

Engineering College Admission Predictor Using Bagging Ensemble Modelling

Mangesh Balpande

Department of Information Technology,
SVKM's Institute of Technology Dhule, India
mangesh.balpande111@gmail.com

Sushrut Patil

Department of Information Technology,
SVKM's Institute of Technology, Dhule, India
sushrutpatil723@gmail.com

Vedant Ranalkar

Department of Information Technology,
SVKM's Institute of Technology, Dhule, India
vedantranalkar7@gmail.com

Gayatri Rana

Department of Information Technology
SVKM's Institute of Technology, Dhule, India
gayatrirana529@gmail.com

Aadil Shaikh

Department of Information Technology
SVKM's Institute of Technology, Dhule, India
shaikhaadilshafiullah5558@gmail.com

Abstract— Currently, a number of students make mistakes when making their college choice list for a number of reasons, such as inaccurate college analysis, insufficient knowledge, and nervous projection. The efficiency of admissions processes has a major impact on how well higher education institutions perform. In this study, we offer a data mining method for forecasting engineering program acceptance of students. Our focus is on obtaining critical information from the student's academic record, including test results, grades, and other relevant data. In order to predict each student's acceptance chances, we next feed these data into a machine learning model, such as a Random Forest or decision tree. By using the power of data mining, we can predict the engineering admissions result more accurately and effectively than we could with earlier methods. We also compare the performance of our model with other, more traditional methods like logistic regression. According to our conclusion, our method performs better at predicting admission than conventional methods. In addition, we examine the implications of our results and suggest possible enhancements to the accuracy of our model.

Keywords— Data Mining, Data Analysis, Cutoff, Preference List, Accuracy Score, Precision, Recall.

I. INTRODUCTION

Machine learning techniques are employed in the "Engineering College Admission Prediction" project to estimate a student's probability of being admitted to an engineering college. This forecast is based on a different number of variables, such as extracurricular activity participation, academic achievement, and personal history.

This assignment is very important since it can help students understand how they can be admitted to particular colleges. They can take the required actions to improve their applications thanks to it. Additionally, it helps universities find applicants who have a better chance of succeeding in their courses, making admissions easier. Using a range of machine learning algorithms, the project includes data gathering, pre-processing, and analysis [1][2]. To find patterns and trends in the admissions process, historical admission data from engineering institutions will be collected first, then preprocessed and analyzed.

The ultimate aim of this research is to develop a prediction model for engineering college admission using machine learning. The probability of admission for each student will be predicted by the model using a range of data sources, such as extracurricular activities, academic performance, and student

background information [3][4].

The model will be used to forecast the likelihood of admission for fresh applicants after it has been trained.

An intuitive interface will display the results, enabling students to input their information and obtain an approximation of their chances of being admitted. This project has substantial potential benefits [5]. The application can assist students in making well-informed decisions about where to apply and how to increase their chances of being accepted [6]. Colleges will be able to enhance their admissions process and find the best applicants for their programs with the aid of this tool [7].

The suggested system compiles a list of colleges that an individual is eligible to attend using Machine Learning algorithms. Neural networks, Random Forest, Decision Trees, and Linear Regression are among the algorithms [8].

After comparing these algorithms, the prediction system will be developed using the algorithm with the best key performance indicators [9]. Additionally, the system groups universities according to a profile and assigns a likelihood of acceptance—high, low, etc.—to each group. All things considered, this engineering college admission prediction project has the power to transform the admissions procedure and enhance outcomes for colleges and applicants alike.

II. LITERATURE REVIEW

The goal of the research is to forecast, from the perspective of the college, the probability that a student will enroll after verifying the course offerings [10]. The probability of students enrolling in college was predicted using the K-Means calculation based on a number of variables, including inspiration, guardian capacity, family income, and occupation. The representation aims to increase college enrollment by taking into account how well each student's attributes match the requirements. Using AI and foresight showing, we created a representation to assess Tennessee Tech University's validation processes and benchmarks. The representation is based on J48, a well-known C4.5 computation variation. They used the numerous understudy components listed to assess their chances of admission to the college, just like the representations previously mentioned [11]. When the student had a strong outline to support the affirmation, the illustration performed a great job of predicting the real positive scenarios. However, it completely failed to identify the real negative scenarios when the student did not meet the specified requirements. In this study, led by Jamison (2017), the outcome of college affirmation was

predicted using AI approaches. The frequency with which students obtain permission from their institution to enroll in a program is known as the yield rate. Among the AI methods applied to create the representation were SVM, Random Forest, and Logistic Regression [12]. The educational au-courant methodology for business processes has made extensive use of automated techniques. Techniques based on artificial intelligence and conventional methods can be distinguished from one another. Conventional techniques have many multivariate research methods as characteristics, but artificial intelligence methods often have the adept schemes methodology [13]. They were combined and used in a number of instructional and teaching apps, such as admission prediction.

The description describes a web application that addresses issues with the admissions process and is intended to help students choose engineering colleges in Maharashtra. The program includes comprehensive college comparisons, a recommendation system based on student interests, and data analysis to forecast cutoffs. Technical information on the recommendation system and the standards utilized for in-depth comparisons is missing from the abstract. An outline of the techniques or algorithms used would give the abstract more substance. The abstract refers to "comparisons between institutions," however in order to give readers a better idea of the functionality of the application, it would be beneficial to include the important factors that were taken into account in these comparisons.[23]A complete approach designed to address issues with engineering college admissions is described in the abstract. It presents a predictor for college admission that takes into account academic, personal, and admission requirement characteristics. It does this by using machine learning and historical data. The focus on database administration and the correctness of the expected results demonstrates a dedication to accuracy. The items to take into account include giving more precise information about the machine learning model, the kinds of difficulties that students encounter, and the criteria that the predictor uses [24].

A comprehensive dataset comprising 41359 institution applications has been generated by researchers in order to estimate four-year bachelor's graduation rates using a generalizable methodology. This method takes into account various factors such as sociodemographic information, test scores, work experience, academic achievement, participation in extracurricular activities, and teacher evaluations. But there are drawbacks to the strategy [16]. First off, the detailed illustration only defines 41359 out of a possible 278201 applicants, even though the proportions of the studied data cluster are significant.

The students will gain insight into how likely it is that a college will accept their application. As a result, it will help students categorize the colleges that best fit their profiles and give them information about those universities [17]. This analysis's limitation was that it only took into account a limited number of universities with varying rankings. A web application could be used to advance the system, but more details about more universities and specializations would need to be provided. Consequently, our research has attempted to address this issue that the student admission community is facing. Furthermore, our Android application development helps the client choose the right institutions for option entry.

III. PROPOSED WORK

A. System Overview

A piece of software called the engineering admissions predictor was developed to estimate a student's likelihood of being admitted into an engineering school based on a number of factors. In order to produce a predicted probability of admission, the system would receive a set of input variables pertaining to the student's academic performance, test results, extracurricular activities, and other pertinent information. The system can be divided into three main components:

Data Input and Preprocessing: The system's input data collection and processing are handled by this component. The student's academic record, test results, personal statements, and other pertinent papers would all be input into the system. Preprocessing would be performed on the data to eliminate any unnecessary or incomplete information and convert it into a format that the prediction model could use.

Prediction Model: This part of the system is in charge of forecasting the student's chances of getting admitted. Numerous machines learning algorithms, including logistic regression, decision trees, and neural networks, can be incorporated into the prediction model. A probability score that represents the likelihood of admission would be the model's output after it had been trained on historical data on student admissions.

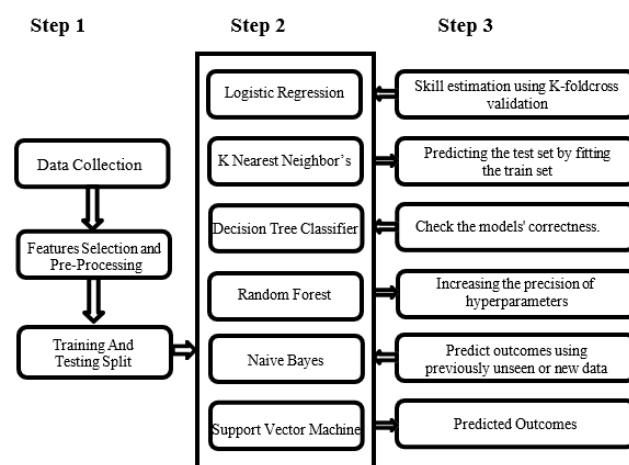


Figure 1. System Architecture

Figure 1: System begins by Data Collection, depicts the suggested system architecture. The engineering admissions predictor project typically consists of several components, including:

Data Collection: The system gathers information on a range of variables, including academic standing, test results, extracurricular involvement, and other pertinent details, that could impact a student's chances of admission.

Data Preprocessing: To remove any errors, outliers, or missing values, the collected data is preprocessed. By doing this, you can be sure the data is appropriate for analysis and forecasting. Numerous machines learning techniques, including logistic regression, neural networks, and decision trees, are trained and tested on preprocessed data. These algorithms look for patterns in the data and create a prediction model that can be used to gauge an applicant's likelihood of being admitted.

User Interface: Through the system's user interface, users can input data and get results according to their likelihood of admission. Additionally, the interface might provide advice and suggestions to help students increase their chances of being admitted.

Prediction Model: The system's central component is the prediction model, which processes the input data and produces an admissions probability score. The patterns and trends found in previous admissions data are used to calculate this score. *Past Admissions Information:* To train and evaluate the prediction model, the system makes use of past admissions information. Past admissions records, academic standing, and other pertinent data are included in this data.

B. Methodology

The workflow diagram for the planned work is shown below. Gathering and organizing data for the machine learning model's training and testing is the first step. To make the data suitable for the model, this might involve cleaning and modifying it. Features are the variables or inputs that the model will use to make predictions. In this step, the most pertinent features are chosen, and new features that could increase the accuracy of the model may be created. The steps to selecting a machine learning model suitable for the task at hand, like regression or classification, are as follows. After that, the models are trained on the prepared data, and their parameters are adjusted to enhance the model's functionality.

A test set of data that was not used during training is used to evaluate the model after it has been trained. The accuracy, precision, recall, and other metrics of the model are calculated to ascertain how well it works with precision in the field. The model can be used to predict data after it has been verified and found to perform satisfactorily. For this model to continue being accurate and useful, it needs to be improved and maintained continuously. This entails keeping an eye on the model's output over time, retraining it with fresh data, and adjusting the model's parameters or architecture as necessary.

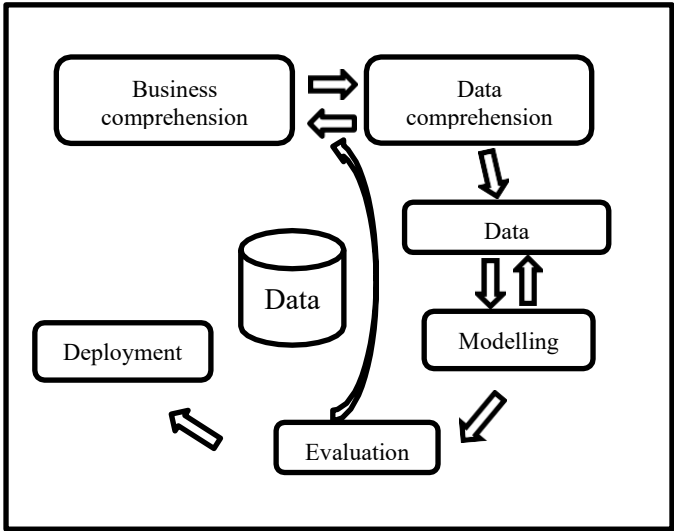


Fig 1.1 Workflow Diagram

During the first phase of data collection, we used the official MAHACET website, which provided information on every student admitted into multiple Maharashtra colleges over the previous five years. We prepared a dataset as pre-records and cleaned the dataset by removing any unnecessary columns in order to retrieve this data using a data scraping technique.

Se. no	Merit No	MHT-CET	Application ID	Name of the student	Gender	Candidate Category	Seat Type	College Name	Location
1	18898	89.889223	EN22182921	SHEGOKAR SAMIKSHA VISHNU	F	SC	LOPENS	Government College of Engineering ,Amravti	Amravti
2	21857	88.241971	EN22169336	GODBOLE HARSHAL RAVI	M	OBC	GOPENS	Government College of Engineering ,Amravti	Amravti
3	22005	88.137472	EN22244230	KAYANDE KALYANI PRAMOD	F	NTS(NT-D)	LOPENS	Government College of Engineering ,Amravti	Amravti
4	22128	88.091617	EN22135944	RATHOD RITESH VILAS	M	DT/VJ	GOPENS	Government College of Engineering ,Amravti	Amravti
5	22579	87.899633	EN2209220	CHAUDHARI AVINASH SHIRIRAM	M	OBC	GOPENS	Government College of Engineering ,Amravti	Amravti

Fig 1.2 Dataset

The following dataset consists 9 columns in which Merit Number, CET Score, Name, Gender, Category, Seat Type, Institute Name and Location are considered respectively. The dataset has 31000 rows containing information of single student per row.

After loading the dataset into a Pandas Data Frame, superfluous columns are removed to leave only relevant features. One-hot encoding of the 'Gender', 'Seat Type', and 'Location' columns is then performed to represent categorical data as binary vectors. The Data Frame that is produced is solely composed of numerical data, making it appropriate for training a machine learning model.

To train the data, a random forest classifier is selected as the machine learning algorithm. To arrive at a final prediction, this algorithm builds several decision trees and adds together their predictions. The number of decision trees (estimators) and the random state that are employed to guarantee reproducibility are the model's hyperparameters. After that, the model is fitted to training data, which enables it to identify patterns in the data that can be utilized to make predictions.

Based on the input data, the model is then used to predict the name of the college. The model produces a numerical label as its output, which must be decoded in order to obtain the college's name. A message stating that no matching college was found is displayed if no college that meets the input criteria is found.

C. Prediction Models

This model frequently makes use of supervised machine learning algorithms, like decision trees or logistic regression, to study historical admission data and find associations and patterns that could be utilized to predict admissions in the future. These models may take into account things like extracurricular activities, articles, letters of recommendation, academic achievement, and results on standardized tests. The model is continuously adjusted and improved to improve its performance, and metrics like precision and recall are frequently used to gauge its accuracy. The model can be used to determine the likelihood of admittance for new applications after it has been trained and evaluated.

Predictive models for student admissions have their uses. Students can use the forecasts to assess their chances of acceptance and choose their application locations with knowledge. It is crucial to stress that admission prediction models should be used to support human decision-making rather than to replace it, as many criteria cannot be adequately represented by statistics alone.

D. Algorithms

1. Linear Regression

It is a kind of supervised machine learning technique that makes use of one or more input characteristics to forecast a continuous target variable. The model is referred to as "linear" because it

presumes a linear relationship between the target variable and the input data.

This algorithm's primary goal is to find the best-fit line that minimizes the variation between the values that are predicted and those that are observed. To achieve this, find the line coefficients (slope and intercept) that work to minimize the squared differences between the actual and predicted values.

2. Decision Tree Algorithm

A supervised machine learning method called a decision tree can be applied to both regression and classification problems. Every node symbolizes an attribute or characteristic, every branch denotes a decision rule, and every leaf node denotes the expected result or class label.

Using the most important feature or attribute that best divides the data, the decision tree method seeks to separate the data into subsets. Until a stopping condition is met, like reaching a maximum depth or a minimum number of samples in a leaf node, this process is repeated endlessly. The splitting process is affected by a number of metrics, including entropy, information gain, and the Gini index.

3. KNN Algorithm

KNN is a simple supervised machine learning (ML) method for imputation of missing values, regression, and classification. It can be used to categorize unexpected points based on the values of the closest existing points. It is based on the idea that the observations closest to a specific data point are the most comparable. By choosing K, the user can determine how many neighboring observations to use in the algorithm. When classifying objects, the KNN algorithm is applied, and the outcomes vary depending on the value of K.

How does it operate?

The class to which a new observation should belong is determined by counting the number K of nearest neighbors. Greater K values produce stronger decision boundaries and are more robust against outliers than minuscule K values (K=3 is better than K=1).

4. Random Forest Algorithm

One popular machine learning method for solving regression and classification issues is random forest. Several decision trees are combined in this ensemble learning method to create a strong model that performs better than any one tree alone. In a random forest model, each decision tree is built independently using a randomly chosen set of characteristics and a randomly sampled subset of training data. This promotes greater model generalization and prevents overfitting. The final forecast of the random forest model is then produced by adding together the forecasts of each individual tree. Using random forest models over other machine learning methods has several benefits. They can handle massive datasets with multiple characteristics and are very scalable. They also perform well with both numerical and categorical data, and they can handle missing data without imputation. Among other uses, random forest models are frequently used in image classification, fraud detection, and consumer segmentation. Additionally, they are extensively employed in bioinformatics to predict gene expression levels and protein-protein interactions. For implementation, following steps can be used.

Step 1: Bringing in Libraries

The NumPy, a multi-linear regression classifier, and Pandas are the three most crucial prediction libraries. Pandas are used in information frame-related activities. Additionally, NumPy will be used to perform the required scientific procedures.

Step 2: Reading the dataset

Knowing that the data is pertinent, of excellent quality and of sufficient amount is all that is necessary to choose the proper variables from a fundamental comprehension of the dataset. The finest predictor factors for our representation are being sought for as part of our representation-building activities. CSV file

reading: `df = pd. Read_csv(body).`

Step 3: Splitting the data into training and test sets, The dataset should divide into two subsets:

- A set used to train a representation is a training set.
- Testing the trained representation using a test set.

Step 4: Importing Algorithms:

Sklearn's formulae After importing a linear representation, multilinear regression lowers the sum of squares separating the dataset's observed foci from the predictions made by the linear estimate.

Step 5: Create a representational item.

Step 6: Using our technique to fit the training data.

Step 7: Determining the test result.

As the all things in this tech world has some limitations and has some good potential for the helpful ness of the community as our model has a potential to help the student more effectively by the predicting the college, they can get by their performance in the entrance exam and have less confusion in the admission process and have right college according to their marks and intelligence. The limitations should be discussed clearly the model is based on the previous year cutoff the college regarding the casts and different categories, the information is need to updated after every year, no of total vacancy after every round is need to updated it. These results are totally based on previous records so the cutoff gets changes after every round. It may not adapt the fast changes in the education policy and admission criteria.

IV. RESULTS

The dataset can be used to forecast future college admission chances, analyze student demographic trends, and enhance educational experiences and results. Educators, researchers, and anybody else with an interest in comprehending and enhancing the educational experience of students can use it.

Table I. Accuracy Table

Algorithm	Precision	Recall	F1 Score	Support
Linear Regression	0.71	0.62	0.64	28
Decision Tree	0.67	0.54	0.74	24
Random Forest	0.93	0.89	0.87	47
KNN	0.83	0.73	0.62	35

The information is organized into columns that reflect the merit number, MHT-CET score, application ID, name, gender, category, seat type, college name, and location of college applicants.

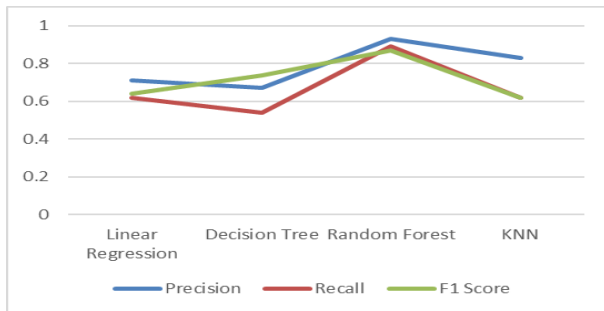
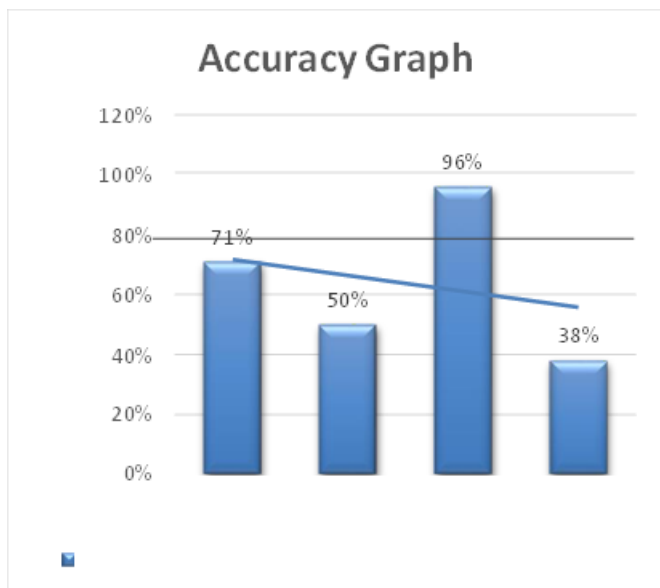


Fig 1.3 Accuracy Graph

This dataset can be utilized for preprocessing, data analysis, and the creation of a classification model to predict a candidate's likelihood of admission to a specific college based on their qualifications.

We can conclude about the various machine learning algorithms based on the data's precision, recall, F1 score, and support value by looking at graph 1.3. From the graph 1.3 we get to conclusion of the various, Machine learning algorithm regarding their Precision, Recall, F1 score and support value the data is being shown. Table 1. Where Random Forest Network is having high precision value as well as the Recall value and F1 Score, this proves that the Random Forest has the more effective than other mentioned machine learning algorithms.



	Decision Tree	Linear	Random Forest	KNN
Accuracy	71%	50%	96%	38%

Fig 1.2 Accuracy Graph in different classifiers

The accuracy of each trained model is displayed in the above graph, with Random Forest having the highest accuracy (96%), and KNN having the lowest accuracy (38%). The Random Forest Algorithm is taken into consideration to produce accurate results because it maximizes accuracy during model training. Our findings indicate that the random forest network Random Forest often produces better results than a single Decision Tree when working with large and complex datasets. It combines the capacity to forecast results from multiple decision trees, which when integrated yields a more accurate and powerful model. Random Forest is a useful technique for handling large datasets because it splits the data and adds multiple trees at once. KNN becomes more computationally expensive as the dataset size increases because it must determine the distances between each data point. Random Forest can be used to handle problems

related to classification and regression. It performs admirably when handling a range of problems.

The practical implementation of this model is in the college prediction can be used by the government as well as the private institute and the student itself by just entering their marks into the system and the system will predict the college can be chosen by the student for their graduation the system will predict the college for their study according to their performance in the entrance example. As the student is not aware about the various courses, cutoffs, for specific course, the system will assist the student or person to choose better college for them.

V. CONCLUSION

In order to make informed decisions, predicting engineering admissions in this position is a critical and difficult problem that requires the use of multiple data sources. In addition, the work needs to be updated and verified on a regular basis as new data sources become available. Machine learning techniques have the potential to simplify and improve the accuracy of engineering admission prediction. More advanced machine learning techniques, such as deep learning and neural networks, will probably be applied to engineering admission prediction in the future in order to improve forecast accuracy. Furthermore, the prediction process can incorporate data from other sources, such as demographics, social media, and even student achievement.

It is demonstrated in this AI college predictor that using a random forest network yields better and more accurate results, as well as more best cases and easier data handling. While other models provide accuracy of 71%, 80%, or 38%, the random forest model yields 96%.

REFERENCES

- [1] M. Injadat, A. Moubayed, A. B. Nassif, and A. Shami, "Multi-split Optimized Bagging Ensemble Model Selection for Multi-class Educational Data Mining," *Appl. Intell.*, vol. 50, pp. 4506–4528, 2020
- [2] Mishra, S. and Sahoo, S. (2016) "A Quality-Based Automated Admission System for Educational Domain," pp. 221 - 223.
- [3] Eberle, W., Simpson, E., Talbert, D., Roberts, L., and Pope, "A. (n.d.). Using Machine Learning and Predictive Representing to Assess Admission Policies and Standards."
- [4] M. N. Injadat, A. Moubayed, A. B. Nassif, and A. Shami, "Systematic ensemble model selection approach for educational data mining," *Knowledge-Based Syst.*, vol. 200, p. 105992, Jul. 2020
- [5] M. S. Acharya, A. Armaan, and A. S. Antony, "A Comparison of Regression Models for Prediction of Graduate Admissions," *Kaggle*, 2018
- [6] G. K. Uyanik and N. Güler, "A Study on Multiple Linear Regression Analysis," *Procedia - Soc. Behav. Sci.*, vol. 106, pp. 234–240, 2013
- [7] C. López-Martín, Y. Villuendas-Rey, M. Azzeh, A. BouNassif, and S. Banitaan, "Transformed k-nearest neighborhood output distance minimization for predicting the defect density of software projects," *J. Syst. Softw.*, vol. 167, p. 110592, Sep. 2020.
- [8] A. B. Nassif, O. Mahdi, Q. Nasir, M. A. Talib, and M. Azzeh, "Machine Learning Classifications of Coronary Artery Disease," in *2018 International Joint Symposium on Artificial Intelligence and Natural Language Processing (ISAI-NLP)*, 2018, pp. 1–6
- [9] A. B. Nassif, M. Azzeh, L. F. Capretz, and D. Ho, "A comparison between decision trees and decision treeforest models for software development effort estimation," in *2013 3rd International Conference on Communications and Information Technology, ICCIT 2013*, 2013, pp. 220–224
- [10] A. B. Nassif, "Software Size and Effort Estimation from Use Case Diagrams Using Regression and Soft Computing Models," *University of Western Ontario*, 2012.

- [11] A. B. Nassif, "Software Size and Effort Estimation from Use Case Diagrams Using Regression and Soft Computing Models," University of Western Ontario, 2012.
- [12] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," MIT Press. Cambridge, MA, vol. 1, no. V, pp. 318–362, 1986.
- [13] M. S. Acharya, A. Armaan, and A. S. Antony, "A comparison of regression models for prediction of graduate admissions," ICCIDS 2019 - 2nd Int. Conf. Comput. Intell. Data Sci. Proc., pp. 1–5, 2019
- [14] N. Chakrabarty, S. Chowdhury, and S. Rana, "A Statistical Approach to Graduate Admissions' Chance Prediction," no. March, pp. 145–154, 2020.
- [15] N. Gupta, A. Sawhney, and D. Roth, "Will I Get in? Modeling the Graduate Admission Process for American Universities," IEEE Int. Conf. Data Min. Work ICDMW, vol. 0, pp. 631–638, 2016
- [16] A. Waters and R. Miikkulainen, "GRADE: Graduate Admissions," pp. 64–75, 2014.
- [17] Jamison, J. (2017). Applying Machine Learning to Predict Davidson College's Admissions Yield, pp. 765–766
- [18] An Improved KNN algorithm Based on Kernel Methods and Attribute Reduction by Wang Xueli, Jiang Zhiyong, Yu Dahai presented in 2015 Fifth International Conference on Instrumentation and Measurement, Computer.
- [19] Communication and Control Prediction Analysis using Random Forest Algorithms to Forecast the Air Pollution Level in a Particular Location by Puli Dilliswar Reddy, L. Rama Parvathy presented in Third International Conference on Smart Electronics and Communication (ICOSEC 2022) ISBN: 978-1-6654-9764-0
- [20] A Predictive Analysis Model of Customer Purchase Behavior using Modified Random Forest Algorithm in Cloud Environment cited by Soumi Ghosh, Chandan Banerjee presented in 2020 IEEE International Conference for Convergence in Engineering published in IEEE in 2020.
- [21] The tree based linear regression model for hierarchical categorical variables by Emilio Carrizosa a,b, Laust Hvas Mortensen c,d, Dolores Romero Morales e, M. Remedios Sillero-Denamiel f,b, published in the springers direct journal with name Expert Systems With Applications 203 (2022) 117423
- [22] An improved decision tree algorithm based on variable precision neighborhood similarity Caihui Liu a, Bowen Lin a, Jianying Lai a, Duoqian Miao published in springers direct journal with name information Sciences 615 (2022) 152–166
- [23] Personalized College Recommender and Cutoff Predictor for Direct Second Year Engineering cited by Abdul Majeed Inamdar, Tanmay Mhatre, Pravin Nadar, Supriya Joshi published in conference 2022 IEEE 7th International conference for Convergence in Technology (I2CT) Pune, India. Apr 07-09, 2022.
- [24] Machine Learning Based Prediction Model for College Admission by Prof. Priya N. Parkhi, Amna Patel, Dhruvraj Solanki, Himesh Ganwani, Manav Anandani published in 2023 11th International Conference on Emerging Trends in Engineering & Technology - Signal and Information Processing (ICETET - SIP)