

STUDENT CAREER PREDICTION SYSTEM USING MACHINE LEARNING

Mr. Gsai Krishna¹, K.ajay², raghavendrasrinivas³, prem kiran⁴, vinayreddy CH⁵.

Associate Professor¹, Students^{2,3,4,5}.

DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING^{1,2,3,4,5}.
MALLA REDDY INSTITUTE OF TECHNOLOGY AND SCIENCE, HYDERABAD, INDIA.

*Email ID: gasikrishna5555@gmail.com, ajaykyadaveni2001@gmail.com², remellaraghavendra6@gmail.com³, premkiran3107@gmail.com⁴,
vinayreddyvinay4518@gmail.com⁵.*

1. Abstract:

There are many good schools and colleges in India. But most of the students are dropping their education because of various reasons. There are a lot of reasons, some of the students have some financial problem with their family, some of the students don't have interest towards their next level of education, and some students think about gender and some about rural areas don't have good schools and educators. In today's world choosing an appropriate career path is one of the most important decisions but with the increase in the number of career options and opportunities, it makes this decision even more difficult for the students. Different people suggest different career paths to the student but at last the student has to select their career themselves. According to the survey conducted by the Council of Scientific and Industrial Research's (CSIR), about 40% of students are confused about their career selections. This may lead the students to wrong career selection and then working in an area which was not meant for them, this leads their career in wrong path, and this may not be good for their future career. Therefore, it is very important to take the right decision regarding your future career at the right time. So, this proposed method deals with whether the students will be going to the next level of higher education or not. This can be decided with the concepts of machine learning which is the subset of artificial intelligence. Machine learning is made up of the concepts of Mathematics and Science. This paper deals with the student's career prediction by using various machine learning algorithms like Decision Tree (DT), Random Forest (RF), Support Vector Machine (SVM) and Adaboost. Machine learning algorithms are implemented by using Python programming language.

Keywords: Machine Learning Classifier, Support Vector Machine, Adaboost, Random Forest, Decision Tree.

1. Introduction:

Traditionally students' careers were predicted by questions and answers. But this method takes a lot of time. Now, computing technologies play an important role in various fields. Machine learning is one of the newest computing techniques. In this digital world machine learning is used in various fields and industries such as image processing, classification, clinical analysis, regression and more and more. It has the capability of developing and studying automation without being explicit. Machine learning is of three types i.e supervised machine learning, unsupervised machine learning and reinforcement machine learning algorithms. In simple words Machine learning is the science of learning and behaving like humans. It is very important to analyze the ability of the students and they should be directed in the right pathway. This research work includes the concepts of machine learning which are applied to detect the next level of education of the students. This prediction is very important for almost all types of educational institutions, recruiters and so on. Based on the result of this prediction accuracy, the educational institutions find the people with low performance and provide the proper training to them to improve their performance. Job providing companies also spend a lot of amount for selecting a qualified candidate. The output of the prediction model is also used to find the status of the students, if they are interested in going to the job or they are interested to do their higher studies. This research work mainly focuses on the career prediction of undergraduate level students. Machine learning algorithms such as SVM, DT, RF and Adaboost classifiers are used to construct

the model. Among the above classifiers RF produces better results. These classifiers are implemented with the help of python programming language, because most of the real time problems are easily implemented by using this programming language. Next section deals with the views and approaches that are used by various authors in career prediction research domain.

2. Literature Review:

A recent fact provides information using student's data based on their behavioral aspects to forecast the career path. Min Nie et al. proposed a novel model known as ACCBOX (Approach Cluster Centers Based On XGBOOST) to forecast student's career.

The final result clearly states that the current method is better than other methods of prediction. This model uses 13 behavioural data which are collected from 4000 students [1]. Mining student's educational data is also one of the important tasks in the education field. In the beginning days data mining methods were used in the education field by using a smaller number of arguments, because low record maintenance in concern institutions. Recently the large volume of data can be stored on the basis of student. In India 0.3 % people only move forward from their PG level to research level. This prediction task evaluates performance of the students by using various arguments and the students are classified as low, high and medium type. To execute this process the authors K. B. Eashwaret al., combined SVM and K-means methods. A SVM concept is used for classification purposes and K-mean technique is mainly used for clustering the student's data [2]. Predicting the student performance level is one of the important tasks in education domain.

Data mining concepts are used to predict the student's performance by using various parameters. Ankita Kadambande et al., uses

semantic rules and SVM concepts for the predictions. Semantic rules are used to improve the quality of educational content and convey educational action to every student. Here the authors help the students by providing better suggestions and recommendations for improving student's performance level in forthcoming exams. This system will provide help to low level and high-level students and also to increase the student's interest in their education. The main purpose of this research work is to increase the quality of learning measures and support the students by forecasting their academic level which will help the students to a great extent [3].

3. Previous Research:

In the past two decades, we have seen a large number of high-quality works using students' academic performance and learning behavioural data to predict outcome variables, such as standard test score, dropout from school, college enrollment, and major subject choices. For example, Feng, Heffernan, and Koedinger (2009) investigated how students' interaction data extracted from the ASSISTment platform can be used to reliably evaluate students' math 19 Journal of Educational Data Mining, Volume 12, No 2, 2020 proficiency. They were especially interested in building features related to student help seeking behaviours and used the Bayesian Information Criterion (BIC) to compare linear regression models with different groups of predictors. They presented that students' end-of-year exam scores can be better predicted by leveraging the interaction data that reflect assistance requirement, effort, and attendance.

Instead of using traditional descriptive variables in college enrollment research, such as family background, career aspiration, and assessment scores, San Pedro, Baker, Bowers, and Heffernan (2013) studied how student online learning behaviours observed in middle school related to their college choice. They built a logistic regression model using automatically generated affect and

engagement features to achieve decent accuracy at predicting college attendance. Their study was further extended to predict STEM and Non-STEM college major enrollment by San Pedro, Ocumpaugh, Baker, and Heffernan (2014).

Pardos, Baker, San Pedro, Gowda, and Gowda (2014) also studied the ASSISTments system, but they focused on the correspondence between student affect and behavioural engagement and scores on a high-stakes math exam. Using eight machine learning models, they created a set of affect and engagement behaviour detectors to estimate the probability that a student is in a state of boredom, engaged concentration, confusion, and so on. They further created a model to predict students' math exam scores and showed that the constructed detectors helped the model achieve high prediction accuracies.

Knowles (2015) described how to create a state-wide dropout early warning system that can accurately predict the probability of graduation for high school students in the State of Wisconsin. The paper systematically demonstrated the workflow of the whole system, from data cleansing to model training and searching. To balance the trade-off between the correct classification of dropouts and false alarms, the receiver-operating characteristics (ROC) metric is used to identify the best models from a large collection of aspirants, from linear logistic regression models to complex nonlinear models, such as support vector machines. This work was also implemented in the open-source R package, EWStools (Knowles, 2014).

Baker, Berning, Gowda, Zhang, and Hawn (2019) presented a case study on automatically identifying students that have a high risk of dropping out of high school, using data on students' discipline, attendance, course-taking, and grades. The logistic regression

model used in the study helped the authors not only select students at risk, but also found which factors played the largest roles in prediction, which provided information to educators that can be used in individualized involvements.

4. Problem Identification:

Selecting an appropriate career is one of the most important decisions and with the increase in the number of career paths and opportunities, has made this decision quite difficult for them. According to the survey conducted by the Council of Scientific and Industrial Research's (CSIR), about 40% of students are confused about their career options. This may result in the selection of the wrong career and then working in a field which was not meant for them, thus reducing the productivity of human resources. Therefore, it is quite important to take the right decision regarding the career at an appropriate age to prevent the consequences that result due to wrong career selection. This system is a web application that will help students studying in high schools to select a course for their career. The system will suggest the student, a career option based on their personality, skills, interest and their capacity to take up the course.

5. Methodology:

Machine Learning is a technique in which the machines are trained in such a way that it achieves the ability to respond to specific input or scenario based on the past knowledge it has learnt. Simply, it gives computers the ability to learn by using statistical techniques. With the help of Machine learning the computers gain ability to act in spite of being explicitly programmed. This aims at reducing the human involvement in machine dependable problems and scenarios. This helps in solving very difficult tasks and problems very easily and without involving much human effort. Various applications of machine learning include NLP, classification, image recognition, prediction, medical diagnosis, algorithm building, self-driving cars and much more. This paper involves classification and prediction. Let us see the details of classification and prediction. Most of the problems in machine learning can be resolved using supervised and unsupervised learning. If the final class labels are previously known and all the other data items are to be assigned with one of the available class labels, then it is called supervised. And if the final output classes and sets are not known and it is done by identifying the similarity between data point and their characteristics and then they are made into groups based on these characteristics then it is called unsupervised. Classification falls under supervision. On the basis of the regression as well. Based on the type of problem the appropriate model is chosen.

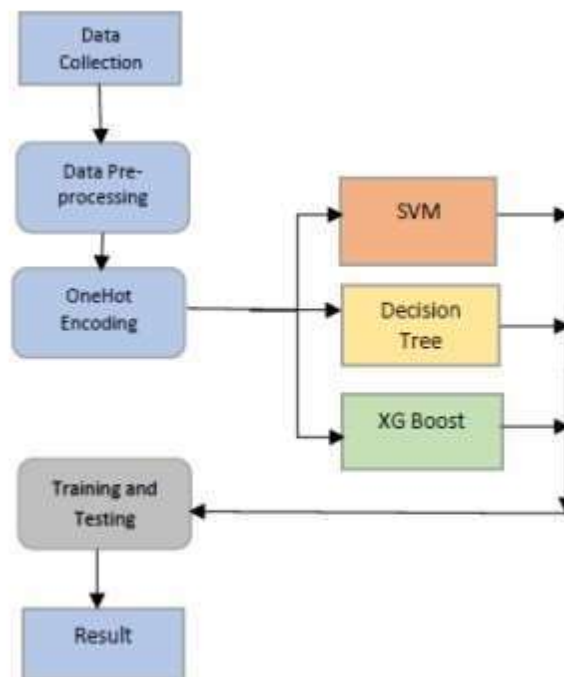


Fig.1: Process Flow Diagram of Proposed System.

Here algorithms like SVM, OneHot encoding, Decision tree and XG boost are used. After training and testing the data with these we consider the most accurate results given by the algorithm for our further processing. So, the initial task done is predicting the output using all algorithms proposed above and later analyzing the results and there after we continue with the most accurate algorithm. So, this paper deals with various advanced machine learning algorithms that involve classification and prediction and are used to improve the accuracy for better prediction, reliability and analyzing these algorithms performance.

6. Implementation:

A. Data Collection:

Data Collection is one of the most important tasks for any project of machine learning. Because the input we feed to the machine learning algorithms is data. So, the algorithms efficiency and accuracy depend on the correctness and quality of data collected. Accurate the data accurate will be the output. For student career prediction many parameters like students' academic scores in various specializations, subjects, programming and analytical capabilities, memory, personal details like relationship, interests, competitions, hackathons, sports, workshops, certifications, books interested and many more are required. All these factors play an important role in deciding a student's progress towards a career area and all these are taken into consideration. Data is collected in many ways. Some data is collected from employees working in different organizations, some amount of data is randomly generated and other from college alumni database. Totally nearly 15 thousand records with 35 columns of data are collected.

B. Data Pre-processing:

Collecting the data is one task and making that data useful is another important task. Data collected from various sources will be in an unorganized format and there may be a lot of null values, invalid data values and unwanted data. Cleaning of all these improper data and replacing them with appropriate or approximate data and removing null and missing data and replacing them with some fixed alternate values are the basic steps in pre-processing of data. Even data collected may completely contain garbage values. It may not be in the exact format or way which is meant to be. All such data must be verified and replaced with alternate values to make data meaningful and useful for further processing. Data must be kept in an organized format.

C. OneHot Encoding:

OneHot Encoding is a technique by which categorical values present in the collected data are converted into numerical or other ordinal format so that they can be provided to machine learning algorithms and get better results of prediction. OneHot encoding simply transforms categorical values into a form that best fits as input to be fed to various machine learning algorithms. This algorithm works good with almost all machine learning algorithms. Few algorithms like random forest handle categorical values properly. In such cases this encoding is not required. The process of OneHot encoding seems to be difficult but most modern day machine learning algorithms take care of that. The process is easily explained here: For example, in a set of data if there are values like yes and no, integer encoder assigns values like 1 and 0 to them. This process can be followed as long as we continue the fixed values as 1 for yes and 0 for no. As long as we allocate or assign these fixed numbers to these particular labels, this is called integer encoding. But here consistency is very important because if we invert the 1 2 | P a g e encoding later, we should get back the labels correctly from those integer values especially in the case of prediction. Next step is to create a vector for each integer value.

Let us suppose this vector is binary and has a length of 2 for the two possible integer values. The label 'yes' encoded as 1 will then be represented with vector [1,1] where the zeroth index is given the value 1. Similarly, label 'no' encoded as '0' will be represented like [0,0] which represents the first index is represented with value 0. For example, [pillow, rat, fight, rat] becomes [0,1,2,1]. This is imparting an ordinal property to the variable, i.e. pillow < rat < fight. As this is ordinal characteristic and is usually not required, OneHot encoding is required for correct representation of distinct elements of a variable. It makes representation of categorical variables to be more expressive and meaningful.

D. Machine Learning Algorithms: -

1. SVM:

SVM denotes Support Vector Machine. It is a supervised machine learning algorithm which is generally used for both regression and classification type of problems. The main applications of SVM can be found in various classification problems. The typical procedure of the algorithm is first each data item is to be plotted in a n-dimensional space, where n is the number of features and the value of each feature being the value of that particular coordinate. Next step is to classify the hyper-plane that separates the two classes very finely.

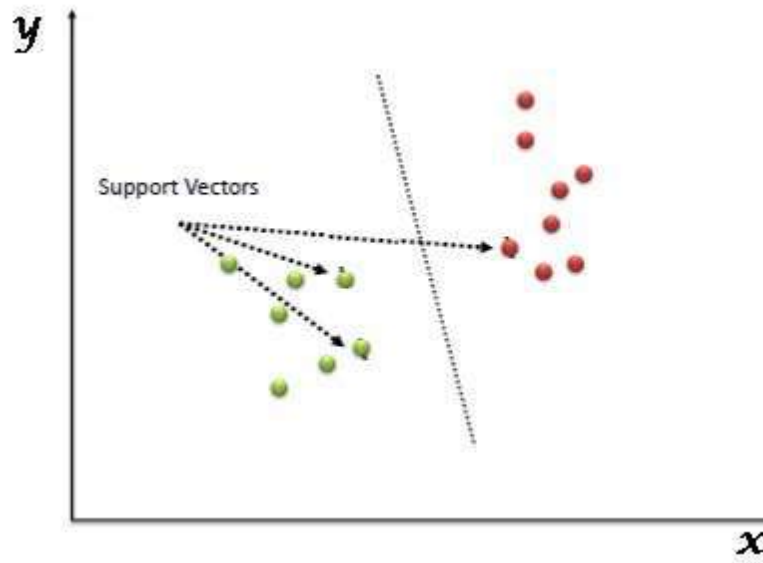


Fig. 2: Support Vector Machine

SVM algorithms are practically implemented using kernels. There are three types of SVM's. In linear SVM hyperplane is calculated or found by transforming the problem using linear algebra. The concept is that SVM can be rephrased by using the inner product of two observations. The sum of the multiplication of each pair of inputs is called the inner product of two vectors. The equation for dot product of a input x_i and support vector x_i is $f(x) = B_0 + \sum(a_i * (x, x_i))$. Instead of using the dot-product, a polynomial kernel can be used, for example: $K(x, x_i) = 1 + \sum(x * x_i)^d$ and not only that a more complex radio kernel is also there. The general equation is $K(x, x_i) = \exp(-\gamma * \sum((x - x_i)^2))$.

2. XG Boost:

XGBoost denotes eXtreme Gradient Boosting. XGBoost is an implementation of gradient boosting algorithms. It is available in many forms like tools, library etcetera. It specifically focuses on model performance and computational time. It reduces the time and lifts the performance of the model greatly. Its implementation has the features of scikit-learn and R implementations and also have newly added features like regularization. Regularized gradient boosting means gradient boosting with both L1 and L2 type regularizations.

The main best features that the implementation of the algorithm provides are: Automatic handling of missing values with sparse aware implementation, and it provides block structure to promote parallel construction of tree and continued training which supports further boost an already fitted model on the fresh data. Gradient boosting is a technique where new models are created that can predict the errors or remains of previous models and then added together to make the final prediction.

They use gradient descent algorithms to reduce loss during addition of new models.

They support both classification and regression type of problems. In the training part generally, an objective function is defined. For example,

$$obj = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^n \Omega(f_i)$$

3. Decision Tree:

Decision Tree is an extremely popular and one of the simple and easy techniques to implement machine learning classification problems. Decision trees are the basic foundation for many advanced algorithms like bagging, gradient boosting and random forest.

The XG Boost algorithm mentioned above is the advanced version of this general decision tree. The commonly used decision trees are CART, C4.5, C5 and ID3. A node denotes a input variable (X) and a split on that variable, assuming the variable is numerical.

The leaf, also called the terminal nodes of the tree, possesses an output variable (y) which is vital for prediction. The typical scenario that a decision tree follows is first selecting a root node. Then calculate information gain or entropy for each of the nodes before the split.

Then select the node that has more information gain or less entropy. Then split the node and reiterate the process. The process is iterated until and unless there is no possibility to split, or the entropy is minimum. Entropy is the measure of uncertainty or randomness of data. Information gain is the measure of how much entropy is reduced before to after split.

$$H(S) = \sum_{x \in X} -p(x) \log_2 p(x)$$

$$IG(A, S) = H(S) - \sum_{t \in T} p(t)H(t)$$

7. Result:

The data was trained and tested with all three algorithms and out of all SVM gave more accuracy with 90.3 percent and then the

XG Boost with 88.33 percent accuracy. As SVM gave the maximum accuracy, for all further data predictions SVM will be followed.

So, finally, we concluded that SVM provides better result. So, this algorithm can be used in the background and the new

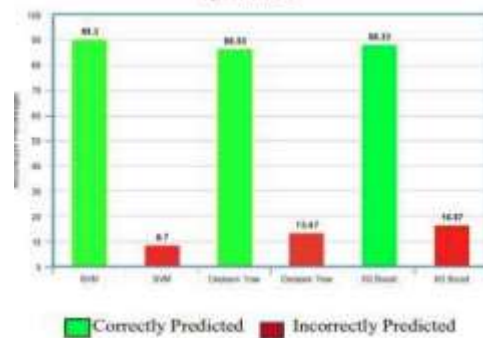


Fig. 3: Final Output Graph

8. Future Scope: A more powerful web application or mobile application can be developed where inputs will not be given directly, instead student parameters will be taken by evaluating students through various evaluations and examining processes. Technical, logical, analytical, psychometry, memory based, and general awareness, interests and skill-based tests can be designed, and parameters can be collected through them so that results certainly will be more accurate, and the system will be more reliable to use.

References:

- [1]. P.KaviPriya, "A Review on Predicting Students' Academic Performance Earlier, Using Data Mining Techniques", *International Journal of Advanced Research in Computer Science and Software Engineering*.
- [2]. Ali Daud, Naif Radi Aljohani, "Predicting Student Performance using Advanced Learning Analytics", *2017 International World Wide Web Conference Committee (IW3C2)*.
- [3]. Mariam-E-Jannat, SaymaSultana, Munira Akther, "A Probabilistic Machine Learning Approach for Eligible Candidate Selection", *International Journal of Computer Applications (0975 – 8887) Volume 144 – No.10, June 2016*. [4] Sudheep Elayidom, Dr.Sumam Mary Idikkula, "Applying Data mining using Statistical Techniques for Career Selection", *International Journal of Recent Trends in Engineering, Vol. 1, No. 1, May 2009*.
- [4]. Dr.Mahendra Tiwari, Manmohan Mishra, "Accuracy Estimation of Classification Algorithms with DEMP Model", *International Journal of Advanced Research in Computer and Communication Engineering Vol. 2, Issue 11, November 2013*.
- [5]. Ms.Roshani Ade, Dr. P. R. Deshmukh, "An incremental ensemble of classifiers as a technique for prediction of student's career choice", *2014 First International Conference on Networks & Soft Computing*.
- [6]. Nikita Gorad, Ishani Zalte, "Career Counselling Using Data Mining", *International Journal of Innovative Research in Computer and Communication Engineering*.

- [7]. *Bo Guo, Rui Zhang, "Predicting Students Performance in Educational Data Mining", 2015 International Symposium on Educational Technology.*
- [8]. *Ali Daud, Naif Radi Aljohani , "Predicting Student Performance using Advanced Learning Analytics". [10] Rutvija Pandya Jayati Pandya, "C5.0 Algorithm to Improved Decision Tree with Feature Selection and Reduced Error Pruning ", International Journal of Computer Applications (0975 – 8887) Volume 117 – No. 16, May 2015.*