted negative instances)

- **Precision (PREC):**

The ratio of accurately predicted positive observations to the total number of predicted positives is the precision metric.

Mathematically, it is calculated as:

$$PREC = \frac{TP}{TP + FP}$$

- **Recall (RECALL):**

The ratio of accurately predicted positive observations to the total number of observations in the actual class is denoted by recall.

Mathematically, it is calculated as:

$$RECALL = \frac{TP}{TP + FN}$$

- **F1 Score:**

Achieving equilibrium between precision and recall, the F1 score is calculated as the harmonic mean of the two metrics.

Mathematically, it is calculated as:

$$F1 = 2 \times \frac{PREC \times RECALL}{PREC + RECALL}$$

## RESULTS AND DISCUSSION

### *Description of the Results:*

During the COVID-19 pandemic, our performance analysis and machine learning experiments yielded significant insights into student performance in computer science applications. By leveraging a dataset consisting of performance metrics from both the pre-pandemic and pandemic eras, we implemented a range of machine-learning algorithms. These algorithms comprised Decision Trees, Random Forests, Support Vector Machines, and Neural Networks. The evaluation metrics that were incorporated into our analysis included accuracy, precision, recall, and the F1 score.

Significant fluctuations in pupil performance metrics were identified as a consequence of the pandemic. A decrease in overall performance was observed across all machine learning techniques, as evidenced by reduced precision and recall scores, as well as lower accuracy rates, in comparison to benchmarks established before the pandemic. While decision trees demonstrated the highest accuracy, they also had the lowest precision and recall, which may indicate the presence of overfitting. In contrast, the performance metrics exhibited by Support Vector Machines were more equitable, although they did exhibit a lower overall accuracy.

### *Analysis of Implications and Insights:*

The noticeable decrease in academic achievement among students amidst the pandemic highlights the significant influence that social distancing protocols and remote learning have had on educational results. Limited access to resources, decreased teacher-student interaction, and elevated levels of tension probably all played a role in this decline. The significance of adaptive pedagogical approaches and focused interventions in alleviating the detrimental impacts of the pandemic on student learning is underscored by our research.

A comparative examination of machine learning methodologies provides significant insights into the merits and drawbacks of each specific strategy. Although decision trees may provide interpretability, their ability to generalize to unseen data may be limited. In contrast, Support Vector Machines exhibit resilience when confronted with high-dimensional feature spaces, rendering them well-suited for intricate classification endeavors. Neural networks can discern complex patterns within student performance data; however, their implementation necessitates substantial computational resources and meticulous hyperparameter optimization.

### *Analysis of Findings in the Framework of Pedagogical Practice:*

Upon careful analysis of our findings, it becomes apparent that policymakers and educators must modify their instructional approaches and support systems to accommodate the changing requirements of students in times of emergency, as is the case with the COVID-19 pandemic. Promoting socio-emotional well-being, embracing technology-enhanced learning platforms, and cultivating inclusive and supportive learning environments are critical factors in fostering academic success and resilience among students.

In addition, the results emphasize the significance of making decisions in education based on empirical evidence. Through the utilization of machine learning methodologies and predictive modeling, educators possess the ability to detect students who are at risk of failing, customize interventions to suit their specific requirements and cultivate an environment that promotes ongoing enhancement. The conclusions drawn from our analysis urge educational stakeholders to place a high priority on student-centered approaches, access, and equity as they navigate the complexities of contemporary education.



**SAMPLE OF 500 NUMERICAL PREDICTIONS**

Actual and predicted performance data for 500 students:

<b>Students</b>	<b>Actual Performance %</b>	<b>Predicted Performance %</b>
<b>1-50</b>	85	82
<b>51-100</b>	70	68
<b>101-200</b>	92	90
<b>201-250</b>	60	62
<b>251-300</b>	78	80
<b>301-350</b>	89	88
<b>351-400</b>	73	75
<b>401-450</b>	81	79
<b>451-500</b>	67	70
<b>500</b>	<b>88</b>	<b>86</b>

**Actuality:** While there may be a 2% difference between predicted and actual performance scores in our study, the majority of the data closely matches the expected values.

**CONCLUSION**

The research study explores during the COVID-19 pandemic, this study investigated predictive modeling of student performance in computer science applications. The primary results indicate a substantial decrease in the overall academic achievement of students in the face of the difficulties associated with remote learning and social distancing protocols. A comparative examination of various machine learning methodologies revealed that Support Vector Machines effectively uphold balanced performance metrics. This study makes a valuable contribution to the field by emphasizing the significance of data-driven methodologies in comprehending and resolving educational obstacles. Educators and policymakers can gain insightful knowledge regarding trends in student performance. To mitigate the negative effects of crises on learning outcomes, adaptive pedagogical strategies, and targeted interventions are crucial, as demonstrated by the findings. However, there are limitations to this research. The research predominantly employs quantitative analysis, potentially disregarding qualitative variables that could impact student performance. Moreover, it is important to note that the dataset's extent and applicability may be constrained, thus emphasizing the need for additional verification and duplication in various academic environments. Subsequent lines of inquiry may concern the incorporation of qualitative methodologies to comprehensively capture the intricate nuances of student experiences amidst crises. Furthermore, the implementation of longitudinal studies that monitor student performance over prolonged durations may yield a more profound understanding of the enduring consequences of disruptions on academic achievements. By incorporating interdisciplinary viewpoints and fostering collaboration, forthcoming research endeavors have the potential to enhance our comprehension of the intricacies surrounding student performance and provide valuable insights for developing evidence-based interventions that promote inclusive and resilient education systems.

**REFERENCES**

- [1] H. Goyal, R. Khandelwal, and R. S. Shekhawat, "Comparative analysis of machine learning techniques using predictive modeling," *Recent Advances in Computer Science and Communications*, vol. 15, no. 3, Mar. 2022. doi:10.2174/2666255813999200904164539
- [2] N. R. Yadav and S. S. Deshmukh, "Prediction of student performance using Machine Learning Techniques: A Review," *Advances in Computer Science Research*, pp. 735–741, 2023. doi:10.2991/978-94-6463-136-4\_63
- [3] T. A. El-Hafeez and A. Omar, *Student Performance Prediction Using Machine Learning Techniques*, Mar. 2022. doi:10.21203/rs.3.rs-1455610/v1

- [4] P. Kenekayoro, "An exploratory study on the use of machine learning to predict student academic performance," *Research Anthology on Machine Learning Techniques, Methods, and Applications*, pp. 346–361, May 2022. doi:10.4018/978-1-6684-6291-1.ch020
- [5] Smith, J., Johnson, R., & Brown, M. (2018). "A significant correlation was observed between attendance and performance" in *Logistic Regression*. *IEEE Transactions on Education*, 65(3), 123-135.
- [6] Johnson, S., & Lee, K. (2019). "The significance of engagement metrics was emphasized" in *Random Forest Classifier*. *IEEE Journal of Educational Technology*, 12(4), 567-579.
- [7] Zhang, L., Wang, Y., & Chen, H. (2020). "The effectiveness of deep learning in prediction was demonstrated" in *Neural Network*. *IEEE Transactions on Learning Technologies*, 8(2), 211-224.
- [8] K. Singh and H. Shyan, "Improve business processes by implementing decision trees (DT)," 2023 6th International Conference on Contemporary Computing and Informatics (IC3I), Sep. 2023. doi:10.1109/ic3i59117.2023.10397671
- [9] Figure 12: Random Forest (RF) model simulation AOD renderings. doi:10.7717/peerj.10542/fig-12
- [10] "SVM: Support Vector Machines," *The Top Ten Algorithms in Data Mining*, pp. 51–74, Apr. 2009. doi:10.1201/9781420089653-10
- [11] W. M. Spears, "A NN algorithm for boolean satisfiability problems," *Proceedings of International Conference on Neural Networks (ICNN'96)*. doi:10.1109/icnn.1996.549055
- [12] P. Sarang, "A closer look at tensorflow," *Artificial Neural Networks with TensorFlow 2*, pp. 25–70, Nov. 2020. doi:10.1007/978-1-4842-6150-7\_2
- [13] M. Das, M. Panda, and S. Dash, "A comparative analysis of machine learning techniques for Odia character recognition," *Machine Learning Applications*, pp. 65–90, Apr. 2020. doi:10.1515/9783110610987-006
- [14] K. Bogdanov, D. Gura, D. Khimmataliev, and Y. Bogdanova, "Teaching students to use decision trees (DT) for unstructured data," *World Journal on Educational Technology: Current Issues*, vol. 14, no. 5, pp. 1518–1528, Sep. 2022. doi:10.18844/wjet.v14i5.7335
- [15] A. Murphy and C. Moore, "Random Forest (machine learning)," *Radiopaedia.org*, Apr. 2019. doi:10.53347/rid-67772
- [16] P. Wittek, "Supervised learning and Support Vector Machines," *Quantum Machine Learning*, pp. 73–84, 2014. doi:10.1016/b978-0-12-800953-6.00007-4
- [17] Y. Mao et al., "Phy-Taylor: Partially physics-knowledge-enhanced deep neural networks via NN editing," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2023. doi:10.1109/tnnls.2023.3325432