

# A Pattern Mining Approach for Prediction of School Student Performance using Classification Algorithms

**Khushbu Agrawal**

Research Scholar  
MATS School of Information Technology  
MATS University, Raipur (C.G.), India  
Email: [khushi9r@gmail.com](mailto:khushi9r@gmail.com)

**Dr. Bhavana Narain**

Professor  
MATS School of Information Technology  
MATS University, Raipur (C.G.), India  
Email: [narainbhawna@gmail.com](mailto:narainbhawna@gmail.com)

**Abstract**—Evaluating student performance is very difficult task for academic area. The educational performance play a vital role in classifies the student in higher education. The student performance effect various factor like learning process, personal and social. This paper demonstrates the impact of student positive or negative performance of student success. Here the most commonly used prediction algorithms were LWL, random forest and bagging. After apply this three algorithm we present a novel model combination of this three algorithm name as ensemble. With the help of data mining algorithm we predict the student dropout rate, which is helpful for academic progress. For this we have to collect big and authentic data which can be done through uniform district information system for education DISE).

**Keywords**— student performance, prediction, classification algorithm, educational data mining (EDM), academic performance

## I. INTRODUCTION

Educational data mining is a scientific research area, it use the multiple algorithm to improve academic result and procedure for further decision making. Predicting student performance in academic data is an important issue in e-learning environments. Student performance is based on various factors such as personal, social, psychological and other issues. Data mining techniques is a promising tool to attain these objectives, data mining techniques are use to bring hidden information, patterns and relationship among the large dataset, which help us in categorization of data into knowledgeable facts. To identify the prediction of risk students with a large no. of student data set, it is very difficult and time consuming to using traditional data mining research methods such as questionnaires. Using traditional method in data mining has some limitations like it cannot properly handle the missing values, requires detailed information about the data, and cannot deal with uncertainty or vagueness in any information domain. Various tools and techniques required for achieving the best result from data ining like data cleansing, AI, association rule mining, clustering, regression, machine learning and classification. So the classification is one of the most useful predictive data mining techniques to solve this problem, and customized traditional method by applying various classification techniques. The prediction of student performance with high accuracy is beneficial for identifying the students with low

academic achievements. It is required that the identified students can be assisted more by the teacher so that their performance is improved in future.

### I.I Contribution and organization of paper

Several empirical studies have been conducted the data mining in academic data I have study 12 research paper; most of the research work is done in small scale and structured dataset that are discussed below:

Romero et al., [1] have studied about data mining techniques and present early prediction of student performance. They search 133 research papers and 82 papers were used to solve their five research question.

Viswanathan and Vengatesh Kumar [2] have found that student performance prediction and define methodology for implementing student prediction. They have proposed a novel model of ensemble support vector machine (esvm) as a tool for data mining algorithm.

Namdeo et al., [3] presented a comparative study on the effectiveness from prediction of student's performance through educational data mining techniques. They describe a various research paper related to prediction of student performance and which use different classification technique.

Mangat & Singh Saini [4] has found that predict student performance in higher educational institutions by using machine learning algorithm. Here Naïve Bayesian algorithm provides more accurate result the study was done on seven hundred student's data set.

Zhang et al., [5] have studied about OLP (online teaching learner) learner scores for data mining technique. They have proposed course score analysis model and EM clustering was adopted for score features and salient features were obtained through PCA.

Ahmad et al., [6] have studied about EDM is the process of transforming raw data obtained from educational systems into useful data that can be used to make data-driven decisions.

Hossain et al. [7] have a proposed a new K-means Cluster Algorithm for data mining. Earlier method of K-means cluster techniques (original) have problem that if the number of clusters is to be chosen small then there is a higher probability of adding dissimilar items into the same group.

Ali et al. [8] investigate k-nearest neighbor classification (k-NN) performance on heterogeneous data sets which include the combination of both numerical and binary data. Earlier traditional k-NN works on numerical data only.

Sankari et al. [9] studied about the educational data mining model and learning analytics. The paper also investigate on how educational data mining (EDM) correlate between the student learning performance with other different learning models.

Fairos et al. [10] data mining algorithm was applied to predict student performance either excellent or non-excellent. The student performance was conducted on selected university in Malaysia.

Majeed and Naaz [11] has reviewed 39 research paper for understand the scope and the importance of academic data mining. The Expectation Maximization Algorithm, C4.5 and Page Ranker.

Qiao and Hong [12] have found that didactic of analyzing process data in log file by data mining methods. They proposed four supervised learning methods- Classification and Regression Tree (CART), gradient boosting, random forest, and SVM are explored to develop classifiers and unsupervised methods self-organizing (SO) and k-means methods are use by students to examine same and different score.

## II. Material and methods

In the study the main problem is to organise the data set to get relevant information about the student performance in the educational industry. This problem is being sorting out with the help of data mining technique using uniform district information system for education (UDISE).

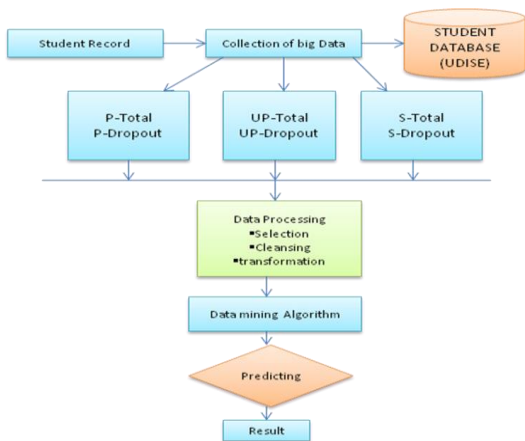


Fig.1 Prediction of student performance using UDISE

## A. Data collection and Description

Table 1. Student database using UDISE

2019-20	P(I-V)Total	P-Dropout	UP(VI-VIII)Total	UP-Dropout	S(IX-X)Total	S-Dropout
Location						
BALOD	64851	0.01	39257	0.83	27857	15.68
BALODABAZAR	148282	0.95	85115	4.15	55936	19.73
BALRAMPUR	91517	2.98	44849	8.16	25516	18.44
BASTER	84203	4.9	42286	5.55	26443	19.45
BEMETARA	89485	0.41	53571	3.32	34811	24.77
BIJAPUR	33402	11.21	11661	12.27	6300	16.57
BILASPUR	226018	0	125801	4.46	83763	19.09
DANTEWADA	30943	9.29	13166	3.76	7624	19.96
DHMTARI	69774	0.71	41245	1.47	30043	19.46
DURG	146365	0	86985	1.52	57084	14.57
GARIABAND	59143	1.9	33011	6.76	19457	17.22
JANJIR - CHAMPA	173628	0.47	100845	1.91	67052	17.53

For our work we have presented sample datasets of year 2019-2020 from UDISE. Sample of this data set present here total number of student enrol\_ in school has been taken. Students are divided in 3 categories they are primary, upper primary and secondary level. In this section total no. of enrol\_ student and dropout percentage of students are considered. Here, location and total number of students are independent variables and dropout percentage of students is dependent variable. By using out location were dropout percentage of student is minimum.

## B. Feature extraction and feature selection

Table 2 .Show year 2019-20 statistical representation

Statistic	P(1-V)Total	P-Dropout	Missing 1(3%)	
Minimum	17412	0		
maximum	2651484	11.21	P(1-V)Total	P-Dropout
Mean	189391.7	2.292	Distinct 28	Distinct 22
stdDev	485407.1	3.429	Unique 28(97%)	Unique 20(69%)

Statistic	UP(VI-VIII)Total	UP-Dropout	UP-Total ,dropout
Minimum	7817	0.83	Missing 1(3%)
Maximum	1481381	12.27	Distinct 28
Mean	105812.9	4.592	Unique

# IJFANS INTERNATIONAL JOURNAL OF FOOD AND NUTRITIONAL SCIENCES

ISSN PRINT 2319 1775 Online 2320 7876

Research Paper © 2012 IJFANS. All Rights Reserved. Journal Volume 11, Iss 11, 2022

			28(97%)
StdDev	271452.2	2.513	

Statistic	S(IX-X)Total	S-Dropout	S-total, dropout
Minimum	4548	12.94	Missing 1(3%)
Maximum	943808	24.77	Distinct 28
Mean	67414.86	18.058	Unique 28(97%)
StdDev	17303037	2.333	

**Table 3 .Show year 2018-19 statistical representation (2018-19)**

Statistic	P(1-V)Total	P-Dropout	Missing 1(3%)	
Minimum	17412	0		
maximum	2651484	11.21	P(1-V)Total	P-Dropout
Mean	189391.7	2.292	Distinct 28	Distinct 27
StdDev	485407.1	3.429	Unique 28(97%)	Unique 2(90%)

Statistic	UP(VI-VIII)Total	UP-Dropout	
Minimum	7805	1.19	Missing 1(3%)
Maximum	1517322	19.22	Distinct 28
Mean	108380.1	0.104	Unique 28(97%)
StdDev	278006	3.816	

Statistic	S(IX-X)Total	S-Dropout	
Minimum	4289	14.72	Missing 1(3%)
Maximum	946427	31.73	Distinct 28
Mean	67601.93	20.695	Unique 28(97%)
StdDev	173494.9	3.78	

**Table 4 .Show year 2017-18 statistical representation (2017-2018)**

Statistic	P(1-V)Total	P-Dropout	P(1-V)Total	P(1-V)Dropout
Minimum	19543	0.005	Missing 1(3%)	
maximum	2671372	4.28	Distinct 26	Distinct 28
Mean	190812.3	1.831	Unique 24(83%)	Unique 28(97%)
stdDev	488913.6	1.116		

Statistic	UP(VI-VIII)Total	UP-Dropout	UP-total, dropout
Minimum	8267	0.09	missing1 (3%)
Maximum	1631532	8.35	Distinct 26
Mean	116538	4.819	Unique 24(83%)
StdDev	298862	2.037	

Statistic	S(IX-X)Total	S-Dropout	S-total, dropout
Minimum	4317	13.44	Missing1(3%)
Maximum	947230	28.87	Distinct 28
Mean	67659.29	19.966	Unique 28(97%)
StdDev	173589	3.566	

**Table 5.Show year 2016-17 statistical representation (2016-2017)**

Statistic	P(1-V)Total	P-Dropout		
Minimum	18761	0	Missing 1(3%)	
maximum	2710696	12.88	P(1-V)Total	P-Dropout
Mean	193621.143	1.855	Distinct 28	Distinct 18
stdDev	496260.866	3.73	Unique 28(97%)	Unique 17(59%)

Statistic	S(IX-X)Total	S-Dropout	S-total,dropout
Minimum	4078	14.78	Missing 1(3%)
Maximum	943703	30.24	Distinct 28
Mean	67407.351	21.066	Unique 28(97%)
StdDev	172971.952	3.824	

Here table 2,3,4 and 5 show statistical representation of primary (I-V) total, dropout and upper primary (VI-VIII) total, dropout and secondary (IX-X) total, dropout values of 3 years student dataset, and show minimum, maximum, mean and stdDev values on each class. And also given the value of missing value distinct value or given unique value. This value is accessed by running data on weka platform.

**C. Proposed methodology**

Research has its special significance in solving various operational problems. Research methodology is the way to systematically solve the research problem.

**Hypothesis of Research Work**

- Hypothesis<sub>1</sub>: Is there no relation between transition rate, promotion rate and dropout rate?
- Hypothesis<sub>2</sub>: Are data gaps deteriorating data quality?
- Hypothesis<sub>3</sub>: Did the model successfully classify the U-DISE data or not, in terms of location wise enrolment and learning performance?

**D. Algorithms used for Classification**

During the intense study of around few contributions, various architecture of machine learning model has been Studied. In most of the contribution authors have suggested different models of machine learning suitable for Student performance prediction. The major contributions are as follows.

- **LWL**

Local weight learning is approximation technique. It is find the underlying relationship between input and output. When we use dataset or Training data were each input is associated with one output and its use to create model that predicts values which come and close to the correct/true function.LWL use local functions and create a local model.

- **Random Forest**

Random forest is a supervised machine learning algorithm. It can be use for both regression and classification problem solving schemes used in machine learning. It follows the

Statistic	UP(VI-VIII)Total	UP-Dropout	Missing 1(3%)	
Minimum	8127	0	UP(1-V)Total	Dropout
maximum	1639555	17.52	Distinct 28	Distinct 26
Mean	117111.071	5.793	Unique 28(97%)	Unique 25(86%)
stdDev	300346.181	3.8		

concept of ensemble learning algorithm which is the combination of multiple classifiers and solves the difficult problem with a great accuracy and also improves the model performance. It is use in Banking, Medicine, Land use or Marketing. Random forest contains a number of decision tress on various subsets of given dataset and predicting the majority of higher voting. It works with two phase first it create the random forest by combining N decision tree and second phase is to make predictions for each tree created in the phase.

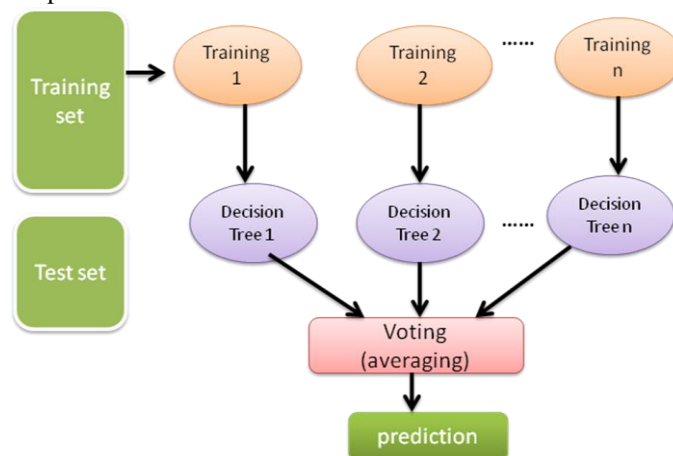


Fig.2 The above diagram explains the working of Random forest.

- **Bagging**

Bagging is also known as bootstrap aggregation, is the ensemble learning technique, which is generally use to improve the stability and accuracy of machine learning algorithms and reduces variance within a noisy dataset. It is help to avoid over fitting and it can be use in different type of method like regression or classification specially decision tree method.



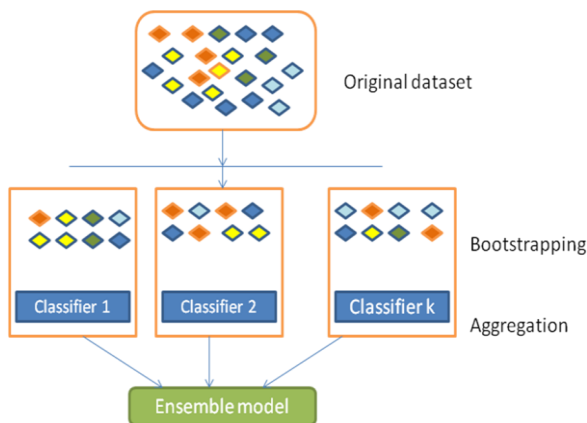


Fig. 3 This diagram presents the Bagging process.

### • Ensemble

Ensemble is a machine learning algorithm that combines more than two or more models and makes one optimal predictive model. There is Bagging, Staking and Boosting are three main classes of ensemble technique. It produces more accurate solution compare to a single model would.

### E. Performance Evaluative method

In our work we have used Weka tool for evaluating the dataset.

(TP) and true negative (TN), respectively; the false positive (FP) and false-negative (FN) denotes the misclassification of normal and infected images, respectively;  $P = TP + FN$  and  $N = TN + FP$ .

(TP) and true negative (TN), respectively; the false positive (FP) and false-negative (FN) denotes the misclassification of normal and infected images, respectively;  $P = TP + FN$  and  $N = TN + FP$ .

Accuracy (ACC) =  $(TP + TN) / (P + N) \times 100$

Specificity =  $TN / N \times 100$

Precision =  $TP / (TP + FP) \times 100$

Recall =  $TP / P \times 100$

F1 – Measure =  $(2 \times \text{precision} \times \text{recall} / (\text{precision} + \text{recall})) \times 100$

Area under Curve (AUC) =  $1/2 (TP / P + TN / N)$

Matthews Correlation Coefficient (MCC) =  $(TP \times TN - FP \times FN) / \sqrt{(TP + FP) \times P \times N \times (TN + FN)}$

Finally, the obtained result is statistically validated using z-test and Friedman average ranking and Holm (Holm, 1979) and Shaffer (Shaffer, 1986) post-hoc multiple comparison methods.

Table The Prototype of the proposed automatic sequential Model ((Abbreviations: LWL, Random Forest, Bagging, Ensemble)

### III. Experimental result and discussion

Instruction	Algorithm	2019-2019-droptout				2019-2019-droptout				2019-2019-droptout				
		use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set
	LWL	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
	Random forest	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
	META	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
	Bagging	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01

Algorithm	(2019-2019-Total)				(2019-2019-total)				(2019-2019-totat)				
	use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set
LWL	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Random forest	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
META	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Bagging	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01

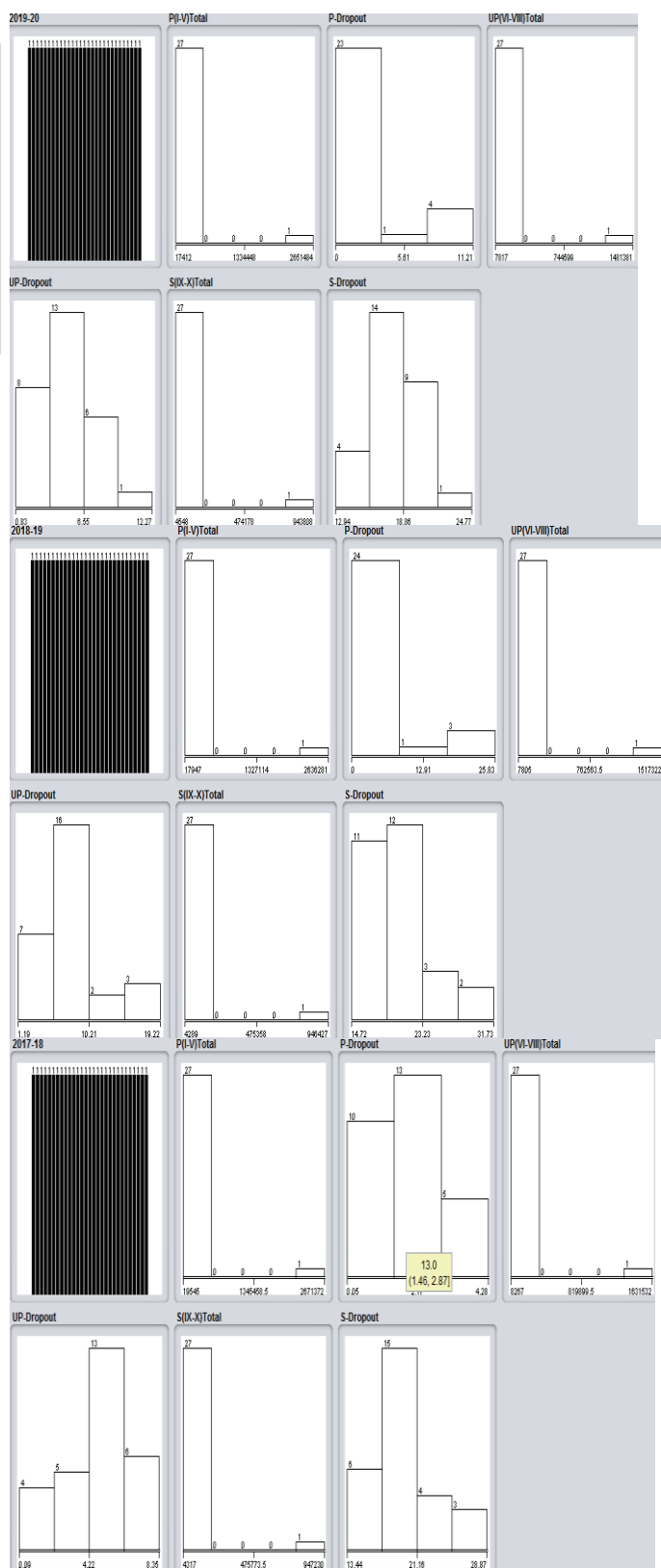
Table 6.Represented data in the year 2019-20 total and dropout value by running weka platform.

Classification/Algorithm	(2018-19SP-droptout)				(2018-19sp-droptout)				(2018-19)-droptout				
	use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set
LWL	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Random forest	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
META	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Bagging	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01

Classification/Algorithm	(2018-19SP-TOTAL)				(2018-19sp-TOTAL)				(2018-19)-TOTAL)				
	use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set	supplied test set	cross validation(10 split(66%))	percentage use training set
LWL	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Random forest	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
META	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01
Bagging	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01	0.01

Table 7.Represented data in the year 2018-19 total and dropout value by running weka platform.

Algorithm	[2017-18]P-dropout				[2017-18]UP-dropout				[2017-18]S-dropout			
	use training set	supplied test set	cross validation(10 fold)	percentage split(66%)	use training set	supplied test set	cross validation(10 fold)	percentage split(66%)	use training set	supplied test set	cross validation(10 fold)	percentage split
LWL	Time taken to test model	0.03	0.01	0	0	0.01	0.01	0	0	0.01	0.01	0.02
	Correlation coefficient	0.7842	0.7842	-0.1137	0.2463	0.7581	0.7581	-0.1897	0.8849	0.7581	0.7581	-0.1897
	Mean absolute error	0.543	0.543	1.0079	1.2621	0.9979	0.9979	1.5727	2.3155	0.9979	0.9979	1.5727
	Root mean squared error	0.8092	0.8092	1.2113	1.3938	1.2465	1.2465	2.0907	2.8181	1.2465	1.2465	2.0907
	Relative absolute error	62.70%	62.70%	108.61%	94.12%	66.26%	66.26%	101.76%	135.21%	66%	66%	101.76%
	Root relative squared error	63.81%	63.81%	105.74%	94.69%	66.22%	66.22%	100.61%	131.71%	66.22%	66.22%	100.61%
	Total Number of Instances	28	28	28	28	28	28	28	28	28	28	28
Ignored Class Unknown Instances	1	1	1	1	1	1	1	1	1	1	1	
Random forest	Time taken	0.01	0.02	0.09	0	0	0.01	0.02	0	0	0.02	0.06
	Correlation coefficient	0.8646	0.8646	-0.1102	0.1581	0.9935	0.9935	0.176	0.5849	0.9932	0.9932	0.248
	Mean absolute error	0.2944	0.2944	0.9381	1.2811	0.5651	0.5651	1.5312	1.84	0.566	0.566	2.5671
	Root mean squared error	0.3875	0.3875	1.2177	1.4436	0.7562	0.7562	1.9766	2.0952	0.7562	0.7562	2.762
	Relative absolute error	34.00%	34.00%	99.11%	95.68%	37.52%	37.52%	99.07%	95.76%	36.22%	36.22%	94.60%
	Root relative squared error	35.17%	35.17%	97.69%	98.00%	37.81%	37.81%	95.61%	97.98%	35.54%	35.54%	95.00%
	Total Number of Instances	28	28	28	28	28	28	28	28	28	28	28
Ignored Class Unknown Instances	1	1	1	1	1	1	1	1	1	1	1	
Bagging	Time taken	0	0	0.02	0	0	0.02	0	0	0	0	0
	Correlation coefficient	0.97	0.97	-0.6628	0	0.9735	0.9735	-0.4011	0	0.9411	0.9411	-0.2581
	Mean absolute error	0.1359	0.1359	0.9139	1.3297	0.6004	0.6004	1.547	1.9016	0.5443	0.5443	2.7146
	Root mean squared error	0.4283	0.4283	1.4657	1.8557	0.8756	0.8756	2.029	2.3957	1.4773	1.4773	3.5441
	Relative absolute error	36.47%	36.47%	98.66%	98.89%	39.86%	39.86%	100.09%	111.04%	35.78%	35.78%	100.04%
	Root relative squared error	39.10%	39.10%	98.38%	98.89%	43.78%	43.78%	99%	112.14%	42.19%	42.19%	112.69%
	Total Number of Instances	28	28	28	28	28	28	28	28	28	28	28
Ignored Class Unknown Instances	1	1	1	1	1	1	1	1	1	1	1	



Algorithm	[2017-18]P-TOTAL				[2017-18]UP-TOTAL				[2017-18]S-TOTAL			
	use training set	supplied test set	cross validation(10 fold)	percentage split(66%)	use training set	supplied test set	cross validation(10 fold)	percentage split(66%)	use training set	supplied test set	cross validation(10 fold)	percentage split
LWL	Time taken to test model	0.01	0.02	0	0	0.02	0.03	0	0	0.02	0.02	0.02
	Correlation coefficient	0.999	0.999	0.8857	0.8857	0.9988	0.9988	0.8827	0.9987	0.9987	0.9987	0.8827
	Mean absolute error	18294.1553	18294.1553	12362.5969	19413.0568	11459.3151	11459.3151	49648.9016	13131.7943	4796.8239	39101.9424	7556.1314
	Root mean squared error	23164.1598	23164.1598	472545.9015	25618.5871	14370.7987	14370.7987	288059.2839	16472.8841	8607.1068	168646.2167	8811.3599
	Relative absolute error	10.03%	10.03%	60.36%	34.43%	10.30%	10.30%	60.84%	16.12%	11%	60.24%	18.86%
	Root relative squared error	4.56%	4.56%	44.44%	17.72%	4.90%	4.90%	94.76%	18.69%	5.05%	5.05%	17.87%
	Total Number of Instances	28	28	28	28	28	28	28	28	28	28	28
Ignored Class Unknown Instances	1	1	1	1	1	1	1	1	1	1	1	
Random forest	Time taken	0	0	0.01	0	0	0	0	0	0	0	0
	Correlation coefficient	1	1	0.3053	0	1	0.2242	0	1	1	0.2268	0
	Mean absolute error	0	0	18453.8027	40327.9412	0	0	98161.1295	26125.5471	0	0	57943.4311
	Root mean squared error	0	0	478610.0901	23704.9871	0	0	288394.024	34186.5418	0	0	167486.0159
	Relative absolute error	0.00%	0.00%	98.62%	29.98%	0.00%	0.00%	85.75%	32.32%	0.00%	0.00%	86.10%
	Root relative squared error	0.00%	0.00%	98.25%	32.14%	0.00%	0.00%	84.87%	30.79%	0.00%	0.00%	84.84%
	Total Number of Instances	28	28	28	28	28	28	28	28	28	28	28
Ignored Class Unknown Instances	1	1	1	1	1	1	1	1	1	1	1	
Bagging	Time taken	0.02	0	0.02	0	0	0.02	0	0	0.9933	0	0.02
	Correlation coefficient	0.9938	0.9938	-0.5206	38178.8934	0.9935	0.9935	-0.5162	0	0.9933	0.9933	-0.5101
	Mean absolute error	73956.0988	73956.0988	198789.4829	172108.2674	44940.7407	44940.7407	121272.6126	99578.9916	26119.327	26119.3268	70461.2722
	Root mean squared error	265418.3095	265418.3095	499938.3753	121.731	149961.1161	149961.1161	305423.744	305111.1217	47038.140	47038.1476	177254.3125
	Relative absolute error	40.67%	40.67%	108.25%	119.03%	46.18%	46.18%	105.14%	112.39%	40.13%	40.13%	112.44%
	Root relative squared error	51.12%	51.12%	100.47%	10.00%	51.08%	51.08%	119.17%	51.07%	51.07%	100.47%	110.84%
	Total Number of Instances	28	28	28	28	28	28	28	28	28	28	28
Ignored Class Unknown Instances	1	1	1	1	1	1	1	1	1	1	1	

**Table 8.** Represented data in the year 2017-18 total and dropout value by running weka platform.

Table 6,7 and 8 show here three years big data classify represented here from 2019-2020, 2018-2019, 2017-2018 and three data mining techniques LAZY, TREE, META present 3 data mining algorithm like LWL, Random forest, Bagging. They process primary total, dropout, upper primary total dropout and secondary total dropout data in four test option like use training data set, supplied test set, 10 fold cross validation and 66% percentage split. After running this test option they all show different results as time taken to test model, correlation coefficient, mean absolute error, Root mean squared error, Relative absolute error, Root relative squared error, Total Number of Instances, Total Number of Instances.

### A. Analysis And Discussion

Given below are the result based on Statistical Parameters of the Year 2019 – 2020, 2018-2019, 2017-2018, 2016-2017.

Fig.4 Graphical representation on P-total, dropout, UP- total, dropout and S-total, dropout three year dataset Year 2019 – 2020, 2018-19, 2017-2018, 2016-2017.

## IV. CONCLUSION

Educational data mining play an important role in higher education system, the use of rising technology need to largest dataset. With the help of U-DISE (unified district information system for education) we get overall type of big data as related to school information like students, faculty members and dropout student etc. it is government authorized nationalized platform, which is use in data mining concept to predict student performance. By this study provide improvement in public and private field to improve academic performance. It can be concluded that classification techniques provide better accuracy compared to other approach.

## References

- [1] Zambrano J.L., Lara Torralbo J.A., and Romero C., "Early Prediction of Student Learning Performance Through Data Mining: A Systematic Review", ISSN 0214 - 9915 CODEN PSOTEG , Vol. 33, No. 3, 456-465 doi: 10.7334 /psicothema 2021.62 (2021)
- [2] Viswanathan S. , Vengatesh S.K. , "Study Of Students' Performance Prediction Models Using Machine Learning", Turkish Journal of Computer and Mathematics Education , Vol.12 No.2 (year-2021), 3085 – 3091
- [3] Ingale N.V., Dr. Sivakkumar M., Dr. Namdeo V., "Survey on Prediction System for Student Academic Performance using Educational Data Mining", Turkish Journal of Computer and Mathematics Education Vol.12 No.13 (2021), 363-369
- [4] Mangat P. and Saini K., "Educational Data Mining Tools and Framework for Predicting Students Academic Performance", International Journal of Advance Science and Technology Vol. 29, No. 10S, (2020), pp. 2525-2533.
- [5] Zhou N., Zhang Z, Li J., "Analysis on Course Scores of Learners of Online Teaching Platforms Based on Data Mining", Ingenierie des Systemes d'information vol 25, No.5, pp.609-617. October 2020
- [6] Ahmad N.B.B, Ajibade M.S, Shamsuddin M.S., "Educational data mining : Enhancement of Student Performance model using Ensemble methods" IOP Conf. Series: Materials Science and Engineering 551 (2019) 012061, doi:10.1088/1757-899X/551/1/012061
- [7] Hossain Z., Akhtar N., Ahmad R.B., Rahman M., "A dynamic K-means clustering for data mining". Indonesian Journal of Electrical Engineering and Computer Science. Vol. 13, No. 2, pp. 521~526 ISSN: 2502-4752, February 2019.
- [8] Ali N., Neagu D., Trundle P., "Evaluation of k-nearest neighbour classifier performance for heterogeneous data sets", SN Applied Sciences (2019) 1:1559.
- [9] Sankari A., Masih S., Ingle M., "A Review On Research Areas In Educational Data Mining And Learning Analytics", International Journal of scientific & technology research volume 8, ISSUE 12, December 2019 ISSN 2277-8616.
- [10] Yaacob W.F.W., Nasir S.A.M., Yaacob W .F. W., Sobri N.M., "Supervised data mining approach for predicting student performance" Indonesian Journal of Electrical Engineering and Computer Science Vol. 16, No. 3, December 2019, pp. 1584~1592 ISSN: 2502-4752, DOI: 10.11591/ijeecs.v16.i3.pp1584-1592.
- [11] Majeed I, Naaz S., "Current State of Art of Academic Data Mining and Future Vision", Indian Journal of Computer Science and Engineering". e- ISSN: 0976-5166, p-ISSN: 2231-3850 April 2018.
- [12] Qiao X. and Jiao H., "Data Mining Techniques in Analyzing Process Data: A Didactic", original research published: 23 November 2018, doi: 10.3389/fpsyg.2018.02231.
- [13] More SS, Narain B, Jadhav BT (2017) A comparative analysis of unimodal and multimodal biometric systems. In: International conference on innovative trends in engineering science and management (ITESM-2017).
- [14] Narain B, Zadgaonkar AS, Kumar S (2013) Impact of digital image processing on research and education. Natl Semin Work.
- [15] More S.S., Narain B., Jadhav B.T. (2021) Advanced Encryption Standard Algorithm in Multimodal Biometric Image. In: Rizvanov A.A., Singh B.K., Ganasala P. (eds) Advances in Biomedical Engineering and Technology. Lecture Notes in Bioengineering. Springer, Singapore. [https://doi.org/10.1007/978-981-15-6329-4\\_7](https://doi.org/10.1007/978-981-15-6329-4_7)
- [16] Data encryption standard algorithm in multimodal biometric image, SS More, B Narain, BT Jadhav - Int. J. Comp. Sci. Eng. 2018.
- [17] Hybrid Support Vector Machine and Distance Classifier in Breast Tumor Detection U Sharma, B Narain, V Nohria - SPAST Abstracts, 2021
- [18] Student Satisfaction in Educational Organization using Machine Learning, P Singh, B Narain - NOVYI MIR, 2021.
- [19] Impact of digital image processing on research and education B Narain, AS Zadgaonkar, S Kumar Natl Semin Work, 2013
- [20] The Online Retail Market Analysis for Social Development with Machine Learning M Nayak, B Narain - SPAST Abstracts, 2021.
- [21] On-Line Big Data Analysis using K-Mean and Modified K-Mean Algorithm with Machine Learning Techniques\* B Narain 2021.
- [22] Big Data Mining Algorithms for Predicting Dynamic Product Price by Online Analysis\* M Nayak, B Narain - Computational Intelligence in Data Mining, 2020.