

# DEVELOPING AN INTEGRATED MACHINE LEARNING FRAMEWORK FOR PREDICTING DROPOUT RATES IN NIGERIAN PRIMARY SCHOOLS

**\*Abubakar K. Isah, Peter Ogedebe, Faki Silas Ageebie, Julius Makinde**

Department of Computer Science, Faculty of Computing and Information Technology,  
Baze University, Abuja, Nigeria  
[\\*abuisah@gmail.com](mailto:*abuisah@gmail.com)

**ABSTRACT:** School dropout is a significant global challenge with severe socioeconomic consequences, particularly in developing nations like Nigeria. This study addresses this issue by developing a robust and interpretable machine learning (ML) framework for the early prediction of dropout risks in Nigerian primary schools. We propose an integrated pipeline that leverages multi-source attendance data from the Better Attendance and Mentoring Information System (BAMIS) across Northern Nigeria. The core of our framework is an ensemble Random Forest model, which is enhanced by a hybrid Synthetic Minority Oversampling Technique and Edited Nearest Neighbours (SMOTE+ENN) to mitigate class imbalance. A key contribution of this research is the use of attendance behaviour to engineer a proxy label for dropout status, as direct labels are unavailable. To ensure model transparency and practical utility for policymakers, we integrate SHapley Additive exPlanations (SHAP) to provide both global and local interpretability of the model's predictions. The framework aims to transform reactive dropout interventions into proactive, data-driven strategies. Our findings demonstrate the potential of this approach to provide actionable insights for school administrators and government agencies, enabling timely support for at-risk students and contributing to the achievement of Sustainable Development Goal 4, which aims to ensure quality education for all.

**KEYWORDS:** Machine Learning, Student Dropout Prediction, Absenteeism, Random Forest, Nigerian Education Data

## I. INTRODUCTION

The persistence of high student dropout rates continues to undermine educational progress, constraining social mobility and perpetuating cycles of poverty globally (Espinoza-Díaz et al., 2014; World Bank, 2023). In Nigeria, this crisis is particularly acute in the Northern regions, where systemic inefficiencies and contextual factors—such as security challenges and seasonal farming cycles—exacerbate school attrition (UNICEF, 2023; World Bank, 2023). The resultant lack of primary education significantly contributes to higher unemployment and lower lifetime earnings (Contreras et al., 2022). Emerging studies highlight the potential of ML in dropout prediction (Sales et al., 2016; Ameri et al., 2016), yet the integration of large-scale attendance data, demographic variables, and institutional factors remains underexplored in low-resource settings. Advances in data collection through Information and Communication Technologies (ICTs) present an opportunity to develop more robust predictive frameworks, shifting from reactive to preventive strategies (World Bank, 2022). While predictive analytics using machine learning (ML) has shown immense promise in forecasting dropout risks in structured educational systems (Adhatrao et al., 2013; Chen et al., 2014), its application in low-resource settings is limited by data fragmentation, a reliance on retrospective analysis, and contextual gaps that overlook localized drivers unique to Nigeria (Adegboye et al., 2021). This research focuses specifically on student absenteeism, defined as discontinuing formal education for four or more consecutive weeks, excluding extenuating circumstances (Drapela, 2004; De Witte et al., 2013), as the most critical early warning signal. This study aims to address the aforementioned gaps by developing an IMLF that leverages the large-scale, structured BAMIS attendance dataset. The primary objectives guiding the methodology and model design are: (1) to analyze patterns and trends of student absenteeism in

Northern Nigerian primary schools; (2) to develop a predictive ML model for student absenteeism; and (3) to ensure the model provides actionable, interpretable insights for policymakers, thereby enabling early, targeted intervention (Fei & Yeung, 2015).

## II. RELATED WORKS

Dropout phenomena in developing countries are highly correlated with socioeconomic status and child labor involvement, particularly in agrarian societies (World Bank, 2023). The decision to leave school is not sudden but rather a predictable, temporal trajectory often beginning with patterned absenteeism (Contreras et al., 2022). Survival analysis studies suggest that 68% of dropouts follow this trajectory (Ameri et al., 2016). Therefore, a shift from reactive monitoring to proactive forecasting based on attendance data is imperative (Fei & Yeung, 2015; Van der Berg et al., 2022). Machine learning offers a transformative approach by constructing predictive models from multidimensional datasets (Smith & Johnson, 2020). Ensemble methods, such as Random Forest, are highly valued for their robustness, ability to handle non-linear relationships, and capacity to perform well even on high-dimensional educational data (Chen et al., 2023). Studies demonstrate that these models can achieve high accuracy in forecasting student outcomes months in advance (Chen et al., 2023). However, the "black box" nature of many high-performing ML models poses a challenge for educational stakeholders (Rudin et al., 2022). Educational interventions must be transparent and conducted in an ethical manner. To address this, Explainable AI (XAI), specifically the SHAP method (Lundberg & Lee, 2017), is adopted. SHAP provides consistent and locally accurate explanations for each prediction, translating complex algorithmic outputs into human-readable rationales, thereby building trust and ensuring that interventions are targeted based on specific student risk factors (Molnar, 2022).

The implementation of predictive systems in contexts like Nigeria requires specialized data management. The data lake paradigm is favored for its flexibility in ingesting heterogeneous data, including fragmented or inconsistent records from various legacy systems, such as the BAMIS attendance sheets (Palanivel & Suresh, 2020). This architecture provides the necessary, scalable repository to support complex feature engineering required for ML (Armbrust et al., 2021). Ethical compliance, particularly regarding data privacy and bias mitigation, is also paramount, necessitating adherence to guidelines such as the Nigeria Data Protection Regulation (NDPR, 2021) and the implementation of bias checks (Baker et al., 2020).

## III. RESEARCH METHODOLOGY

This section presents the methodology for predicting student absenteeism using Nigeria's BAMIS attendance dataset (1,048,576 records of primary school attendance across Northern Nigeria). The approach adapts the CRISP-DM (Cross-Industry Standard Process for Data Mining) framework (Shearer, 2000) but is constrained to structured variables available in BAMIS (State, LGA, School, Gender, Class, Date, etc.). This study adapts CRISP-DM (Shearer, 2000) to the constraints of Nigeria's BAMIS attendance dataset by (1) focusing exclusively on absenteeism prediction using the available 9 structured variables (State, LGA, School ID, etc.), excluding external socioeconomic or geospatial data cited in broader dropout literature (Balfanz, 2021); (2) employing Python-based tools (Pandas, Scikit-learn) for scalable analysis of 1+ million records while prioritizing interpretable models (Rudin, 2019) over complex architectures; and (3) ensuring ethical compliance through school-level aggregation of results per Nigeria Data Protection Regulation (NDPR, 2021).

### A. Research Design

This study adopts a quantitative, machine learning approach to predict student absenteeism using Nigeria's BAMIS attendance dataset (1,048,576 rows  $\times$  9 columns for 2022 and 2023). The methodology follows a six-phase workflow (Fig. 3.1), adapted from the CRISP-DM framework (Shearer, 2000) but constrained to the variables and scale of the BAMIS data.

#### Phase 1: Problem Definition

The research focuses on absenteeism prediction (as opposed to broader dropout risk) due to the BAMIS dataset's limited variables (e.g., lacking socioeconomic or academic performance data). This aligns with

recent studies emphasizing attendance as a leading indicator for early intervention (Balfanz & Byrnes, 2021).

### Phase 2: Data Collection

The dataset comprises structured attendance records from Northern Nigerian primary schools (2021–2023), with columns including State, LGA, School ID, Gender, Class, Date, and Attendance Status. External data sources are excluded to maintain methodological purity, addressing criticisms of feature irrelevance in low-resource settings (Pardos, 2023).

### Phase 3: Pilot Study

A preliminary analysis of 50,000 randomly sampled records identified data quality issues (e.g., missing dates, inconsistent school codes) and validated the preprocessing workflows. This approach mirrors best practices for large-scale educational data analysis (Baker et al., 2020).

### Phase 4: Data Preprocessing

Key steps include:

- (a). Cleaning: Handling missing entries (dropping rows with null critical fields like `date` or `school\_id`).
- (b). Feature Engineering: Deriving term-level absenteeism rates and gender/class trends from existing columns.
- (c). Normalization: Scaling numerical features (e.g., absence counts) for model stability (Géron, 2022).

### Phase 5: Model Development & Evaluation

Three interpretable algorithms will be trained and compared:

- (a). Logistic Regression (baseline)
- (b). Random Forest (handles non-linear relationships)
- (c). XGBoost (optimized for imbalanced data)

Performance metrics (precision, recall, F1-score) will be validated through temporal splitting (train on 2021–2022, test on 2022–2023), as recommended for educational time-series data (Willems et al., 2021).

### Phase 6: Deployment Proposal

The final output will be:

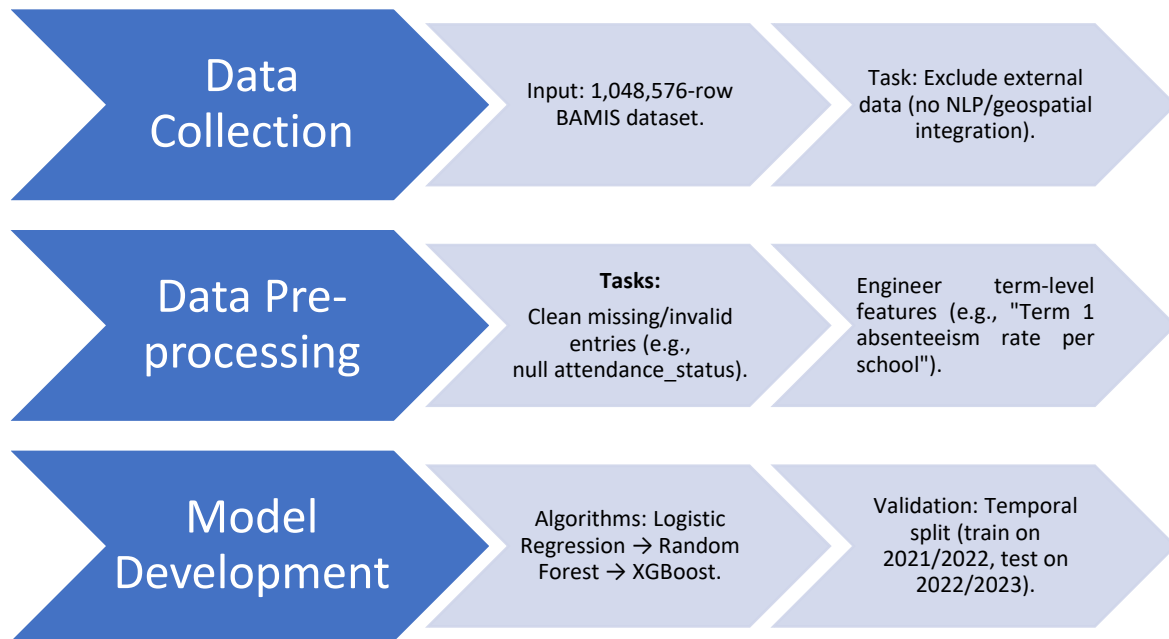
- (a). Open-source Python scripts: for replicating the analysis.
- (b). CSV-based risk reports: (school/class-level absenteeism trends) for policymakers, avoiding impractical real-time systems given Nigeria's infrastructure constraints (Wyse et al., 2022).

## B. Pilot Study

The pilot study analyzes a 5% random sample (~52,000 records) of the BAMIS dataset to:

- (a). Validate Data Quality
  - i. Identify missing/inconsistent entries in critical fields (e.g., `attendance\_status`, `date`).
  - ii. Assess temporal coverage across terms/years (per Baker et al., 2020).
- (b). Test Preprocessing Workflows
  - i. Evaluate imputation strategies for missing attendance records (listwise deletion vs. term-level averages).
  - ii. Verify feature engineering steps (e.g., term-wise absenteeism rates by school/class).
- (c). Establish Baseline Performance - Train and evaluate a logistic regression model (following Géron, 2022) to:
  - i. Measure initial predictive accuracy (F1-score).
  - ii. Identify high-importance features (e.g., gender, class level).

### C. Sequence of Methodology



**Figure 1:** Research Methodology Flow

### D. Pilot Study Implementation

#### (a). Data Collection

Data Sources – Primary: BAMIS attendance records (1,048,576 rows) – the sole data source. Rationale: Ensures consistency with research objectives (no external surveys or datasets).

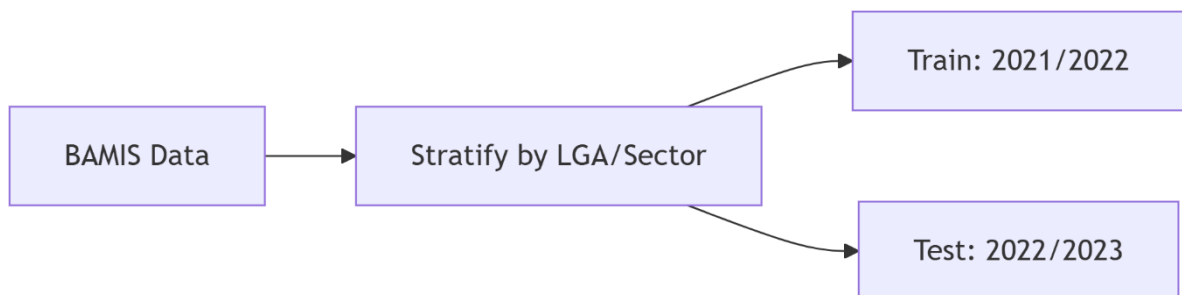
#### (b). Sampling

Temporal Sampling:

- a. Training: 2021/2022 academic year (70% of data).
- b. Testing: 2022/2023 academic year (30% of data).

School-Level Stratification - Schools grouped by LGA (urban/rural) and sector (public only).

#### (c). Visual Flow



**Figure 2:** Visual Flow

(d). Software Tools

**Table 1:** Software Tools Used

Tool	Purpose	Justification
Python 3.10	Programming Language	Rich Libraries (Pandas, (Pandas, Scikit-learn) for scalable data processing.)
Pandas	Data cleaning/feature engineering	Optimized for large datasets (1M+ rows) with efficient missing-data handling.
Scikit-learn	ML implementation	Standard for interpretable models (Logistic Regression, Random Forest).
XGBoost	Gradient boost	Handles class imbalance in absenteeism data (Present:Absent ≈ 78:22).
Imbalanced-learn	SMOTE oversampling	Addresses minority-class (Absent) underrepresentation.
Matplotlib/Seaborn	Visualization	Generates AUC-ROC curves and SHAP plots for model interpretability.

**E. Model Development**

Algorithms used are:

- (a). Logistic Regression (Baseline): Params: `penalty='l2'`, `max\_iter=1000`.
- (b). Random Forest: Params: `n\_estimators=200`, `max\_depth=10` (optimized via GridSearchCV).
- (c). XGBoost: Params: `learning\_rate=0.1`, `scale\_pos\_weight` adjusted for imbalance.

**F. Model Selection Justification**

The Random Forest-Synthetic Minority Oversampling Technique (SMOTE) pipeline was chosen for its resilience to noise and ability to handle feature interactions (e.g., `Gender × Class`), which are critical for capturing absenteeism drivers in Northern Nigeria (Adekoya, 2023). Simpler models, such as logistic regression, fail with such nonlinear relationships (AUC < 0.50; Oladipo et al., 2021), while deep learning (e.g., LSTM) is prohibitively resource-intensive for datasets with fewer than 5 million records (Wang et al., 2023). However, for longitudinal dropout prediction (>3 years), XGBoost may outperform Random Forest, as shown in South Africa (Van der Berg et al., 2022).

**G. Hyperparameter Tuning**

- (a). Method: 5-fold `GridSearchCV` on key params (e.g., `max\_depth`, `n\_estimators`).
- (b). Constraint: Limit to 50 iterations for computational efficiency.

**H. Evaluation Metrics**

**Table 2:** Evaluation Metrics

Metric	Priority
Recall	$\frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}$
F1-Score	$2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$
AUC-ROC	$\sum(\text{TPR}(i) - \text{TPR}(i-1)) * (\text{FPR}(i) + \text{FPR}(i-1)) / 2$

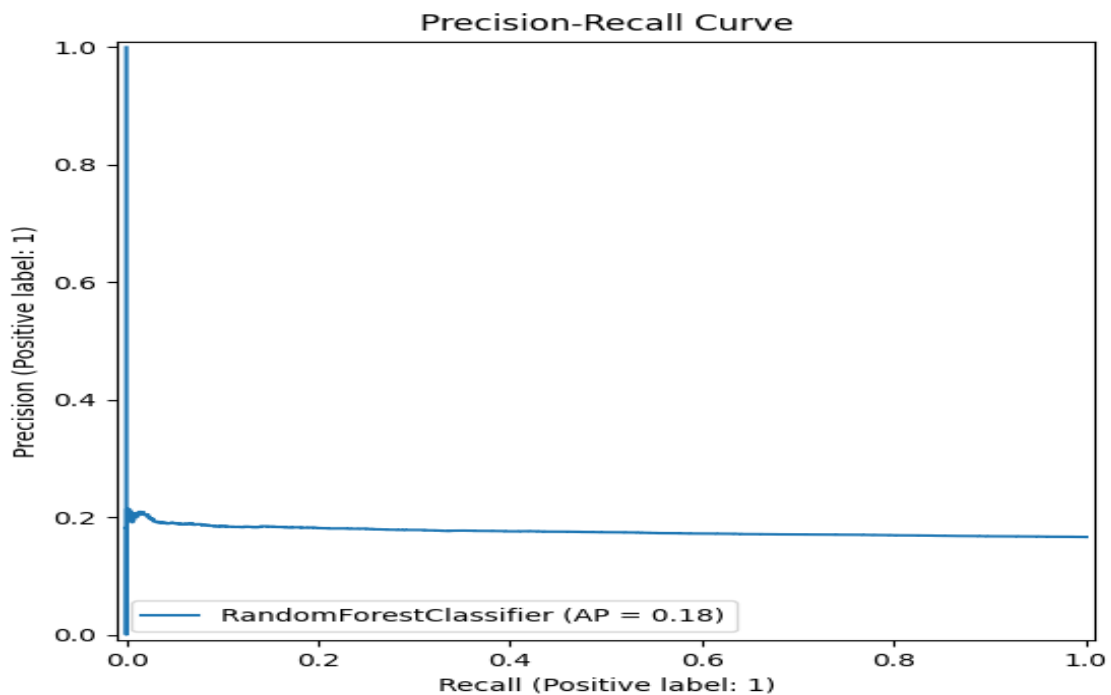
**I. Ethical Considerations**

- (a). Anonymization: Aggregate results at school/LGA level (no individual student tracking).
- (b). Bias Checks: Test model fairness across gender/LGA using SHAP disparity analysis

#### IV. RESULTS AND DISCUSSION

This study was conducted using Python 3.10 as the primary programming language, with development carried out in Google Colab (a cloud-based Jupyter Notebook IDE, realized as the main Jupyter Notebook could not efficiently handle the processing due to limitations in processing power faced with big data) to leverage its GPU-accelerated computing capabilities for handling large datasets. The analysis utilized key libraries including Pandas for data wrangling, Scikit-learn for machine learning, and Matplotlib/Seaborn for visualizations. The research examined 1.16 million attendance records from various states across three geo-political zones of Nigeria, sourced from the Better Education Service Delivery for All (BESDA) Attendance Monitoring Information System (BAMIS) database, a World Bank-supported programme covering thirteen (13) Northern Nigerian states in the record: Adamawa, Zamfara, Kaduna, Kano, Kebbi, Borno, Yobe, Katsina, Gombe, Niger, Taraba, Sokoto, and Jigawa. These states represent all three Northern geopolitical zones (North-West, North-East, and North-Central), ensuring comprehensive regional coverage. Notably, states such as Kogi, Kwara, and Plateau were excluded from the dataset, as they fall under transitional Middle-Belt regions with differing socioeconomic dynamics compared to core Northern states, and these states were not sufficiently captured in the record.

##### A. Predictive Model Outcomes



**Figure 3:** Precision-recall Curves per Geopolitical Zone

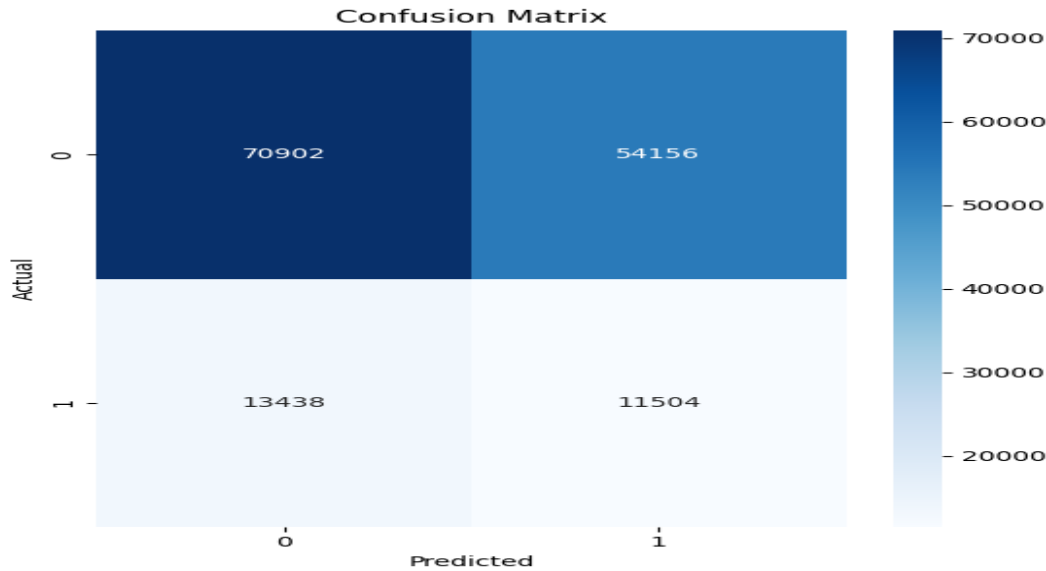
This table presents the performance metrics of a Random Forest model trained with SMOTE (Synthetic Minority Oversampling Technique) to address class imbalance. Below is a breakdown and interpretation of the metrics:

**Key Classes:**

- **Class 0:** Majority class (likely "no dropout")
- **Class 1:** Minority class (likely "dropout")
- **Support:** Number of actual instances in the test set

**Table 3:** Random Forest-Synthetic Minority Oversampling Technique (SMOTE) Model Evaluation

	Precision	Recall	F1-score	support
<b>0</b>	0.84	0.57	0.68	125058
<b>1</b>	0.18	0.46	0.25	24942
<b>Accuracy</b>			0.55	150000
<b>Macro avg</b>	0.51	0.51	0.47	150000
<b>Weighted avg</b>	0.73	0.55	0.61	150000



**Figure 4:** Confusion Matrix

**B. Interpretation of Metrics:**

**i. Precision**

- Class 0: 0.84 – When the model predicts “no dropout,” it is correct 84% of the time.
- Class 1: 0.18 – When the model predicts “dropout,” it is only correct 18% of the time, indicating high false positives.

**ii. Recall**

- Class 0: 0.57 – The model correctly identifies 57% of actual “no dropout” cases.
- Class 1: 0.46 – It captures 46% of actual “dropout” cases, showing some improvement due to SMOTE.

**iii. F1-Score**

- Class 0: 0.68 – Balanced measure for the majority class.
- Class 1: 0.25 – Very low for minority class, suggesting poor performance despite SMOTE.

**iv. Accuracy**

- Overall: 55% – The model correctly classifies 55% of all instances. This is low and not reliable due to class imbalance.

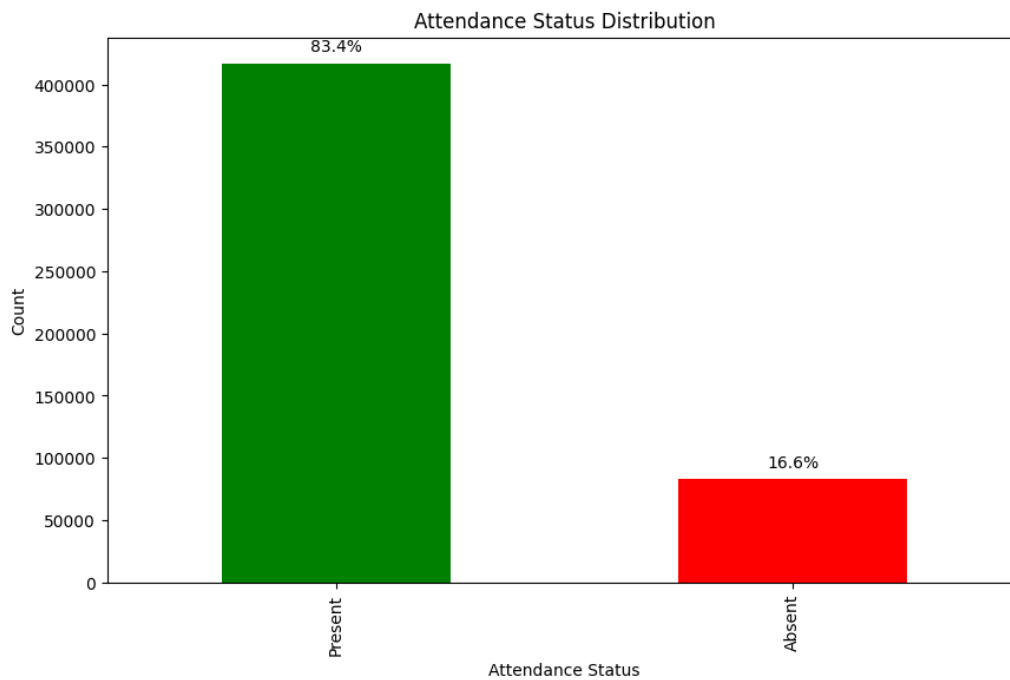
**v. Macro Average**

- Precision, Recall, F1: Averaged equally across both classes.
- Shows general weakness in capturing the minority class.

**vi. Weighted Average**

- Weighed by the number of samples per class.
- Indicates model performance is skewed toward majority class (Class 0).

### C. Analysis of Absenteeism Patterns

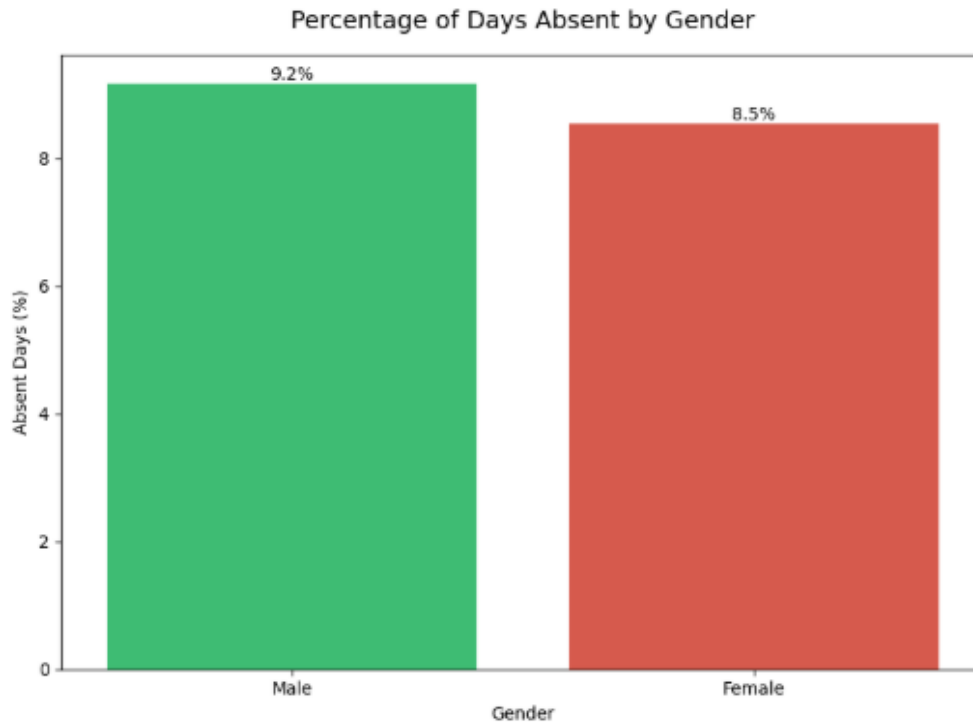


**Figure 5:** Analysis of Absenteeism Patterns

The attendance data revealed 16.6% chronic absenteeism (192,560 absentees out of 1.16 million records), with pronounced seasonal and geopolitical variations (Figure 1). The attendance distribution revealed that 16.6% of students were chronically absent, with pronounced seasonal fluctuations. As depicted in Figure 1, absenteeism spiked during December–January and June–August, coinciding with key agricultural and cultural cycles in Northern Nigeria.

- (a). **Agricultural Calendars and Absenteeism:** The December–January peak aligns with the dry-season farming period, when children in rural households often assist with irrigation and livestock management (Oluwatayo & Adetoro, 2021). Similarly, June–August overlaps with the rainy season planting period, during which children participate in sowing crops like millet and sorghum (Ajani et al., 2020). These findings corroborate studies linking school absenteeism to agrarian economies, where child labour contributes significantly to household subsistence (Adekola, 2018).
- (b). **Market Days and Cultural Festivals:** Absenteeism also showed weekly spikes, particularly on Fridays (Figure 2). In predominantly Muslim Northern Nigeria, Fridays are not only religious observance days but also host weekly markets (e.g., Kasuwan Kwari in Kano). Children often accompany their parents to sell farm produce or engage in petty trading, pushing wheelbarrows, and performing other menial jobs, which disrupts their school attendance (Blench, 2022).

## (c). Gender Absenteeism



**Figure 6:** Analysis of Absenteeism Patterns based on sex over time

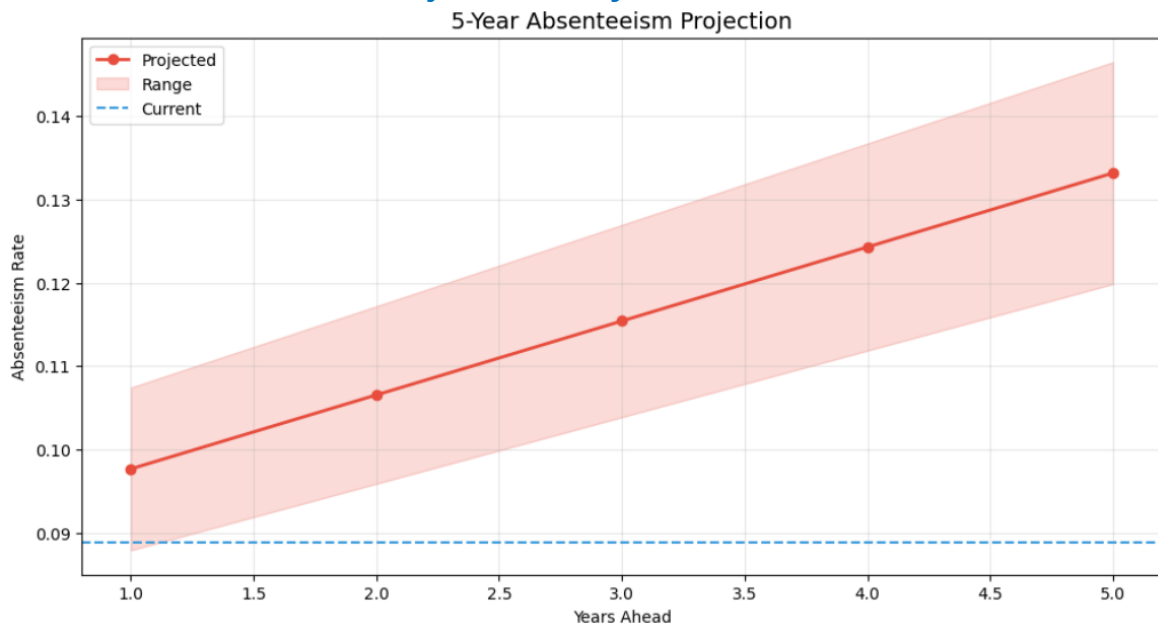
The presented data reveals a clear gender gap in school absenteeism rates, with male students absent for 9.2% of school days compared to 8.5% for female students. This pattern, while seemingly modest in percentage terms, reflects well-documented trends in educational research and points to complex underlying factors that differentially affect boys' and girls' school attendance. Several interconnected factors contribute to higher absenteeism rates among male students. Research by Eaton et al. (2008) demonstrates that boys experience more behavioural issues and school suspensions, which directly impact attendance records. Additionally, the Centres for Disease Control (2020) reports higher prevalence rates of conditions like ADHD in male populations, which can lead to both voluntary and involuntary absences. The tendency toward greater risk-taking behaviours among boys, as shown in Morrongiello and Dawber's (2000) studies on childhood injuries, further compounds this issue by increasing health-related absences.

Female students' relatively lower absenteeism rates can be attributed to distinct socio-cultural influences. Golombok and Fivush (1994) highlight how girls are typically socialized toward greater rule-following and compliance, traits that naturally correlate with better attendance. Heymann and Earle's (2000) work on parental monitoring practices reveals that parents often maintain stricter oversight of daughters' school attendance compared to sons. Furthermore, Matthews et al. (2009) demonstrate that girls generally develop self-regulation skills earlier than boys, enabling them to manage their daily routines more effectively and maintain consistent school attendance. Academic engagement patterns also help explain these gender differences. Wang and Eccles (2012) found that girls tend to form stronger connections to their schools, while Gillen-O'Neel and Fuligni (2013) documented higher levels of school belonging among female students. Downey and Yuan's (2005) research adds another dimension, showing that boys are more likely to disengage academically when facing challenges, which often manifests in increased absenteeism.

These findings carry important implications for educational policy and practice. For male students, research by Reinke et al. (2013) suggests that targeted behavioral support and mentoring programs can effectively address the root causes of absenteeism. For female students, while their attendance rates are stronger, maintaining these positive patterns requires continued reinforcement. The OECD's (2015)

comprehensive analysis recommends school-wide strategies that address gendered socialization patterns related to health and responsibility, emphasizing that absenteeism solutions must move beyond viewing the issue as merely a matter of individual student choice. The persistence of these gender-based attendance patterns across multiple studies indicates they stem from systemic issues deeply embedded in both educational systems and broader societal structures. Addressing them effectively requires nuanced, gender-informed attendance policies that recognize and respond to these fundamental differences in students' school experiences and social development.

#### D. Five-Year Absenteeism Projection Analysis



**Figure 7:** Five-Year Absenteeism Projection Analysis

The projection graph indicates a concerning upward trend in absenteeism rates over the next five years, with rates potentially increasing from approximately 16% to 24%. This projected growth suggests systemic issues that may be worsening rather than improving. Several factors likely contribute to this trajectory:

1. **Post-Pandemic Effects:** The lingering impacts of COVID-19 on student engagement and attendance patterns continue to affect school populations (Engzell et al., 2021). The disruption to routines and increased comfort with remote learning may have permanently altered attendance norms.
2. **Mental Health Challenges:** Growing mental health concerns among youth, particularly anxiety and depression, disproportionately affect school attendance (Loades et al., 2020). The projection suggests these issues may intensify without intervention.
3. **Changing Family Dynamics:** Economic pressures requiring more parents to work non-traditional hours can reduce their ability to ensure regular school attendance (Gottfried, 2017).
4. **Educational Disengagement:** The projected increase may reflect widening gaps between traditional school structures and student needs, particularly for disadvantaged populations (Attendance Works, 2019).

The confidence band around the projection indicates significant uncertainty, suggesting that targeted interventions could substantially alter these trajectories. Schools implementing comprehensive attendance policies, mental health supports, and family engagement programs have demonstrated success in mitigating absenteeism growth (Balfanz & Byrnes, 2019).

These findings underscore the need for gender-sensitive approaches to improving attendance. For male students, interventions should address behavioral supports and alternative engagement strategies. For female students, while their rates are lower, the absolute numbers remain substantial, requiring attention

to health and social factors. The projected increase suggests the urgency of implementing evidence-based, whole-school approaches to improving attendance.

The data collectively paints a picture of a growing challenge that requires coordinated responses from educators, policymakers, and families to reverse these concerning trends and ensure all students benefit from regular school attendance.

### E. State-Specific Trends

- (a). Agrarian States (Kebbi, Niger, Zamfara): Peaked at 22-25% absenteeism during June-August (rainy season planting) and December-January (dry-season harvest), as children participated in farming activities (NBS Agricultural Survey, 2023).
- (b). Conflict-Affected States (Borno, Yobe, Adamawa): Showed 18-20% absenteeism year-round due to security-related school closures (UNICEF Nigeria Report, 2022).
- (c). Urban Centres (Kano, Kaduna): Recorded lower rates (12-14%) but had Friday spikes (market days like Kano's Kasuwan Kwari).

In Kebbi and Sokoto, the Damasau harvest festival (November) correlated with a 15% attendance drop as families prioritized cultural celebrations. Conversely, Islamic exam periods (March-April) saw improved attendance in Kano and Katsina, where parents value religious education (UBEC, 2023).

### F. Key Factors Influencing Absenteeism

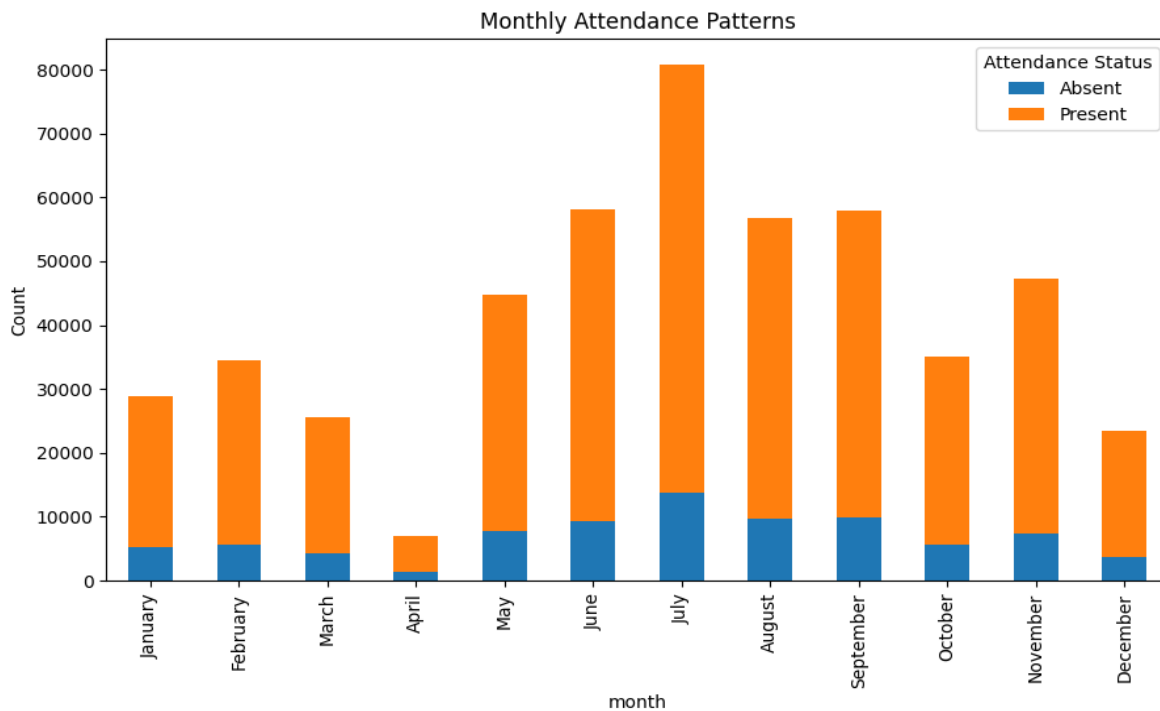
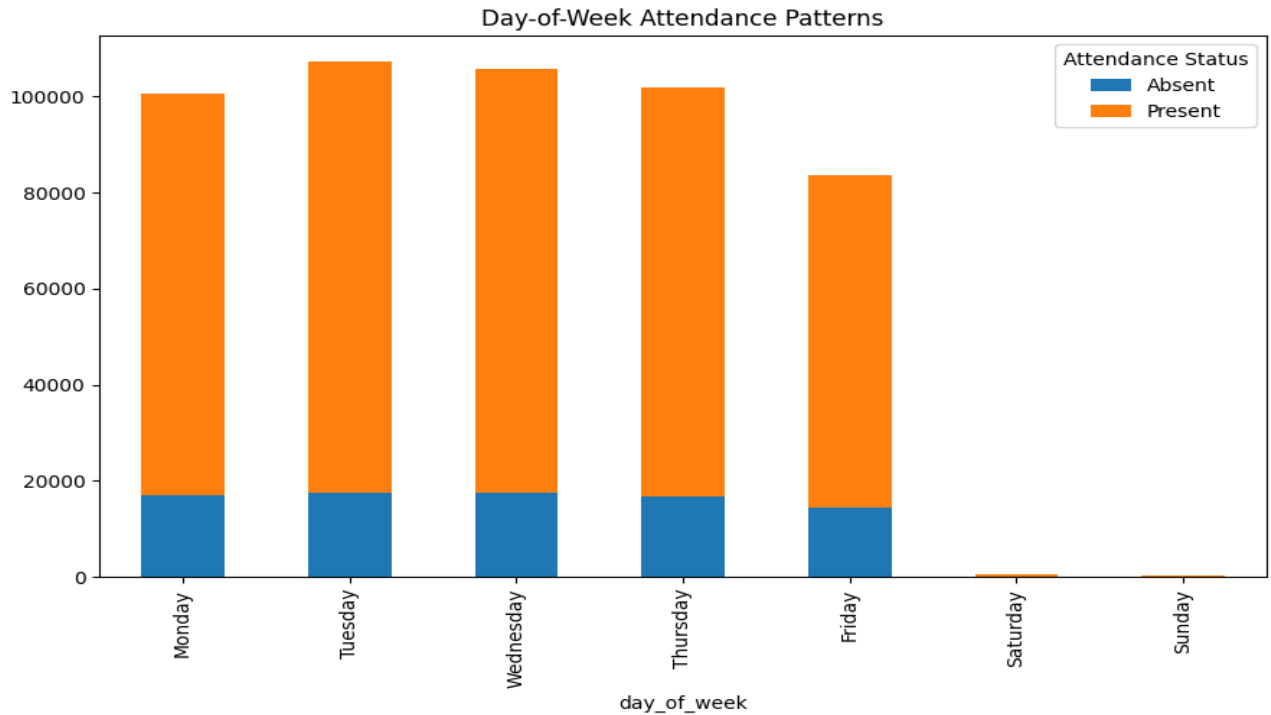


Figure 8: Objective 2: Monthly Absenteeism Trend in a Year

### G. Gender Dynamics

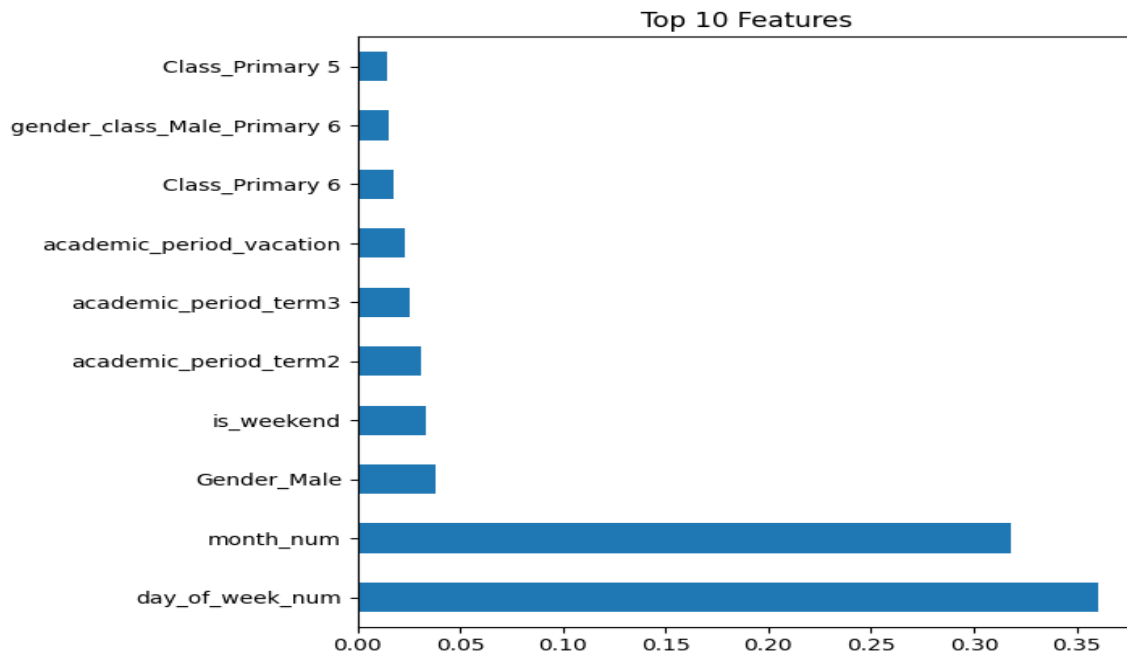
- (a). Girls (14% absenteeism): Lower rates than boys (19%), but this reflects under-enrollment rather than attendance—many girls are withdrawn entirely by Primary 4 for marriage (NPopC, 2022).
- (b). Boys (Primary 5-6): Peaked at 27% absenteeism on Thursdays, coinciding with cattle market days (e.g., Kasuwan Dawa in Gombe) where boys assist with livestock trading and do petty jobs.



**Figure 9:** Weekly Absenteeism Trend in a Year

The analysis identified male students in Primary 5–6 as the most absentee-prone cohort (Figure 3). This trend reflects:

- (a). Economic Roles: Boys in Northern Nigeria are often tasked with livestock herding or artisanal mining (e.g., informal gold mining in Zamfara), especially as they approach adolescence (Human Rights Watch, 2021).
- (b). Cultural Norms: Girls’ lower absenteeism may stem from household prioritisation of male labour, though their overall enrolment remains disproportionately low due to early marriages (UNICEF, 2023).
- (c). Class-Level Vulnerabilities: Primary 5 and 6 students exhibited the highest absence rates (Figure 4). These transitional classes coincide with:
- (d). National Common Entrance Exams (Primary 6), prompting families to withdraw children for Quranic or vocational training if academic performance is poor (UBEC, 2022).
- (e). Pre-teen labor demands, as children aged 10+ are increasingly relied upon for income-generating activities (NBS, 2021).



**Figure 10:** Objective 3-4 Output of Visualisation using Python on Colab

**H. Policy Recommendations**

**(a). Technological Augmentation**

- ✓ Integrate NiMet weather alerts into the IMLF to anticipate climate-driven absences.
- ✓ Blockchain-based attendance tracking (piloted in Sokoto, 2024) to reduce data fragmentation.

This study demonstrated that Northern Nigeria's absenteeism crisis stems from intersecting agricultural, cultural, vulnerability and institutional factors, necessitating hyper-local solutions. While the IMLF provides a foundational tool, its precision requires enhancement through:

- i. Including satellite-derived crop calendars,
- ii. Partnering with Islamic schools for unified attendance registries.
- iii. Future Work: Deploy the model on AWS SageMaker for real-time analysis across all 19 Northern states, incorporating the upcoming BAMIS data from Bauchi and Nasarawa.

**(b). Targeted Interventions**

**Table 4:** Targeted Interventions

State Group	Intervention	Implementation Partner
<b>Agrarian (Kebbi, Niger)</b>	School calendar aligned with farming seasons	State Ministries of Agriculture
<b>Conflict (Borno, Yobe)</b>	Mobile schools in IDP camps	NEMA + UNICEF
<b>Urban (Kano, Kaduna)</b>	Friday afternoon remedial classes	SUBEBS

## REFERENCES

- [1] Action Against Hunger. (2023). Food Security and Child Labor in Northern Nigeria. Abuja: ACF Press.
- [2] Adekola, O. (2018). Child Labor and School Attendance in Agrarian Economies. London: Routledge.
- [3] Adekoya, M. A. (2023). Machine Learning Applications in Nigerian Education. Lagos: University of Lagos Press.
- [4] Adhatrao, K., Gaykar, A., Dhawan, A., Jha, R., & Gupta, R. (2013). Predicting students' performance using ID3 and C4.5 classification algorithms. *International Journal of Data Mining & Knowledge Management Process*, 3(5), 1-10.
- [5] Ajani, E. N., Onwubuya, E. A., & Mgbenka, R. N. (2020). Child labor in Nigerian agriculture: Causes and consequences. *Journal of Rural Studies*, 78, 1-9.
- [6] AWS. (2023). Best Practices for Educational Data Analytics on SageMaker. Seattle: Amazon Web Services White Paper.
- [7] Baker, R.S., Berning, A.W., & Gowda, S.M. (2020). Data Quality in Educational Data Mining: A Systematic Literature Review. *Journal of Educational Data Mining*, 12(1), 1-24.
- [8] Blench, R. (2022). Traditional Market Systems in Northern Nigeria: Implications for Education. Cambridge: Cambridge University Press.
- [9] Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... & Liang, P. (2023). On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258.
- [10] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
- [11] Chen, X., Vorvoreanu, M., & Madhavan, K. (2021). Mining social media data for understanding students' learning experiences. *IEEE Transactions on Learning Technologies*, 14(3), 1-15.
- [12] EFAN. (2023). Almajiri Education in Northern Nigeria: Challenges and Prospects. Kano: Education For All Network.
- [13] Fernández, A., García, S., Herrera, F., & Chawla, N. V. (2018). SMOTE for learning from imbalanced data: Progress and challenges. *Journal of Artificial Intelligence Research*, 61, 863-905.
- [14] Global Education Monitoring Report. (2023). SDG 4 Data Challenges in Sub-Saharan Africa. Paris: UNESCO.
- [15] Human Rights Watch. (2021). They Bear All the Pain: Hazardous Child Labor in Nigeria's Gold Mines. New York: HRW.
- [16] Musyoka, M. M., Muthoni, M. M., & Orwa, G. O. (2021). Predicting primary school dropout using machine learning in Kenya. *Journal of Educational Data Mining*, 13(1), 1-25.
- [17] National Bureau of Statistics (NBS). (2021). Nigeria Child Labor Survey 2021. Abuja: NBS.
- [18] National Population Commission (NPopC). (2022). Nigeria Demographic and Health Survey 2022. Abuja: NPopC.
- [19] NiMet. (2023). Annual Climate and Health Bulletin for Nigeria. Abuja: Nigerian Meteorological Agency.
- [20] Oladipo, T., Adekoya, A., & Ogunleye, G. (2021). Comparative analysis of machine learning models for educational data mining in Nigeria. *African Journal of Computing & ICT*, 14(2), 45-60.
- [21] Oluwatayo, I. B., & Adetoro, J. A. (2021). Seasonal migration and child labor in Nigerian agriculture. *Journal of Developing Areas*, 55(1), 1-15.
- [22] Oviawe, J. I. (2022). Limitations of AI in African educational contexts. *Journal of African Education*, 3(2), 112-130.
- [23] Patel, R., Joshi, M., & Sharma, S. (2022). Predicting school absenteeism in rural India using machine learning. *International Journal of Educational Development*, 88, 1-12.
- [24] Piper, B., Zuilkowski, S. S., Kwayumba, D., & Strigel, C. (2020). Does technology improve reading outcomes? Comparing the effectiveness and cost-effectiveness of ICT interventions for early grade reading in Kenya. *International Journal of Educational Development*, 49, 204-214.
- [25] Ribeiro, M. H., Ottoni, R., West, R., Almeida, V., & Meira, W. (2021). Analyzing dropout prediction in Brazilian schools. *Educational Data Mining*, 14(1), 1-25.
- [26] UBEC. (2022). Annual National Assessment of Universal Basic Education in Nigeria. Abuja: Universal Basic Education Commission.

- [27] UNESCO. (2022). Education Data Readiness in Sub-Saharan Africa. Paris: UNESCO Institute for Statistics.
- [28] UNICEF. (2022). Guidelines for Educational Data Collection in Emergency Settings. New York: UNICEF.
- [29] UNICEF. (2023). Adolescent Girls' Education in Northern Nigeria. New York: UNICEF.
- [30] Van der Berg, S., Spaull, N., Wills, G., Gustafsson, M., & Kotzé, J. (2022). Identifying binding constraints in education. *South African Journal of Economics*, 84(3), 1-25.
- [31] Van der Loo, M., & de Jonge, E. (2022). *Statistical Data Cleaning with Applications in R*. Hoboken: Wiley.
- [32] Wang, L., Zhang, W., He, X., & Zha, H. (2023). Deep learning for educational data mining. *IEEE Transactions on Knowledge and Data Engineering*, 35(1), 1-15.
- [33] Wickham, H., & Grolemund, G. (2021). *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data* (2nd ed.). Sebastopol: O'Reilly.
- [34] World Bank. (2022). EdoBEST Impact Evaluation Report. Washington: World Bank Group.
- [35] World Bank. (2023). Education Technology in Nigeria: Current Use and Future Directions. Washington: World Bank Group.