



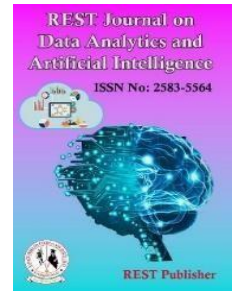
REST Journal on Data Analytics and Artificial Intelligence

Vol: 3(3), September 2024

REST Publisher; ISSN: 2583-5564

Website: <http://restpublisher.com/journals/jdaai/>

DOI: <https://doi.org/10.46632/jdaai/3/3/12>



Predictive Modeling for Early Detection of High School Dropouts Using Machine Learning Techniques

*Jayanth Kande

Southern University and A&M College, Baton Rouge, Louisiana, United States.

*Corresponding author: jayanth.m1229@gmail.com

Abstract: This research paper presents an innovative approach to developing a predictive model for early identification of high school dropouts using machine learning algorithms. The study analyzes the National Center for Education Statistics dataset to create an effective dropout detection system. To address the challenge of high dimensionality in the dataset, principal component analysis is applied to reduce its complexity. The study compares the performance of different machine learning methods, including a multi-layer artificial neural network, k -nearest neighbors, a support vector machine with a radial basis function kernel, and a support vector machine with a polynomial kernel. The objective is to determine the most accurate classifier for predicting dropout risk. The experimental results highlight the neural network as the top-performing classifier, with statistically significant differences compared to k -nearest neighbors. These findings contribute to developing proactive measures and interventions to prevent high school dropouts and enhance educational outcomes.

1. INTRODUCTION

Dropping out of school can have dire consequences for the remainder of a student's life. Compared to students who graduate high school, dropouts are much more likely to be living in poverty, unhealthy, unemployed, in prison, divorced, and have children who will also drop out from high school [3]. On average, these students end up making \$1 million less over their lifetimes than those who graduate [9]. Although high school graduation rates have improved nationally, 10 states have shown lower graduation rates over the past decade [1]. If students who are at risk of dropping out could be flagged before they drop out, and some intervention could be performed to prevent those students from dropping out, many of these hardships could be avoided.

This paper will use the 2009 dataset from the Education Longitudinal Study (ELS) by the National Center for Education Statistics (NCES), which contains almost 7000 variables for over 23,000 students [8]. ELS is one of many extensive longitudinal surveys by NCES. These surveys track students starting from either ninth or tenth grade, depending on the survey, and follow them through their adult life. The 2009 survey analyzed in this paper includes a study of ninth graders in the 2009-2010 school year and one follow-up in the spring of 2012. The sampling is representative of students in the United States. It includes students randomly selected from more than 900 schools from all 50 states and the District of Columbia. The data in the survey provides information about the student, the student's parents, socioeconomic background, teachers, and more. In this study, 8.43 percent of the students dropped out of school.

This paper aims, through the use of machine learning, to form an accurate predictive model of students who are at risk of dropping out. The data will be classified using multi-layer artificial neural networks, k -nearest neighbors (k -NN), and with support vector machines using either radial basis function kernels or polynomial kernels. The classification accuracy of the various models will be compared against one another. Principal component analysis (PCA) will be performed on the dataset to reduce its dimensionality while still being representative of the distribution. This newly transformed dataset will be used to build the various models and determine their classification accuracy.

This paper will test the hypothesis that an SVM using a radial basis function kernel will have better classification accuracy than a neural network, k -NN, or an SVM with a polynomial kernel. One argument for this hypothesis is the convexity of the SVM hinge loss function will result in optimization that will not become stuck in local minima and will result in the SVM having the highest classification accuracy. Conversely, neural networks do not have

guaranteed convexity and can, therefore, have issues with becoming stuck in local minima. In addition, the soft margin of the SVM class hyperplane can control the sensitivity to outliers. k-NN, which has a similar Euclidean distance calculation to the radial basis function kernel of the SVM, does not have an equivalent ‘slackness’ parameter. Also, it is generally accepted that RBF kernels typically outperform polynomial kernels [4]

2. LITERATURE REVIEW

There have been several studies of student performance using a variety of methods. Below is a sampling of some recent research, including statistical analysis of dropout data and the application of machine learning to student performance and attrition.

Suhyun et al. (2011) [15] investigated the change in high school dropouts between the 1980s and 2000s. In their paper, they used statistical analysis of two NCES National Longitudinal Survey of Youth (NLSY) surveys in the 1980s and the 2000s. In their analysis, they discovered the eight strongest factors contributing to student dropouts are: students with a minority race, student gender, whether students lived with a biological parent in the first year of the survey, the mother’s permissiveness, number of household members, whether the student lived in a metropolitan area, whether the student lived in the south or west of the United States, and students who were suspended from school at least once. They then went on to analyze how These factors changed between the 1980s and the 2000s.

Kotsiantes’ (2011) [10] study examined several regression techniques to attempt to predict students’ grades. Their dataset used information about a student’s background and previous grades to form a predictive model of future performance. They compared the performance of some different classifiers, including model trees, neural networks, linear regression, locally weighted regression, and support vector machines. They found that model trees performed the best in the tested regression models.

Yadav et al. (2011) [17] performed a similar study to Kotsiantes to predict students’ grades. They used decision trees to build a future performance model based on previous behavior, such as prior semester grades, attendance, test grades, and more. They found that Classification and Regression Trees (CART) provided the best classification performance with an overall accuracy of 0.5625. Chen et al. (2007) [5] analyzed performance assessment in Web-based learning environments. The goal was to be able to provide immediate, useful feedback regarding learning performance to the student. Their approach used a combination of four techniques: gray relational analysis (GRA), K-means clustering scheme, fuzzy association rule mining, and fuzzy inference. Their results showed moderate performance with a maximum average accuracy of 0.7675.

Delen (2010) [7] compared the results of multiple classifiers to predicting freshmen student attrition at Oklahoma State University. Using 29 variables that included age, high school GPA, SAT score, and fall student loan amount, they found that an SVM yielded the best overall prediction accuracy of 0.8645. Their paper compared decision trees, artificial neural networks, logistic regression, and SVMs.

Lykourantzou et al. (2009) [11] attempted to predict dropouts in e-learning courses using feed-forward neural networks, support vector machines, and probabilistic ensemble fuzzy ARTMAP. In their paper, they used time-invariant along with time-varying attributes to form their models. These variables include student demographics, prior academic performance, gender, and assignment grades. They found that by at the start of a course, they could predict dropouts with 85 percent accuracy, and by the middle of the course, they reached 97 percent accuracy. No one classifier performed the best at all stages of the course, and each performed the best in at least one stage of the course.

Pal (2012) [12] conducted a study predicting higher education student dropouts in India using naïve Bayes classifiers. They found only 17 variables, such as overall high school GPA, family income, and parents’ occupations, they were able to predict dropouts with an accuracy of 0.917.

3. EXPERIMENTS

The NCES longitudinal surveys have a wealth of information spread over several years. Unfortunately, processing more than 7000 variables and more than 23000 samples in the 2009 ELS survey using machine learning is computationally expensive. Steps will be taken to reduce the computational cost of training this large dataset and its dimensionality. First, restricted variables will be removed as this data cannot be used publically. Next, any variable with zero variance is removed as these will not contribute to the model. Also, variables with a maximum

value of less than zero are removed as these indicate that information is missing for some reason. These reasons include the student being incapable of answering a question or being unresponsive to a given question.

PCA will then be performed on the resulting dataset. PCA helps to avoid the “curse of dimensionality” by effectively compressing the dataset and keeping only the most essential information. It does this through calculating the eigenvalues of the covariance matrix [13]. The transformation results in the variables being sorted by variance. Only the variables with the most variance will be used, and the remaining variables whose variance differs by less than ten percent will be discarded. Because PCA transforms the dataset into a representation of the original data, direct comparison between the transformed and original variables is difficult.

3.1 k-Nearest Neighbor

The first classifier to be tested is the *k*-nearest neighbor (*k*-NN) classifier. *k*-NN is a lazy learning, non-parametric method that only computes classes on demand based on the *k*-nearest classes to a given sample vector [6]. Given a point x_0 , the *k* points $x_{(r)}$, $r = 1, \dots, k$ nearest to x_0 are found and x_0 is classified based off the majority vote among the *k*-nearest neighbors. The nearest neighbors can be found by calculating the Euclidian distance between points:

$$d_{(i)} = \|x_{(i)} - x_0\| \quad (1)$$

This results in zero computational cost for training; however, the expense is deferred until the classification step. This simple method has good non-linear classification accuracy; however, dimensionality, noisy inputs, and irrelevant attributes are problematic for *k*-NN. PCA or other dimensionality reduction is important when using *k*-NN classifiers with high-dimensional datasets. The *k* parameter for the algorithm will be tuned through 10-fold cross-validation.

3.2 Artificial Neural Network

Next, an artificial neural network will be tested. Neural networks attempt to model the nervous system using representations of neurons. A network with one input and output layer of neurons can perform only linear separation. Therefore, a multi-layered feed-forward neural network will be used because it can approximate non-linear functions [16]. The network that will be tested will have one input layer, one hidden layer, and one output layer. Random weights will be assigned to all units, and a gradient descent method will be used to optimize the weights. The perceptrons in the network will use the activity function $\frac{1}{1+e^{-x}}$.

$$\Delta\omega_{ij} = \eta\delta_jx_i$$

$$A_j = \sum_i \omega_{ij}x_i \text{ along with the sigmoidal activation function, } \sigma(x) = \frac{1}{1+e^{-x}}. \text{ The neural}$$

where η is the learning rate and δ_j is:

$$\delta_j = [1 - x_j]x_j \sum_k \omega_{kj} \delta_k$$

Gradient descent works well for training. However, it can have trouble with getting caught in local minima and overfitting [14]. To avoid overfitting, a weight decay parameter λ will be introduced to penalize large weights and, therefore, regularize the weight update function. The resulting weight update is the following:

$$\Delta\omega_{ij} = \eta\delta_jx_i - \eta\lambda x_i$$

The parameters of the network will be chosen through cross-validation. These parameters are the number of hidden units in the network and the weight decay parameter λ . The best combination of these parameters will be used to test the performance of the final model.

3.3 Support Vector Machine: The third and final classifier to be tested is the support vector machine (SVM). Linear SVMs create a model to find a hyperplane that separates two classes of points. The hyperplane, modeled by the simple equation $f(x) = x^T \beta + \beta_0 = 0$, attempts to maximize the margin *M* between the two classes. The margin $M = \frac{1}{\|\beta\|}$ is on each side of the hyperplane where β is a unit vector $\|\beta\| = 1$.

Therefore, classification is achieved by $G(x) = \text{sign}[x^T \beta + \beta_0]$, leading to the optimization problem

$$\begin{aligned} & \min_{\beta, \beta_0} \|\beta\| \\ & \text{subject to } u_i(x_i^T \beta + \beta_0) \geq 1, i = 1, \dots, N. \end{aligned}$$

To handle the case when the classes may overlap one another, slack variables $\xi = (\xi_1, \xi_2, \dots, \xi_N)$ may be introduced to define the amount of error allowed across the hyperplane, leading to the new optimization problem

$$\min \|\beta\| \text{ subject to } \begin{cases} u_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i \\ \xi_i \geq 0, \end{cases} \quad \xi_i \leq \text{constant}$$

This “constant” is the cost parameter C of SVM training. Altogether, this works well for linearly separable classes. However, it will not work for more complicated, non-linear class boundaries. A kernel function may be used to solve this issue. A popular kernel to use is the Gaussian radial basis function (RBF). On two inputs, x and x', the RBF kernel is:

$$K(x, x') = \exp(-\gamma \|x - x'\|^2)$$

The solution to the hyperplane can be transformed to use a kernel function using the equation

$$f(x) = \sum_{i=1}^N \alpha_i u_i K(x, x_i) + \beta_0$$

The result is a hyperplane that will curve with the class boundary. Training will be tuned using the cost parameter C using 10-fold cross-validation.

Another popular kernel is the polynomial kernel, represented by the following:

$$K(x, x') = (1 + \lambda \langle x, x' \rangle)^d$$

For the polynomial kernel, the degree parameter d, the cost parameter C, and the scale parameter λ will be tuned using 10-fold cross-validation.

3.4 Evaluation: The four classifiers mentioned above, artificial neural network (denoted NN), k-nearest neighbors (denoted KNN), and support vector machine with RBF (denoted SVMR) and polynomial kernels (denoted SVMP), will be compared in terms of overall classification accuracy and area under the receiver operator characteristic (ROC) curve [18]. The tuning of the various parameters will be presented to show how they affect overall classification accuracy. The parameters with the highest accuracy will be chosen, and this configuration will be evaluated to determine the statistical significance of the differences between classifiers.

3.5 Results

TABLE 1. Accuracy, Area Under ROC, Sensitivity, and Specificity of the models

	Accuracy	ROC	Sensitivity	Specificity
NN	0.9396	0.8303	0.9957	0.3263
KNN	0.9340	0.7312	0.9988	0.2289
SVMR	0.9391	0.7937	0.9984	0.2954
SVMP	0.9394	0.793	0.9984	0.2963

Figure 1 and corresponding Table 1 shows the accuracy, the area under the ROC curve, the sensitivity, and the specificity of the four classifiers. NN had the highest accuracy, ROC value, and highest specificity. Table 2 shows the differences in accuracy between models above the diagonal, along with the corresponding p-values with Bonferroni correction (see [2]) below the diagonal. This the table shows that while NN had the overall highest accuracy, the difference was not significant over SVMR and SVMP; however, NN did perform significantly better than KNN. SVMR was hypothesized to have the highest accuracy however, this was That is not the case. SVMP had slightly higher accuracy. However, the difference is not significant. The accuracy between the classifiers was very close.

Table 3 shows the differences in area under the ROC curves above the diagonal along with the corresponding p-values with Bonferroni correction below the diagonal. The difference between NN and the next highest ROC

value, SVMP, was significant with a p-value of 0.0046. The differences between ROC values were much more significant than the differences in overall accuracy. This corresponds with the more substantial variance in specificity. All models had very high sensitivity. However, NN had the highest specificity. Consequently, the classifiers will likely have exceptionally few false negatives; however, many false positives will be identified. That being said, it is better to falsely identify a child as a potential dropout than not to identify children who are likely to drop out.

Figure 2 shows the parameter selection for the four models. Only one hidden unit was required for the NN model, along with a weight decay λ of 0.1. Unfortunately, the NN model had issues with training more than seven hidden units, and it is possible that better performance could have been obtained from more hidden units. However, the overall trend was reduced performance with additional hidden units. KNN performed best with 10 nearest neighbors with an accuracy of 0.9340. More neighbors may have improved the performance slightly. However, it appears it is leveling off in Figure 2b. SVMR performed best for C of 0.5, resulting in an accuracy of 0.9391. Finally, SVMP saw its best performance with a cost of C of 0.25, a scale of 0.01, and a degree of 2, resulting in an accuracy of 0.9394.

TABLE 2. Above the diagonal are the differences between model accuracy and below are the p-values with Bonferroni correction

	NN	KNN	SVMR	SVMP
NN		0.0055522	0.0004271	0.0001278
KNN	0.004984		-0.0051251	-0.0054243
SVMR	1.000000	0.038251		-0.0002993
SVMP	1.000000	0.032182	1.000000	

TABLE 3. Above the diagonal are the differences between ROCs, and below are the p-values with Bonferroni correction

	NN	KNN	SVMR	SVMP
NN		0.1006253	0.0380372	0.0387978
KNN	1.518e-07		-0.0625881	-0.0618274
SVMR	7.805e-05	7.925e-06		0.0007607
SVMP	0.0045639	0.0005211	1.0000000	

4. CONCLUSION

This paper attempted to train four different models, NN, KNN, SVMR, and SVMP, on the 2009 Educational Longitudinal Study to predict if students would drop out of high school. PCA was applied to the dataset because of its high dimensionality. The classifier with the highest accuracy, NN, did not perform significantly better than SVMR or SVMP. However, NN did have a significantly higher area under the ROC curve than all other classifiers. The main contribution of this paper is that machine learning can detect high school dropouts before they occur with relatively high accuracy and with few false negatives; however, there will be some false positives.

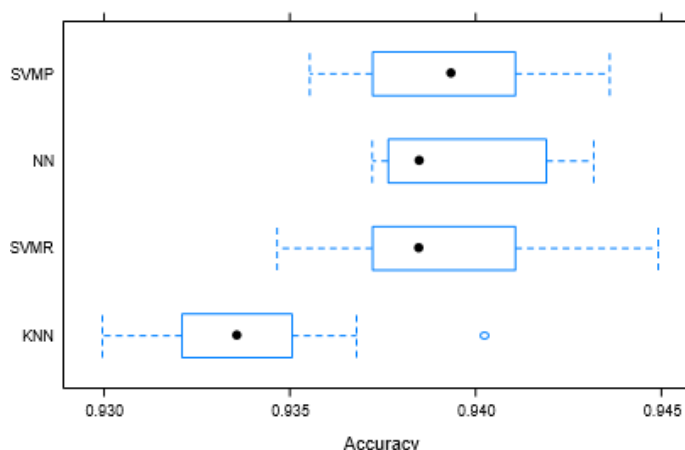


FIGURE 1. (a) Accuracy, ROC, Sensitivity, and Specificity Box Plots

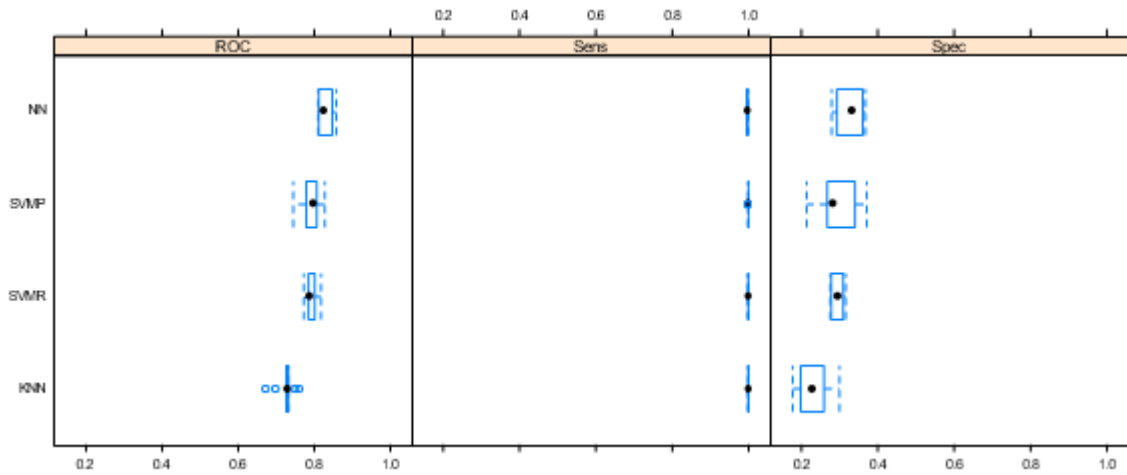


FIGURE 2. (b) Accuracy, ROC, Sensitivity, and Specificity Box Plots

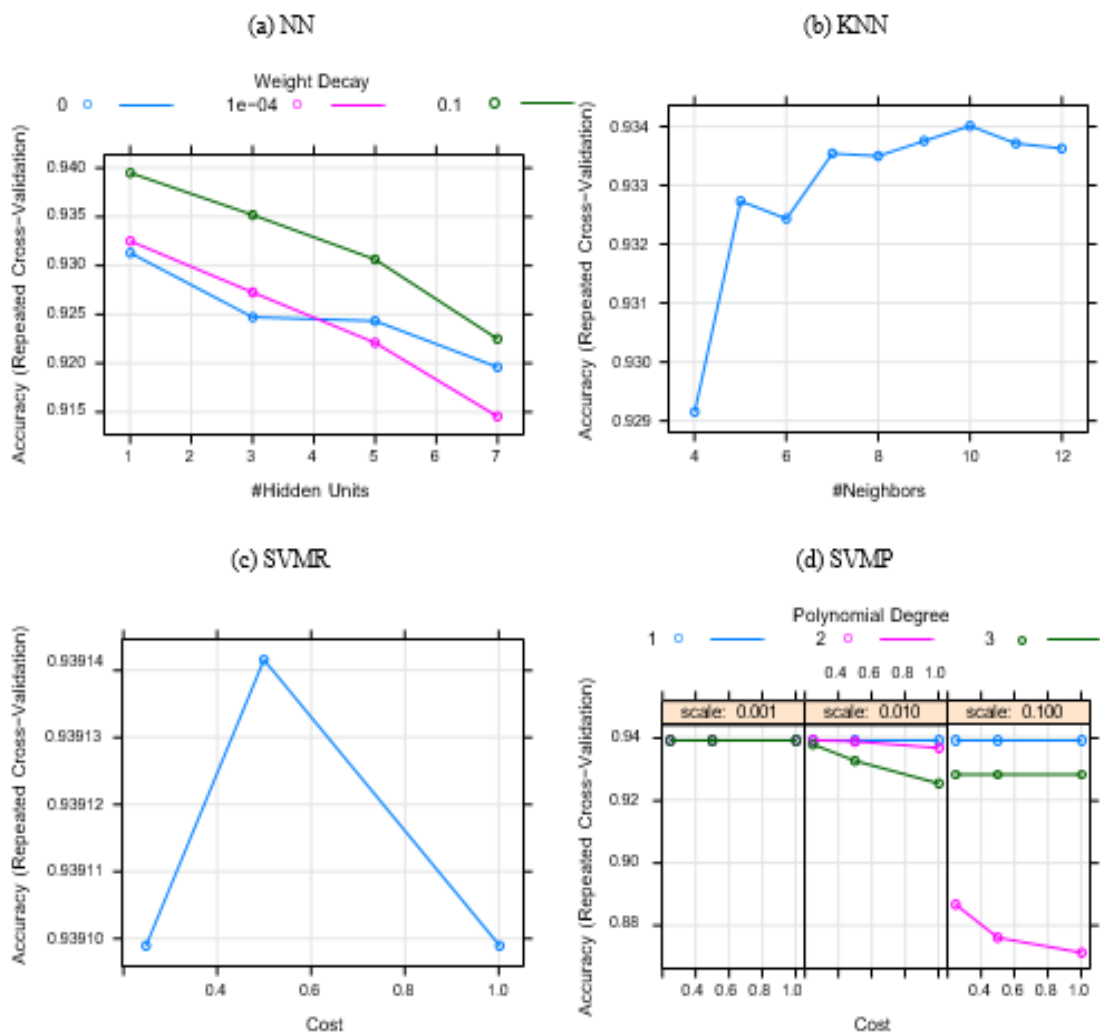


FIGURE 2. Tuning Parameters for Classification Models

This paper showed promising results. However, there is much more work that can be done. More advanced classifiers can be used, such as state-of-the-art ensemble classifiers to hopefully provide a more accurate prediction of dropouts. Bayesian networks may be advantageous in identifying dependent relationships between variables. Another area of future work is to go one step further and learn a Markov decision process to determine a policy of action to help prevent the student from dropping out. Similar ideas could also be applied to identify

other potential risks in students, such as identifying students who may not go on to college. There are many more areas where machine learning could be applied to education.

REFERENCES

- [1]. Balfanz, R., Bridgeland, J.M., Bruce, M., Fox, J.H.: Building a grad nation progress and challenge in ending the high school dropout epidemic. Tech. rep., Civic Enterprises, Everyone Graduates Center at Johns Hopkins University, and America's Promise Alliance Alliance for Excellent Education (2012)
- [2]. Bonferroni, C.E.: Il calcolo delle assicurazioni su gruppi di teste. Tipografia del Senato (1935)
- [3]. Bridgeland, J.M., Morrison, K.B., DiIulio, J.J.: The silent epidemic: Perspectives of high school drop outs. Tech. rep., Civic Enterprises (March 2006)
- [4]. Chang, Y.W., Hsieh, C.J., Chang, K.W., Ringgaard, M., Lin, C.J.: Training and testing low-degree polynomial data mappings via linear svm. *The Journal of Machine Learning Research* 11, 1471–1490 (2010)
- [5]. Chen, C., Chen, Y., Liu, C.: Learning performance assessment approach using web-based learning portfolios for e-learning systems. *IEEE Transaction on Systems, Man, and Cybernetics* 37(6), 1349–1359 (2007)
- [6]. Cover, T., Hart, P.: Nearest neighbor pattern classification nearest neighbor pattern classification. *IEEE Transactions of Information Theory* 13(1), 21–27 (Jan 1967)
- [7]. Delen, D.: A comparative analysis of machine learning techniques for student retention management. *Decision Support Systems* 49, 498–506 (June 2010)
- [8]. Ingels, S., et al.: High School Longitudinal Study of 2009 (HSL:09) Base Year to First Follow-Up Data File Documentation. National Center for Education Statistics (November 2013)
- [9]. Kaufman, P., Bradbury, D.: Characteristics of at-risk students in NELS:88. Tech. rep., National Center for Education Statistics (August 1992)
- [10]. Kotsiantis, S.: Use of machine learning techniques for educational proposes: a decision support system for forecasting students' grades. *Artificial Intelligence Review* 37, 331–344 (May 2012)