# A survey of Machine Learning Approaches and Techniques for Student Dropout Prediction

Neema Mduma [a,*], Khamisi Kalegele [b] and Dina Machuve [c]

[a] *Department of Information and Communication Science and Engineering, The Nelson Mandela African Institution of Science and Technology, NM-AIST, Tanzania*
*E-mail: mduman@nm-aist.ac.tz; ORCID: https://orcid.org/0000-0002-4364-3124*
[b] *Commission for Science and Technology, COSTECH, Tanzania*
*E-mail: kalegs03@gmail.com; ORCID: https://orcid.org/*
[c] *Department of Information and Communication Science and Engineering, The Nelson Mandela African Institution of Science and Technology, NM-AIST, Tanzania*
*E-mail: dina.machuve@nm-aist.ac.tz; ORCID: https://orcid.org/*

**Abstract.** School dropout is absenteeism from school for no good reason for a continuous number of days. Addressing this challenge requires a thorough understanding of the underlying issues and effective planning for interventions. Over the years machine learning has gained much attention in addressing the problem of students dropout. This is because machine learning techniques can effectively facilitate determination of at-risk students and timely planning for interventions. To this end, several machine learning algorithms have been proposed in literature. This paper presents an overview of machine learning in education with focus on approaches and techniques for student-dropout prediction. Furthermore, the paper discusses the state of student dropout in developing countries and several performance metrics used by researchers to evaluate machine learning techniques in the context of education with experimental examples. Finally, the paper highlights challenges and future research directions.

Keywords: Machine Learning (ML), Imbalanced learning classification, Secondary education, Evaluation metrics

## 1. Introduction

Reducing student dropout rates is one of the challenges faced in the education sector globally. This problem brought a major concern in the field of education and policy-making communities [8]. A growing body of literature indicates high rates of students dropout of school, especially pronounced in the developing world; with higher rates for girls compared to boys in most parts of the world [68]. In Tanzania, for example, student dropout is higher in lower secondary compared to higher level where girls are much less likely than boys to complete secondary education; 30 percent of girls dropout before reaching form 4 compared to 15 percent for boys [12]. Finding and implementing solutions to this problem has implications well beyond the benefits to individual students. Moreover, enabling students to complete their education means investing in future progress and better standards of life with multiplier effects. To effectively address this problem, it is very crucial to ensure that all students finish their school on time

---

*Corresponding author. E-mail: mduman@nm-aist.ac.tz.

through early intervention on students who might be at risk of dropping classes. This require data-driven predictive techniques that can facilitate determination of at-risk students and timely planning for interventions [24].

Machine learning approaches are one of the well sought solutions to addressing school dropout challenge. Various studies have been conducted in developed countries on developing student predictive algorithms [2, 16, 21]. Moreover, there exist quite a significant body of literature on machine learning based approaches associated with fighting dropouts [5, 42, 65]. The knowledge embodied in literature has the potential to transform the fight against dropout from reactive to proactive. This is a more reality now than ever because the ICTs have already transformed the way we collect and manage data, which is a key ingredient to any intelligent harnessing of useful patterns of recorded events. Despite several efforts done by previous researchers, there are still challenges which need to be addressed. Most of the widely used datasets are generated from developed countries. However, developing countries are facing several challenges on generating public datasets to be used on addressing this problem. The study conducted by [52] used the primary data collected in Kenya, although the dataset is not public available. Besides, Uwezo data on learning [1] is the publicly available dataset which was collected countrywide for primary schools in Tanzania. The dataset focused on individual household data, including education. This work presents an up to date overview of predictive techniques and approaches for addressing the problem of student dropout.

In developing countries, prospects of dropout-free education system are still slim considering the scale of socio economic challenges, which are deemed central to the retention of students in schools especially girls. Increasingly, communities of practitioners and researchers are looking at machine learning approaches as a likely solution for achieving dropout-free schools. In this article, a review of how machine-learning techniques have been used in the fight against dropouts is presented for the purpose of providing a stepping-stone for students, researchers and developers who aspire to apply the techniques. Key intervention points that were identified during our preliminary survey guided the herein presented review. The intervention points included issues related to data preprocessing, choosing an algorithm to predict dropouts, and evaluation metrics. In this article, potential machine learning techniques for the three intervention areas are summarized and also the results of demonstration experiments are presented.

## 1.1. The state of student dropout in developing countries

The issue of student dropout is a serious problem which adversely affects the development of education sector, this is due to a complex interplay of socio-cultural, economic and structural factors [55]. Schooling, according to the human capital theory, is an investment that generates higher future income for individuals [59]. Many developing countries are experiencing high dropout rate of secondary school students as a big challenge which has been considered as a problem for the individual and society [28]. However, less attention is paid to improve quality of education to people belongs to any class. In this regard, a [75] report points out, that about one thirty million children in the developing world denied their right to education through dropping out [18].

In responding to this problem of dropping out and other challenges facing secondary schools, Tanzania as one among developing countries introduced an Education Training Policy (ETP) and Education Sector Development Plan (ESDP) [73]. These were established to focuses on access, quality improvement, capacity development and direct funding to secondary schools. The combined effort was expected to improve the overall status of secondary education, but still the problem is far from over.

---

[1]http://www.twaweza.org/go/uwezo-datasets

In addition to that, gender plays a role on addressing this problem of student dropout. Across the world, females are more likely than males to be out of school, and the poorest girls from the most disadvantaged rural areas tend to have the lowest educational attainment levels. The prevalence of unequal distribution of education in male and female students hinders the development at every stage of a nation. Though, insignificant attention has been dedicated to examining the effects of girl child dropout in schools especially in the developing countries where the problem is widespread, literature has found that girlsâĂŹ dropout rate is significantly higher in rural schools compared to urban schools [68]. Yet, the issue of girl child dropout is a serious problem that dramatically impact on national development.

Furthermore, [68] observes that though the enrollment in school is almost same for girls and boys, boys have a higher likelihood of continuing school compared to girls. Moreover, girls overall attain less education and tend to drop out earlier as compared to boys. Thus, when dropout rate varies by gender and if girls tend to drop out earlier compared to boys, it manifests that there are some unique factors contributing to the increase in the dropout rate, particularly for girls.

The reasons why females are more likely than males to be out of school relate to social power structures and socially-constructed norms that define the roles that boys and girls should play. These gender roles affect the rights, responsibilities, opportunities and capabilities of males and females, including their access to and treatment in school. Mainly because of gendered perceptions of adolescent girl's roles and responsibilities, in most developing countries, girl's enrollment rates fall when they reach lower secondary school age and then decline further when they reach upper secondary school age [72].

In many developing countries including Tanzania, more than half of the school dropouts are largely attributed to healthy challenges during adolescence. [76], estimated that about 10 percent of school-age African girls do not attend school during menstruation, or drop out at puberty because of the lack of clean and private sanitation facilities in schools. While the government has to ensure that secondary education remains to be free and compulsory [30], it should also include the commitment to guarantee that intervention programs are informed by the collected statistics so as to contributes in reducing the dropout rates.

### 1.2. Machine learning in education

Over the past two decades, there has been significant advances in the field of machine learning. This field emerged as the method of choice for developing practical software for computer vision, speech recognition, natural language processing, robot control, and other applications [34]. There are several areas where machine learning can positively impact education. The study conducted by [15], reported on the growth of the use of machine learning in education, this is due to the rise in the amount of education data available through the digitization. Various schools have started to create personalized learning experiences through the use of technology in classrooms. Furthermore, Massive open on-line courses (MOOCs) have attracted millions of learners and present an opportunity to apply and develop machine learning methods towards improving student learning outcomes, leveraging the data collected [44].

Owing to the advancement of the amount of data collected, machine learning techniques have been applied to improve educational quality including areas related to learning and content analytics [43, 80], knowledge tracing [83], learning material enhancement [3] and early warning systems [11, 14, 77]. The use of these techniques for educational purpose is an promising field aimed at developing methods of exploring data from computational educational settings and discovering meaningful patterns [58].

One of the first applications of machine learning in education has been to help quizzes and tests move from multiple choice to fill in the blank answers [2]. The evaluation of students free form answers was based on Natural Language Processing (NLP) and machine learning. Various studies on efficacy of automated scoring show better results than human graders in some cases. Furthermore, automated scoring provides more immediate scoring than a human, which helps for use in formative assessment.

A few years ago, prediction has been observed as an application of machine learning in education [3]. A research conducted by [38], presented a novel case study describing the emerging field of educational machine learning. In this study, students key demographic characteristic data and grading data were explored as the data set for a machine learning regression method that was used to predict a studentâĂŹs future performance. In a similar vein, several projects were conducted including a project that aims to develop a prediction model that can be used by educators, schools, and policy makers to predict the risk of a student to drop out of school [4]. Springboarding from these examples, IBM Chalapathy Neti shared their vision of Smart Classrooms using cloud-based learning systems that can help teachers identify students who are most at risk of dropping out, and observe why they are struggling, as well as provide insight into the interventions needed to overcome their learning challenges[5].

Certainly, machine learning application in education still face several challenges which need to be addressed. There is lack of available open-access datasets especially in developing countries; more data-sets need to be developed, however cost must be acquired. Apart from that, several researchers ignore the fact that evaluation procedures and metrics should be relevant to school administrators. According to [42], the evaluation process should be designed to cater the needs of educators rather than only focused on common used machine learning metrics. In addition to that; the same study reveal that, many studies focused only on providing early prediction. While, a more robust and comprehensive early warning systems should be capable of identifying students at risk in future cohorts, rank students according to their probability of dropping and identifying students who are at risk even before they drop. Therefore, there is need to focus on facilitating a more robust and comprehensive early warning systems for students dropout. Also, there is need to focus on school level datasets rather than only focusing on student level datasets; this is due to the fact that school districts often have limited resources for assisting students and the availability of these resources varies with time. Therefore, identifying at risk schools will help the authorities to plan for resource allocation before the risk.

The power of machine learning can step in building better data to help authorities draw out crucial insights that change outcomes. When students drop out of school instead of continuing their education, both students and communities lose out on skills, talent and innovation [6]. On addressing student dropout problem, several predictive models were developed to process complex data sets that include details about enrollment, student performance, gender and socio-economic demographics, school infrastructure and teacher skills to find predictive patterns. Despite the fact that, evaluation of developed predictive models tend to differ but the focus remain on supporting administrators and educators to intervene and target the most at-risk students so as to invest and prevent dropouts in order to keep young people learning.

---

## 2. Approach

During last few years several works have been done on machine learning in education such as student dropout prediction, student academic performance prediction, student final result prediction etc. The findings of these studies is very useful on understanding the problem and improving measures to address solution.

In this paper we searched the following databases: ResearchGate, Elsevier , Science Direct, Springer Link, IEEE Xplore, and other computer science journals. In searching sentences and keywords we used: Predicting student Dropout, Predicting student dropout using machine learning techniques, Application of machine learning in education and student dropout prediction using machine learning techniques.

Publication periods taken into consideration is 2013 to 2017. On types of text searched we use PDF, Documents and Full length paper with abstract and keywords. Furthermore, in search items we used journal articles, conferences paper, workshop papers, topics related blogs, expert lectures or talks and other topic related communities such as educational machine learning community.

Besides, preliminary survey was conducted to stakeholders and identified that most of the collected data are not in the direct form to support machine learning approach. Most of the data collected are not clean and contains missing values, this is unavoidable problem in dealing with most of the real world data sources. Furthermore, on addressing the problem of student dropout; most of the datasets are facing imbalance problem which needs special attention on developing predictive algorithm. Therefore, this paper provides a suggested approach on addressing the problem of student dropout with focus on three important aspects of data preprocessing, machine learning techniques and evaluation measures to be considered.

## 3. Material and methods

### 3.1. Preprocessing of data

Data preprocessing includes data cleaning, normalization, transformation, feature extraction and selection, etc, and the product of data preprocessing is the final training set. In selection, relevant target data is selected from retained data (typically very noisy) and subsequently preprocessed. This goes hand in hand with the integration from multiple sources, filtering irrelevant content and structuring of data according to a target tool [36]. On developing a generalized algorithm, data preprocessing can often have a significant impact. Based on the nature of datasets in many domains, it is well known that data preparation and filtering steps take considerable amount of processing time in ML problems.

Various approaches have been identified in handling missing values, outliers data and numeric values [70]. On addressing the problem of student dropout, one of the common problem which must be considered during preprocessing is data imbalance [74]. Several re-sampling techniques such as under-sampling, over-sampling and hybrids methods can be applied [50].

Under-sampling is a non-heuristic method that aim to create a subset of the original dataset by eliminating instances until the remaining number of examples is roughly the same as that of the minority class. Over-sampling method create a superset of the original dataset by replicating some instances or creating new instances from existing ones until the number of selected examples plus the original examples of the minority class is roughly equal to that of the majority class. While, hybrids method such as SMOTE (Synthetic Minority Oversampling Technique) combines both under-sampling and over-sampling approaches [74].

## *3.2. Machine learning techniques on addressing student dropout*

In the context of education on addressing student dropout prediction, the techniques for learning can be supervised or unsupervised.

Supervised learning based on learning from a set of labeled examples in the training set so that it can identify unlabeled examples in the test set with the highest possible accuracy [23]. The paradigm of this learning is efficient and it always finds solutions to several linear and non-linear problems such as classification, plant control, forecasting, prediction, robotics and so many others [67].

Several existing works have focused on supervised learning algorithms such as Naive Bayesian Algorithm, Association rules mining, ANN based algorithm, Logistic Regression, CART, C4.5, J48, (BayesNet), SimpleLogistic, JRip, RandomForest, Logistic regression analysis, ICRM2 for the classification of the educational dropout student [41]. However, under the classification techniques, Neural Network and Decision Tree are the two methods highly used by the researchers for predicting students performance [35, 69]. The advantage of neural network is that it has the ability to detect all possible interactions between predictors variables [27] and could also do a complete detection without having any doubt even in complex nonlinear relationship between dependent and independent variables [6], while decision tree have been used because of its simplicity and comprehensibility to uncover small or large data structure and predict the value [57].

Unlike supervised, unsupervised learning algorithm is used to identify hidden patterns in unlabeled input data. It refers to provide ability to learn and organise information without an error signal and be able to evaluate the potential solution. The lack of direction for the learning algorithm in unsupervised learning can sometime be advantageous, since it lets the algorithm to look back for patterns that have not been previously considered [67].

Several techniques have been proposed on addressing this problem of student dropout using different approaches such as Survival Analysis [4, 5], Matrix Factorization [9, 22, 29, 32, 33], and Deep Neural Network [24, 79]. Other approaches such as time series clustering [31, 54] were presented to perform clustering, which are extensively used in recommender systems [81].

On addressing the problem of student dropout, machine learning techniques have been applied in various platforms such as Massive Open On-line Course (MOOC) [17, 24, 47, 62] and other Learning Management System (LMS) such as Moodle [22, 31, 66]. These platforms generated datasets which contain information that can be categorized into academic performance, socio-economic and personal information [45]. MOOC platforms such as Coursera and edX is among popular used platforms for student dropout prediction [17]. While, Moodle as a popular Learning Management System [66], provides public datasets such UMN LMS [22]. Furthermore, on identifying at risk students for early interventions, other researchers collected data from an on-line graduate program in the United States and validation was conducted by using Fall 2014 data set [31].

### *3.2.1. Survival Analysis*

Survival analysis is used to analyze data in which the time until the event is of interest [37]. It provides various mechanisms to handle such censored data problems that arise in modeling such longitudinal data (also referred to as time-to-event data when modeling a particular event of interest is the main objective of the problem) which occurs ubiquitously in various real-world application domains [78].

In the context of education, the use of survival analysis modeling to study student retention was developed. [5] developed a survival analysis framework for early prediction using Cox proportional hazards model (Cox) and applied time-dependent Cox (TD-Cox), which captures time-varying factors and can leverage those information to provide more accurate prediction of student dropout. Certainly, in survival

analysis subjects are usually followed over a specified time period and the focus is on the time at which the event of interest occurs [46]. Thus, the benefit of using survival analysis over other methods is the ability to add the time component into the model and also effectively handle censored data. In spite of the success of survival analysis methods in other domains such as health care, engineering, etc., there is only a limited attempt of using these methods in student retention problem [10].

### 3.2.2. Matrix Factorization

Matrix factorization is a clustering machine learning methods that can accommodate framework with some variations [82]. The study presented by [22, 29], described matrix factorization. In [22] study, two classes of methods for building the prediction models were presented. The first class builds these models by using linear regression approaches and the second class builds these models by using matrix factorization approaches. Regression-based methods describe course-specific regression (CSpR) and personalized linear multi-regression (PLMR) while matrix factorization based methods associate standard Matrix Factorization (MF) approach. One limitation of the standard MF method is that it ignores the sequence in which the students have taken the various courses and as such the latent representation of a course can potentially be influenced by the performance of the students in courses that were taken afterward.

Furthermore, the work present in [32] study, proposed a new data transformation model, which is built upon the summarized data matrix of link-based cluster ensembles (LCE). Like several existing dimension reduction techniques such as Principal Component Analysis (PCA) and Kernel Principal Component Analysis (KPCA), this method aims to achieve high classification accuracy by transforming the original data to a new form. However, the common limitation of these new techniques is the demanding time complexity, such that it may not scale up well to a very large dataset. Whilst worst-case traversal time (WCT-T) is not quite for a highly time-critical application, it can be an attractive candidate for those quality-led works, such as the identification of those students at risk of under achievement.

### 3.2.3. Deep Neural Network and Probabilistic Graphical Model

Deep neural network (DNN) is an approach based on Artificial Neural Networks (ANN) with multiple hidden layers between the input and output layers [20]. While, Probabilistic Graphical Model (PGM) combine probability theory and graph theory so as to offer a compact graph-based representation of joint probability distributions exploiting conditional independences among the random variables [60]. Similar to shallow ANNs, DNNs can model complex non-linear relationships [56, 63]. Recently, different deep learning architecture such as Recurrent Neural Network (RNN) and other probabilistic graphical model such as Hidden Markov Model (HMM) have been employed to the problem of student dropout.

The study presented by [24] considered two temporal models which are state space models and recurrent neural networks. State space models describe two variants of Input Output Hidden Markov Model (IOHMM) with continuous state space while recurrent neural networks describe vanilla RNN and RNN with Long Short Term Memory (LSTM) cells as hidden units. IOHMM was proposed by for learning problems involving sequentially structured data. While it is originated from HMM, it is more general that it can learn to map input sequences to output sequences. Moreover, unlike the standard discrete-state HMM, the state space in described IOHMM formulation is continuous so that the state space can in principle bear more representation power compared with enumerating discrete states. Furthermore, Vanilla Recurrent Neural Network (Vanilla RNN), unlike feed forward neural networks such as the Multi Layer Perceptron (MLP), allows the network connections to form cycles.

The limitation of that conducted study was vanishing gradient problem. While an important property of RNNs is their ability to use contextual information in learning the mapping between the input and output sequences, a subtlety is that, for basic RNN models, the range of temporality that can be accessed

in practice is usually quite limited so that the dynamic states of RNNs are considered as short term memory. This is because the influence of a given input on the hidden layer, and therefore on the network output, either decays or blows up exponentially as it cycles around the recurrent connections. To handle short-term memory of RNNs last for longer so as to tackle the vanishing gradient problem Long Short-Term Memory RNN (LSTM Network) was introduced.

*3.3. Evaluation metrics for student dropout prediction*

Many researchers use various evaluation metrics to measure the performance of student dropout algorithms. On measuring percentage of subjects that are classified correctly, several researchers use accuracy metric [5, 42]. Accuracy is a statistical measure for quantifying the degree of correctness with which a prediction model is able to label the data points [42]. Though accuracy is a very widely used metric and is very useful in practice, it is also a very conservative metric in this context. Further, the metric does not distinguish between the magnitude of errors and it might not be appropriate when the data is imbalance [48, 50].

Classification of imbalanced class size data is where one class is under-represented relative to another [1, 13, 25, 39, 40, 48, 50, 51, 71]. According to [26], the imbalanced ratio is about at least 1:10. Since the minority class usually represents the most important concept to be learned, it is very difficult to identify it due to exceptional and significant cases [50].

In the context of education, data imbalance is very common classification problem in the field of student retention, mainly because there is large number of students who are registered but there are few number of dropout students [74]. Since accuracy has less effect on minority class than majority class [49], several researchers applied other metrics such as F-measure, Mean Absolute Error (MAE) and Area Under the curve (AUC) on addressing this problem of student dropout.

F-measure is defined as a harmonic mean of precision and recall [5]. Several researchers [7, 53, 64] used this metric to evaluate algorithms when predicting student dropout. A high value of F-measure indicates that both precision and recall are reasonably high as defined in equation 1.

$$F_m = \frac{2 \cdot Precison \cdot Recall}{Precision + Recall} \tag{1}$$

where $Precision = \frac{TP}{TP+FP}$ and $Recall = \frac{TP}{TP+FN}$.

From the formula above; $TP$ stand for True Positive, $FP$ for False Positive and $FN$ for False Negative.

The study conducted by [42], presented precision at top K and recall at top K as metrics which are far more informative to educators than traditional precision recall curves. The metrics were used so as to provide precision and recall values of various algorithms at different values of K. Furthermore, the metrics were more informative and help educators to infer the precision and recall of various algorithms at a threshold K of their choice.

Other researchers [5, 22, 42, 64] applied frequently used metric in regression problem such as Mean Absolute Error (MAE) on addressing student dropout, with consideration of time to dropout prediction. MAE is a quantity used to measure how close the predictions are to the actual outcomes as stated in

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i| \tag{2}$$

where $\hat{y}_i$ is the predicted value and $y_i$ is the true value for subject *i*.

The limitation of this metric is based on how it treats both underestimating and overestimating of actual value in the same manner. However, in student retention problem, these types of errors have different meaning.

Furthermore, several studies observed AUC on measuring the performance of algorithms used on addressing student dropout problem. AUC is expressed as area under the receiver operating characteristic (ROC) where the curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) under various threshold values [8, 24, 28, 47, 53, 62]. In order to plot a ROC curve, we need to vary the discrimination or classification threshold to generate the corresponding TPR and FPR, so that the AUC measure is invariant to the classification threshold. Numerically speaking, an AUC score is a number in the range [0,1], and the closer the number is to 1, the better the classification performance [24].

Other metrics used to evaluate the performance of models are mean squared error [32, 81], Root-Mean-Square Error (RMSE) [22], error residuals [61], and misclassification rates [31].

## 4. Experimental framework

### 4.1. Dataset description and procedures

In this paper, uwezo data on learning [7] at the country level in Tanzania which was collected in 2015 was used. The dataset consists of 18 features and approximately 61340 samples, which were collected with aim on assessing childrenâĂŹs learning levels across hundreds of thousands of households. The dataset were cleaned by replacing the missing values with medians and zeros. Since our target variable is dropout, we checked the distribution of this variable in the dataset and observed that there was imbalance for target variable with only 1.6% dropout as shown in Fig. 1.
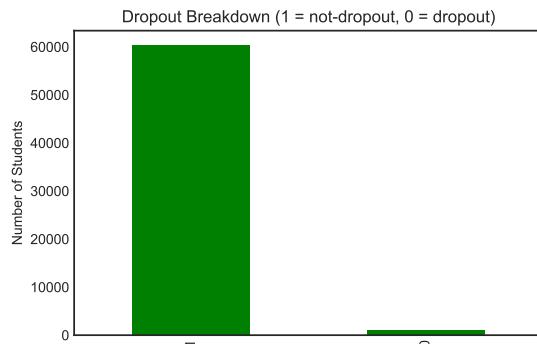


Fig. 1. Dropout distribution training data

Several approaches such as data re-sampling and generating synthetic samples, just to mention a few, can be used to address this problem. For this problem we opt to use a variety of SMOTE [50]; a popular

---

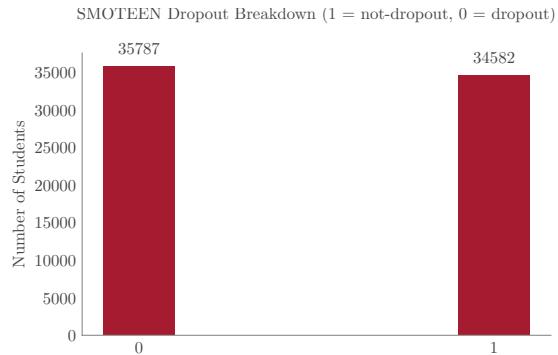[7]http://www.twaweza.org/go/uwezo-datasets
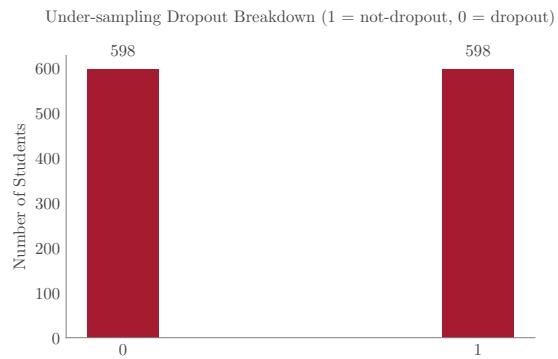
Fig. 2. Dropout-distribution (SMOTEEN)



Fig. 3. Dropout-distribution (Under-sampling)

technique for generating synthetic samples from minority class in order to reinforce its signal. Specifically, SMOTEEN[8] and RandomUnderSampler technique as implemented in Imbalanced-Learn[9] were used. SMOTEEN combine over- and under-sampling using SMOTE and Edited Nearest Neighbour (EN) to generate more minority class where RandomUnderSampler is a fast and easy way to balance the minority class by randomly selecting a subset of data for the targeted classes as observed in Fig. 2, we also observed dropout distribution (Under-sampling) as shown in Fig. 3. The dataset was separated into training (60%), test (20%) and validation (20%).

### 4.2. Experimental procedures

The experiment procedure is summarized in Fig. 4. In each experiment stratified $k$-fold cross validation was used. We use $k = 5$ fold out-of-bag overall cross validation instead of averaging over folds. The entire process was repeated 5 times and then averaged to get the training and validation results.

---

[8]combine over- and under-sampling using SMOTE and Edited Nearest Neighbour (EN)

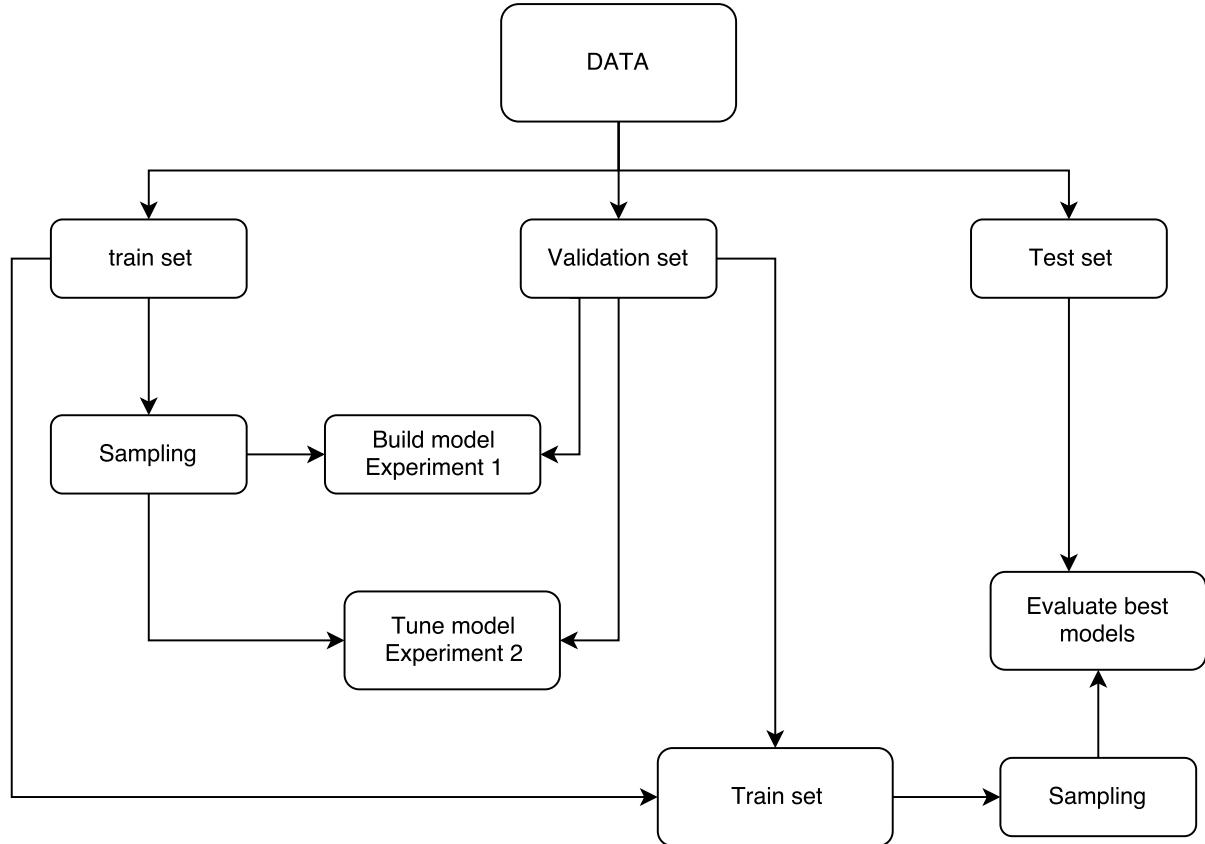[9]A Python library containing various algorithms to handle imbalanced data sets as well as producing imbalanced data sets:http://contrib.scikit-learn.org/imbalanced-learn/stable/index.html

Fig. 4. Experiment procedure

## 4.3. Evaluation metrics

To evaluate the model, Geometric Mean ($G_m$), F-measure ($F_m$) and Adjusted Geometric Mean ($AG_m$) metrics were used. The choice of these metrics is attributed by the fact that in imbalanced domains, the evaluation of the classifiersâĂŹ performance must be carried out using specific metrics in order to take into account the class distribution [50].

Therefore, the $G_m$ is a measure of the ability of a classifier to balance TPrate (sensitivity ) and TNrate (specificity) [48] as defined in Fig. 3. This measure is maximum when TPrate and TNrate are equal. Furthermore, in order to ensure TPrate to the changes in the positive predictive value (precision) than in True Positive rate (TPrate), $F_m$ is used as defined in equation 1 above. Besides, $AGm$ as defined in equation 4 was used to obtain the highest TPrate without decreasing too much the (TNrate) [50].

$$G_m = \sqrt{(\text{TPrate} \cdot \text{TNrate})} \tag{3}$$

$$AG_m = \begin{cases} \frac{\text{GM} + \text{TNrate} \cdot (FP + TN)}{1 + FP + TN} & \text{if TPrate} > 0, \\ 0 & \text{if TPrate} = 0 \end{cases} \tag{4}$$

where:

- TN is true negative, TP is true positive, FN is false negative and FP is false positive.
- TPrate = $\frac{TP}{TP+FN}$ the percentage of positive instances correctly classified.
- TNrate = $\frac{TN}{FP+TN}$ the percentage of negative instances correctly classified.

## 4.4. Experiment 1: Model selection

In this phase 13 classifiers are considered: k-nearest-neighbors (KNN), Gaussian Naive Bayes (GNB), Logistic Regression classifier (LR), Linear Discriminant Analysis (LDA), Decision Tree (DTree), Random Forest (RForest), Adaptive Boosting (AdaBoost), Multilayer perceptron (MLP), SGD Classifier which is regularized linear models with stochastic gradient descent (SGD) learning, Extra Tree classifier (EXT), Gradient Boosting Classifier (GBC), Bernoulli Naive Bayes (BNB) and Quadratic Discriminant Analysis (QDA). The aim of this experiment is to identify the classifier with best performance for this problem. The experiment is repeated for three different cases: when no sampling is used, when under-sampling used and when over sampling (SMOTEEN) is used. The experimental results for both cases are presented.
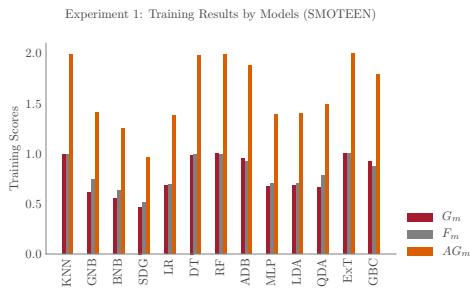


Fig. 5. Experiment 1: Training results (over-sampling)


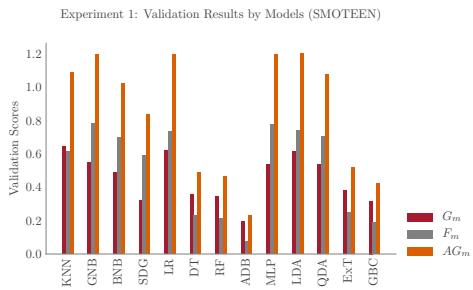
Fig. 6. Experiment 1: Validation results (oversampling)

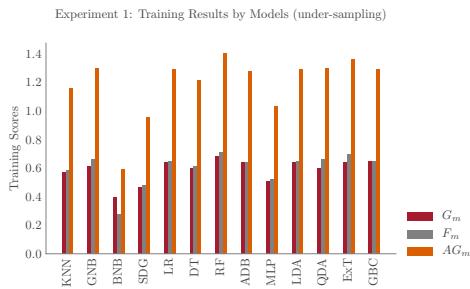Fig. 7. Experment 1: Over sampling



Fig. 8. Experiment 1: Training results (under-sampling)
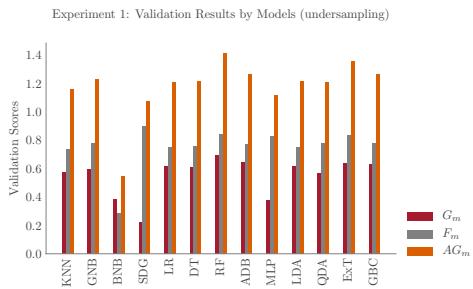


Fig. 9. Experiment 1: Validation results (under-sampling)

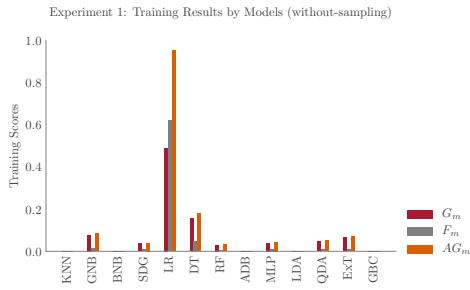Fig. 10. Experiment 1: Under sampling

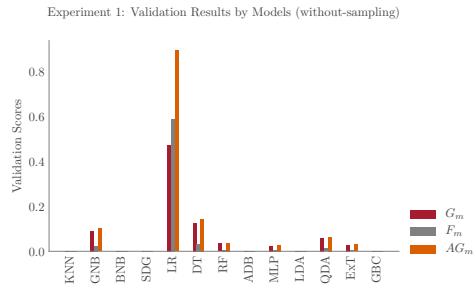Fig. 11. Experiment 1: Training results (no-sampling)



Fig. 12. Experiment 1: Validation results (no-sampling)

Fig. 13. Experiment 1: without sampling

To select the best classifiers, we only consider validation results because it give an estimate on how the classifier will perform on actual dataset which is imbalance. From the result presented in Fig. 7 three classifiers: LDA, LR and MLP show better generalization results. They show better validation result for the three metrics used. Considering the case when under-sampling is used, Fig. 10 all classifiers have considerably the same generalization results for both metrics with exception to BNB, SGD, and MLP which show lower $G_m$. The experiment conducted without sampling revel that, only LR classifier show better performance than others. However, the score rates is less than 1 for $AG_m$ as compared to when LR is used with oversampling case. Therefore, for the next experiment we only consider the following three classifiers; LR, LDA and MLP with oversampling case.

## 4.5. Experiment 2: Hyper-parameter optimization

Most ML algorithms contain several hyper parameters that can affect performance significantly (for example, the number hidden layers in MLP classifier). This experiment aimed at tuning the three selected classifiers; LR, LDA and MLP to further improve their performance. We employed hyper-parameter tuning via cross-validation and identified the best parameters for each classifier as presented in Fig. 1. We

Table 1

Model parameter

| Classfier | Parameter |
|---|---|
| LR | fit_intercept:True, tol:1, C:0.001, Penalty:'l1' |
| LDA | shrinkager:'auto, tol:1e-06, solver:'lsqr' |
| MLP | solver:'adam', learning_rate_int:0.001, shuffle:True, hidden_layer_size:10, alpha:1, early_stoping: True |

then evaluated the models by comparing their validation results as shown in Fig. 2. The experimental results allow us to measure the extent to which hyper parameter tuning improves each algorithm's performance compared to its baseline settings.

To further improve the performance of the models we employed ensemble technique. Ensemble method is one of the popular approach for improving machine learning algorithms. This approach create multiple models and then combine them to produce improved results. Several ensemble techniques such as bagging, boosting and voting have been extensively use in the literature [19]. For this problem, voting ensemble technique was appropriate. We employed voting (stacking) by soft combined the three tuned classifiers LR2, LDA2 and MLP2. The tuned classifiers where then trained on the new training set obtained by combining validation and training set used in previous experiment. To evaluate the generalization performance, the models were tested on unseen tested data. The result for this experiment is presented in Fig. 2.

Table 2

Experiment 2: Results

|  |  | LR | LR2 | MLP | MLP2 | LDA | LDA2 | ENB |
|---|---|---|---|---|---|---|---|---|
| Validation Scores | $G_m$ | 0.616 | 0.617 | 0.568 | 0.606 | 0.614 | 0.613 | **0.623** |
|  | $AG_m$ | 1.191 | **1.265** | 1.110 | 1.225 | 1.199 | 1.198 | 1.262 |
|  | $F_m$ | 0.732 | **0.787** | 0.683 | 0.766 | 0.741 | 0.740 | 0.781 |
| Test Scores | $G_m$ | 0.610 | 0.612 | 0.513 | 0.614 | 0.612 | 0.612 | **0.635** |
|  | $AG_m$ | 1.183 | 1.260 | 1.336 | 1.167 | 1.200 | 1.200 | **1.277** |
|  | $F_m$ | 0.731 | **0.788** | 0.664 | 0.714 | 0.744 | 0.744 | 0.784 |

From 2, it can be seen that the stacking classifier (ENB) show considerably better validation and test results followed by the tuned logistic regression model (LR2).

*4.6. Experiment 3: Feature Importance*

The experiment aimed at identifying the contribution of each features on the prediction performance by automatically selecting features that are most relevant to the dropout predictive modeling. Our dataset consists of 18 features which are: Main source of household income (Income), Boy's Pupil Latrines Ratio (BPLR), Pupil Teacher Ratio (PTR), School has girl's privacy room (SGR), Pupil Classroom Ratio (PCR), Region, District, Parent Teacher Meeting Ratio (PTMR), Girl's Pupil Latrines Ratio (GPLR), Household size (HHsize), Enumeration Area type (EAarea), Village, Student age (Age), Parent who discuss his/her child's progress with teacher last term (PTD), Student who did read any book with his/her parent in last week (SPB), Household meals per day (MLPD), Parent who check his /her childâĂŹs exercise book once in a week (PCCB) and Student gender (Sex). This was accomplished by directly measuring the impact of each feature on the model performance ($G_m$) obtained by permuting the values of each feature and measure how much the permutation decreases the model performance. Thus for unimportant features, the permutation will have little or no effect on model accuracy, while permuting important variables should significantly decrease it. It clearly that student sex have strong contribution on the dropout prediction performance. The results are presented in Fig. 14.
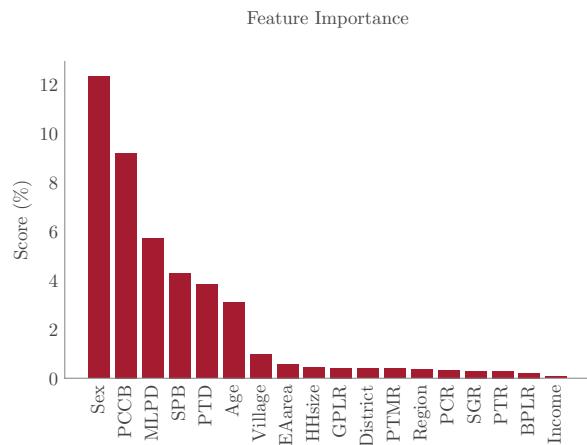
Fig. 14. Feature selection

## 5. Discussion and conclusions

### 5.1. Open challenge and future research direction

In the previous sections we have presented an up to date overview of machine learning techniques on addressing student dropout problem and highlighting the gaps and limitations. Despite several efforts done by previous researchers, there are still some challenges which need to be addressed.

It has been observed that, most of the algorithms have been developed and tested in developed countries using existing datasets generated from developed countries. Furthermore, MOOC and Moodle are among the most used platforms which offer public datasets to be used on addressing the problem of student dropout. The limitation of public datasets from developing countries [53], brought need to develop more datasets from different geographical location. However, cost and time must be acquired to accommodate data collection process. Furthermore, to the knowledge of researchers, there are only few research which has been conducted in developing countries. Thus, further research is needed to explore the value of machine learning algorithms in cubing dropout in the context of developing countries.

Second, most of the presented works have focused on providing early prediction only [42]. Therefore, future work should focus on facilitating a more robust and comprehensive early warning systems for students dropout which can identify students at risk in future cohorts, rank students according to their probability of dropping and identifying students who are at risk even before they drop.

Third, most existing study ignore the fact that dropout rate is often low in existing datasets. This is a serious problem especially in the context of student retention [74], with dropout students significantly less than those who stay and thus future research should consider developing a student dropout algorithm with consideration of data imbalance problem.

Fourth, many studies focus on addressing student dropout using student level datasets. However, developing countries need to include school level datasets on addressing the problem due to the issue of limited resources which face many school districts [42]. This will involve the use of new sources school level data and applying additional machine learning approaches to improve predictive power of the proposed algorithm. The algorithm will enable relevant authorities to effectively and accurately plan, formulate policies, and make decisions on measures to address the problem.

*5.2. Conclusions*

In this work, a review of machine learning techniques on addressing student dropout problem is presented. The review draws several conclusions;

First, while several techniques have been proposed for addressing student dropout in developed countries, there is lack of research on the use of machine learning on addressing this problem in developing countries.

Second, despite the major efforts on using machine learning in education, data imbalance problem has been ignored by many researchers. This facilitate using improper evaluation metrics on analyzing performance of the algorithms.

Third, many research focus on providing early prediction rather than including ranking and forecasting mechanisms on addressing the problem of student dropout.

Fourth, school level datasets must be considered when addressing this problem, in order to come up with the proposed solutions to facilitate the authorities on identifying at risk schools for early intervention.

Fifth, we have empirically assessed 13 supervised classification algorithms on a set of approximately 61340 supervised classification dataset in order to provide a contemporary set of recommendations to researchers who wish to apply machine learning algorithms to their data with consideration of data imbalanced problem. The three classifiers LR, LDA and MLP have proven superior to all the other classifiers by achieving highest performance metrics when over-sampling technique is employed. Furthermore, we show that hyper parameter tuning improves each algorithmâĂŹs performance compared to its baseline settings and stacking these classifiers improve the overall predictive performance. We also show the contribution of each features on prediction performance with student sex being the leading feature.

**Acknowledgments**

**References**

[1] Abdi, L. and Hashemi, S. (2014). An Ensemble Pruning Approach Based on Reinforcement Learning in Presence of Multi-class Imbalanced Data.

[2] Adhatrao, K., Gaykar, A., Dhawan, A., Jha, R., and Honrao, V. (2013). Predicting Students' Performance Using Id3 and C4.5 Classification Algorithms. *International Journal of Data Mining and Knowledge Management Process*, 3(5):39–52.

[3] Agrawal, R. (2014). Mining Videos from the Web for Electronic Textbooks.

[4] Ameri, S. (2015). Survival Analysis Approach For Early Prediction Of Student Dropout.

[5] Ameri, S., Fard, M. J., Chinnam, R. B., and Reddy, C. K. (2016). Survival Analysis Based Framework for Early Prediction of Student Dropouts. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, CIKM '16, pages 903–912, New York, NY, USA. ACM.

[6] Arsad, Pauziah Mohd Buniyamin, Norlida Manan, J.-l. A. (2013). A neural network students' performance prediction model (NNSPPM). *Smart Instrumentation, Measurement and Applications (ICSIMA), 2013 IEEE International Conference on*, (November):1–5.

[7] Aulck, L., Aras, R., Li, L., Heureux, C. L., Lu, P., and West, J. (2017). STEM-ming the Tide : Predicting STEM attrition using student transcript data. (August).

[8] Aulck, L., Velagapudi, N., Blumenstock, J., and West, J. (2016). Predicting Student Dropout in Higher Education.

[9] Babu, A. R. (2015). Comparative Analysis of Cascadeded Multilevel Inverter for Phase Disposition and Phase Shift Carrier PWM for Different Load. *Indian Journal of Science and Technology*, 8(April):251–262.

[10] Bani, M. J. and Haji, M. (2017). College Student Retention: When Do We Losing Them?

[11] Beck, H. P. and Davidson, W. D. (2016). Establishing an Early Warning System : Predicting Low Grades in College Students from Survey of Academic Orientations ... (December 2001).

[12] BEST (2015). Pre-Primary, Primary and Secondary Education Statistics in Brief 2016 The United Republic of Tanzania President's Office Regional Administration and Local Government.

[13] Borowska, K. and Topczewska, M. (2016). New Data Level Approach for Imbalanced Data Classification Improvement. pages 283–294.

[14] Brundage, A. (2014). The use of early warning systems to promote success for all students.

[15] Center for Digital Technology and Management (2015). *THE FUTURE OF EDUCATION TREND REPORT 2015*.

[16] Chen, J. F., Hsieh, H. N., and Do, Q. H. (2014). Predicting student academic performance: A comparison of two meta-heuristic algorithms inspired by cuckoo birds for training neural networks. *Algorithms*, 7(4):538–553.

[17] Chen, Y., Chen, Q., Zhao, M., Boyer, S., Veeramachaneni, K., and Qu, H. (2017). DropoutSeer: Visualizing learning patterns in Massive Open Online Courses for dropout reasoning and prediction. *2016 IEEE Conference on Visual Analytics Science and Technology, VAST 2016 - Proceedings*, pages 111–120.

[18] Choudhary AI, L. A. (2015). Economic Effects of Student Dropouts: A Comparative Study. *Journal of Global Economics*, 03(02):2–5.

[19] Dalvi, P. T. and Vernekar, N. (2016). Anemia Detection using Ensemble Learning Techniques and Statistical Models. pages 1747–1751.

[20] Deng, L. and Yu, D. (2014). Deep Learning: Methods and Applications. *Foundations and Trends® in Signal Processing*, 7(3-4):197–387.

[21] Durairaj, M. and Vijitha, C. (2014). Educational data mining for prediction of student performance using clustering algorithms. *International Journal of Computer Science and Information Technologies (IJCSIT)*, 5(4):5987–5991.

[22] Elbadrawy, A., Polyzou, A., Ren, Z., Sweeney, M., Karypis, G., and Rangwala, H. (2016). –okay–Predicting Student Performance Using Personalized Analytics. *Computer*, 49(4):61–69.

[23] Erik G. (2014). Introduction to Supervised Learning. pages 1–5.

[24] Fei, M. and Yeung, D.-Y. (2015). Temporal Models for Predicting Student Dropout in Massive Open Online Courses. *2015 IEEE International Conference on Data Mining Workshop (ICDMW)*, pages 256–263.

[25] Galar, M., Fernández, A., Barrenechea, E., Bustince, H., and Herrera, F. (2016). New Ordering-Based Pruning Metrics for Ensembles of Classifiers in Imbalanced Datasets.

[26] Gao, T. (2015). Hybrid classification approach of SMOTE and instance selection for imbalanced datasets.

[27] Gray, G., McGuinness, C., and Owende, P. (2014). An application of classification models to predict learner progression in tertiary education. *2014 4th IEEE International Advance Computing Conference, IACC 2014*, pages 549–554.

[28] Halland, R., Igel, C., and Alstrup, S. (2015). High-School Dropout Prediction Using Machine Learning : A Danish Large-scale Study. *ESANN 2015 proceedings, European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, (April):22–24.

[29] Hu, Q. and Rangwala, H. (2016). Enriching Course-Specific Regression Models with Content Features for Grade Prediction.

[30] Human RIghts Watch (2017). âĂIJI Had a Dream to Finish SchoolâĂİ:Barriers to Secondary Education in Tanzania.

[31] Hung, J. L., Wang, M. C., Wang, S., Abdelrasoul, M., Li, Y., and He, W. (2017). Identifying At-Risk Students for Early Interventions - A Time-Series Clustering Approach. *IEEE Transactions on Emerging Topics in Computing*, 5(1):45–55.

[32] Iam-On, N. and Boongoen, T. (2017). Generating descriptive model for student dropout: a review of clustering approach. *Human-centric Computing and Information Sciences*, 7(1):1.

[33] Iqbal, Z., Qadir, J., Mian, A. N., and Kamiran, F. (2017). Machine Learning Based Student Grade Prediction: A Case Study. pages 1–22.

[34] Jordan, M. I. and Mitchell, T. M. (2015). Machine learning: Trends, perspectives, and prospects. *Science*, 349(6245):255–260.

[35] Joseph, H. R. (2014). Promoting education: A state of the art machine learning framework for feedback and monitoring E-Learning impact. *2014 IEEE Global Humanitarian Technology Conference - South Asia Satellite, GHTC-SAS 2014*, pages 251–254.

[36] Kalegele, K., Sasai, K., Takahashi, H., Kitagata, G., and Kinoshita, T. (2015). Four Decades of Data Mining in Network and Systems Management. *IEEE Transactions on Knowledge and Data Engineering*, 27(10):2700–2716.

[37] Kartal, O. O. (2015). *USING SURVIVAL ANALYSIS TO INVESTIGATE THE PERSISTENCE OF STUDENTS IN AN INTRODUCTORY INFORMATION TECHNOLOGY COURSE AT METU*. PhD thesis.

[38] Kotsiantis, S. B. (2012). Use of machine learning techniques for educational proposes: A decision support system for forecasting students' grades. *Artificial Intelligence Review*, 37(4):331–344.

[39] Krawczyk, B. (2015). Combining One-vs-One Decomposition and Ensemble Learning for Multi-class. pages 27–36.

[40] Krawczyk, B. and B, G. S. (2015). Pattern Recognition and Machine Intelligence. 9124:535–544.

[41] Kumar, M., Singh, A. J., and Handa, D. (2017). Literature Survey on Educational Dropout Prediction. *I.J. Education and Management Engineering*, 2(March):8–19.

[42] Lakkaraju, H., Aguiar, E., Shan, C., Miller, D., Bhanpuri, N., Ghani, R., and Addison, K. L. (2015). A Machine Learning Framework to Identify Students at Risk of Adverse Academic Outcomes. *Kdd*, pages 1909–1918.

[43] Lan, A. S., Studer, C., and Baraniuk, R. G. (2014). Time-varying Learning and Content Analytics via Sparse Factor Analysis.

[44] Lee, K. (2017). Large-Scale and Interpretable Collaborative Filtering for Educational Data. pages 1–7.

[45] Lei, C. and Li, K. F. (2015). Academic Performance Predictors. In *Proceedings - IEEE 29th International Conference on Advanced Information Networking and Applications Workshops, WAINA 2015*.

[46] Li, Y., Wang, J., Ye, J., and Reddy, C. K. (2016). A Multi-Task Learning Formulation for Survival Analysis. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16*, pages 1715–1724.

[47] Liang, J., Li, C., and Zheng, L. (2016). Machine learning application in MOOCs: Dropout prediction. *ICCSE 2016 - 11th International Conference on Computer Science and Education*, (Iccse):52–57.

[48] Lin, W. J. and Chen, J. J. (2013). Class-imbalanced classifiers for high-dimensional data. *Briefings in Bioinformatics*, 14(1):13–26.

[49] Longadge, R., Dongre, S. S., and Malik, L. (2013). Class imbalance problem in data mining: review. *International Journal of Computer Science and Network*, 2(1):83–87.

[50] López, V., Fernández, A., García, S., Palade, V., and Herrera, F. (2013). An insight into classification with imbalanced data : Empirical results and current trends on using data intrinsic characteristics. *Information Sciences*, 250:113–141.

[51] Mazumder, R. U., Begum, S. A., and Biswas, D. (2015). Rough Fuzzy Classi fi cation for Class Imbalanced Data.

[52] Mgala, M. (2016). *Investigating Prediction Modelling of Academic Performance for Students in Rural Schools in Kenya*. PhD thesis, University of Cape Town.

[53] Mgala, M. and Mbogho, A. (2015). Data-driven Intervention-level Prediction Modeling for Academic Performance. *Proceedings of the Seventh International Conference on Information and Communication Technologies and Development*, pages 2:1—–2:8.

[54] Młynarska, E., Greene, D., and Cunningham, P. (2016). Time series clustering of Moodle activity data. *CEUR Workshop Proceedings*, 1751:104–115.

[55] Mosha, D. (2014). *Assessment of Factors behind Dropout in Secondary Schools in Tanzania. A Case of Meru District in Tanzania*. PhD thesis, Open University of Tanzania.

[56] Mun, S., Shin, M., Shon, S., Kim, W., Han, D., and Ko, H. (2017). DNN transfer learning based non-linear feature extraction for acoustic event classification. *IEICE Transactions on Information and Systems*, E100D(9):1–4.

[57] Natek, S. and Zwilling, M. (2014). Expert Systems with Applications Student data mining solution âĂŞ knowledge management system related to higher education institutions. *Expert Systems with Applications*, 41:6400–6407.

[58] Nunn, S., Avella, J. T., Kanai, T., and Kebritchi, M. (2016). Learning Analytics Methods, Benefits, and Challenges in Higher Education: A Systematic Literature Review. *Online Learning*, 20(2):13–29.

[59] Patron, R. (2014). Early school dropouts in developing countries: An equity issue? The Uruguayan case. *University of Uruguay*, page P13.

[60] Pernkopf, F., Peharz, R., and Tschiatschek, S. (2013). *Introduction to Probabilistic Graphical Models Introduction*.

[61] Poh, N. and Smythe, I. (2015). To what extend can we predict students' performance? A case study in colleges in South Africa. *IEEE SSCI 2014 - 2014 IEEE Symposium Series on Computational Intelligence - CIDM 2014: 2014 IEEE Symposium on Computational Intelligence and Data Mining, Proceedings*, pages 416–421.

[62] Prieto, L. P., Rodríguez-Triana, M. J., Kusmin, M., and Laanpere, M. (2017). Smart school multimodal dataset and challenges. *CEUR Workshop Proceedings*, 1828:53–59.

[63] Ramachandra, V. and Way, K. (2018). Deep Learning for Causal Inference. (0).

[64] Rovira, S., Puertas, E., and Igual, L. (2017). Data-driven system to predict academic grades and dropout. *PLOS ONE*, 12(2):1–21.

[65] Sales, A., Balby, L., and Cajueiro, A. (2016). Exploiting Academic Records for Predicting Student Drop Out : a case study in Brazilian higher education. 7(2):166–180.

[66] Santana, M. A., Costa, E. B., Neto, B. F. S., Silva, I. C. L., and Rego, J. B. A. (2015). A predictive model for identifying students with dropout profiles in online courses. *CEUR Workshop Proceedings*, 1446.

[67] Sathya, R. and Abraham, A. (2013). Comparison of Supervised and Unsupervised Learning Algorithms for Pattern Classification. *International Journal of Advanced Research in Artificial Intelligence*, 2(2):34–38.

[68] Shahidul, S. M. and Karim, A. H. M. Z. (2015). Factors contributing to school dropout among the girls: a review of literature. 3(2):25–36.

[69] Shahiri, A. M., Husain, W., and Rashid, N. A. (2015). A Review on Predicting Student's Performance Using Data Mining Techniques. *Procedia Computer Science*, 72:414–422.

[70] Shahul, S., Suneel, S., Rahaman, M. A., and Swathi, &. (2016). A Study of Data Pre-Processing Techniques for Machine Learning Algorithm to Predict Software Effort Estimation. *Imperial Journal of Interdisciplinary Research*, 2(6):2454–1362.

[71] Stefanowski, J. (2016). On Properties of Undersampling Bagging.

[72] Subrahmanyam, G. (2016). Gender perspectives on causes and effects of school dropouts.

[73] TAMISEMI (2004). The United Republic of Tanzania Ministry of Education and Culture. pages 2004–2009.

[74] Thammasiri, D., Delen, D., Meesad, P., and Kasap, N. (2014). A critical assessment of imbalanced class distribution problem: The case of predicting freshmen student attrition. *Expert Systems with Applications*, 41(2):321–330.

[75] UNESCO (2011). UNESCO Global Partnership for Girls' and Women's Education- One Year On.

[76] UNICEF (2006). UNICEF Water, Sanitation and Hygiene Annual Report.

[77] US Department of Education (2016). Definition of Early Warning Systems Research on Early Warning Systems Issue Brief: Early Warning Systems. (September):1–13.

[78] Wang, P., Li, Y., and Reddy, C. K. (2017a). Machine Learning for Survival Analysis: A Survey. *ACM Comput. Surv. Article*, 1(1):38.

[79] Wang, W., Yu, H., and Miao, C. (2017b). Deep Model for Dropout Prediction in MOOCs. *Proceedings of the 2nd International Conference on Crowd Science and Engineering - ICCSE'17*, pages 26–32.

[80] Waters, A. E., Studer, C., and Baraniuk, R. G. (2014). Sparse Factor Analysis for Learning and Content Analytics. 15:1959–2008.

[81] Xu, J., Moon, K. H., and van der Schaar, M. (2017). A Machine Learning Approach for Tracking and Predicting Student Performance in Degree Programs. *IEEE Journal of Selected Topics in Signal Processing*, 11(5):742–753.

[82] Yang, D., Piergallini, M., Howley, I., and Rose, C. (2014). Forum Thread Recommendation for Massive Open Online Courses. *Proceedings of the 7th International Conference on Educational Data Mining (EDM)*, pages 257–260.

[83] Yudelson, M. V., Koedinger, K. R., and Gordon, G. J. (2013). Individualized Bayesian Knowledge Tracing Models. pages 1–10.