

RECOMMENDING COLLEGE PROGRAMS TO STUDENTS USING MACHINE LEARNING

AISHA GHAZAL FATEH ALLAH, GHAZALA BILQUISE

Lecturer, Higher Colleges of Technology, Department of Computer Information Science, U.A.E

Lecturer, Higher Colleges of Technology, Department of Computer Information Science, U.A.E

E-mail: aisha.ghazal@gmail.com, gbilquise@hct.ac.ae

ABSTRACT

The decision to choose a program of study is a major one since a student must commit to it for four years, until graduation. Hence it is a crucial decision for academic as well as future career success. Despite this, students often make academic choices without careful thought mainly due to lack of proper advice and support. In this paper, we use four methods that utilize students' data such as their performance in high school, college placement test, and standardized IELTS exam to recommend a college program as well as predict the students GPA in those programs. Using the four methods utilizes the advantages of each of them and provides insight into the reason for the recommendation. The four methods that we used, evaluated, and compared are Decision Trees (DT), Neural Network (NN), K-Nearest neighbor (KNN), and Linear Regression (LR). To the best of our knowledge, this is the first study that utilizes and compares the different approaches.

Keywords: *Machine Learning, Classification, Decision Tree, Neural Network, k-NN, Linear Regression, Collaborative Filtering, Recommender Systems,*

1. INTRODUCTION

Several factors determine a student's academic performance. One of those factors that hinder students' success is selecting a program of study that does not match their abilities and skills. Often students choose their program of study based on the current trends in addition to peers or family influence. This can affect their performance negatively and increase the likelihood of drop out within the first year. New students could benefit from the recommendation of programs based on their abilities, skills, and performance. Al-Lawati et al. [1] stated that the student's choice of specialization is essential to their academic as well as future success. Hence it is essential to guide students and facilitate their decision. In an empirical study conducted by Al-Lawati et al. [1], liking and preference of the specialization and peer influence were the top reported factors that influence the choice of majors among business students. Ray and Sharma [2] also stated that due to lack of awareness students often rely on peer recommendations to make academic choices.

Kusumaningrum et al. [3] highlighted student failure as a significant challenge faced by many universities in Indonesia, where the number of students who graduate on time is only 52%. They

attributed inappropriate selection as one of the causes for the low percentage of students who fail to complete their degree on time or drop out. Deorah et al. [4] suggest that students often choose their program of study based on trends, which may eventually lead to failure as their choice may not match their abilities and skills.

In the college of study, 30% of students did not graduate on time in the year of 2016 [1]. Hence, on-time graduation is one of the KPIs to meet the strategic goal of Enhancing student success. To ensure the success of the student, we use, evaluate, and compare four different approaches to recommend a program to the students, namely, Decision Trees (DT), Neural Network (NN), K-Nearest neighbor (KNN), and Linear Regression (LR). The algorithms utilize high-school data and college acceptance tests to recommend a program and predict the GPA to help students make an informed decision. We could not find any other research that did the same. The rest of the paper is organized as follows.

Section 2 presents related work in the use of recommender systems in higher education for making academic choices, section 3 explores data and describes the preprocessing tasks, section 4 describes the methodology and algorithms used,

section 5 explains our evaluation method, section 6 discusses the results of the study. Finally, in section 7 we present our conclusion and future work.

2. RELATED WORK

Recommendation systems are employed in various domains of educational institutions such as recommending learning resources, making academic choices, advising, and planning, to name a few [9]. Our study is focused on the use of recommendation system for supporting students while making an academic choice, namely, choosing a program of study. Literature review in this domain has shown that several studies used recommendation systems to personalize courses choices ([5], [8], [2], [6], [7], while other studies have recommended specialization [3], [4]) and postgraduate universities [10], however, they have a different approach than in this study.

Stein et al. [11] recommends majors to college students based on their performance in courses during the first few years of college. The study used collaborative filtering using nearest neighbour. However, students in many universities, including the one in our study, choose majors at enrolment and not after a few years. Hence, this approach cannot be generalized. Our paper for instance bases the recommendation on pre-college data.

Park [12] in the paper titled “A Recommender System for Personalized Exploration of Majors, Minors, and Concentrations” proposes a recommender system using collaborative filtering to predict students’ grade, however, the author has not actually used any data to test or evaluate the recommended approach. Our system is based on real data and the system was actually built, tested, and evaluated.

Alshaikh et al. [13] uses collaborative filtering to recommend a specialization to students, however, they only use a single algorithm, namely, nearest neighbor.

Mostafa et al. [14] illustrated how a recommender system could be applied to automate the otherwise manual academic advising process. They use a case-based reasoning system to recommend the most suitable major to students. The concepts studied in each course, as well as the concepts that the students will study in the chosen major, are used to determine the similarity with the cases. However, it is not an easy task to extract

features of courses and match them with the majors. The study by Kusumaningrum et al. [3] investigates the use of association rule to recommend academic majors to students based on their academic history, profile data, as well as their preference. While [4] used case-based reasoning using academic history as well as implicit and explicit interests of students to recommend a suitable major out of 15 available choices. Explicit preferences were gathered using surveys while implicit preferences were gathered by computing the time spent by the student to respond to each question. The study used student input to narrow down the number of generated recommendations.

Hasan et al. [10] investigated the use of collaborative filtering to recommend a graduate school based on the student’s undergraduate program, TOEFL/IELTS score, cumulative GPA, research interest, and profile data. Students are requested to provide all this information to get a recommendation. Several papers use a recommendation system to recommend suitable courses. Engin et al. [5] proposed the use of a rule-based system to recommend courses to students. The purpose of the system is to ease the load on faculty members as well as to provide valid and suitable recommendations. Rules were based on students GPA, pre-requisites, offered courses, specialization of the students, and courses already taken.

Ray and Sharma [2], [6] and [7] applied collaborative filtering technique to recommend courses to students by predicting course grades. Ray and Sharma [2] employed user-based as well as itembased collaborative filtering techniques to predict grades of elective courses and generate a list of recommendations, while [7] used item-based recommendation technique to predict the final score in the elective course. A Root Mean Square Error of 0.5 was reported which indicates a good performance since the predicted values ranged from 0-10.

Ng and Linn [8] did not use any historical data stored in the system to make recommendations for courses. Instead, all the data was gathered using surveys to build a user profile. Also, sentiment information was retrieved from course reviews and teacher feedback. Matrix factorization was used to predict students rating of a given course which was then used to make a suggestion.

3. DATA EXPLORATION AND PREPROCESSING

The original dataset consists of 7,074 students’ enrollment records in the year 2015-2016.

Each record is described by 71 features which include personal data as well as pre-college performance data. The dataset includes students registered in three college programs (Business, IT and Engineering). We excluded year one students since they would not have achieved enough credits to make the GPA trustworthy enough to make valid recommendations. After applying the preprocessing tasks, described in 3.2, the dataset size reduced to 1,892 records out of which 1,052 represents business program students, 339 represented IT students and 501 represents engineering students. Table 1 shows a description of the main features of the dataset that were used in our recommendation algorithm.

Table 1: Main Features of Dataset.

Feature	Description	Values
Coll	Program of study(college)	Business, IT, Engineering.
CGPA	Cumulative GPA	Value between 0 and 4.
HS-AVG	High School Average	A numeric score between 0-100.
CEPA	CEPA English Score	A numeric score between 140 -210.
CEPA-MATH	CEPA Math Score	A numeric score between 140 -210.
IELTS-BAND	IELTS Score	A numeric value between 0-9.

3.1 Descriptive Statistics

To better understand the data, we explored the basics statistics of all the numeric attributes. The minimum and maximum values, average, standard deviation, in addition to a histogram of the data are shown in Figure 1.



Figure 1: Summary Statistics

3.2 Data Preprocessing

After exploring the dataset, we realized that it is noisy and contains missing values. Therefore, we performed preprocessing tasks to smoothen the data and resolve the inconsistencies. The following is a list of all the preprocessing tasks:

- Anonymized the dataset by removing over 20 attributes that contained personal details of students such as Student ID, Student Name, ID, contact numbers, and addresses.

- Removed records with data entry errors
- Removed records with missing values

We performed the following two additional preprocessing to extract only the required data that is relevant to our study.

- Removed records of first-year students.
- Filtered students for each program (Business, IT and Engineering)

Figure 2 shows the preprocessing tasks in RapidMiner.

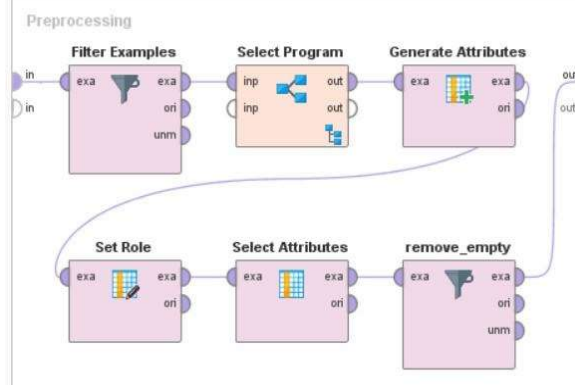


Figure 2: Preprocessing in Rapid Miner

4. METHODOLOGY

We investigated the use of four approaches to recommend programs of study to newly enrolled students to increase their chance of academic success. The four approaches are: Decision Trees (DT), Neural Network (NN), K-Nearest neighbor (KNN), and Linear Regression (LR). The algorithms utilize high-school data and college acceptance tests to recommend a program and predict the student's GPA in the program.

The Three programs of study are Engineering, Business and Information Technology (IT). The student profile is built using their pre-college performance which includes high school average, the college placement tests for English and Math as well as the IELTS band. The Cumulative GPA (CGPA) is used as the predictor to make recommendations. To make valid recommendations, we have excluded year one students since they would not have achieved enough course credits to make the Cumulative GPA trustworthy.

The resultant number of records in the dataset for each program is shown in Table 2.

The four approaches are described in the coming sections

Table 2: Number of records in each program.

Program	Number of Records
Business	1,052
Information Technology	339
Engineering	501

4.1 Decision Tree, Neural Network, and k-Nearest Neighbor

We tested three machine-learning algorithms in our system: Decision tree, Neural Networks, and k-NN (k-Nearest Neighbor). To make a recommendation, the following approach was used:

1. A new attribute called Recommendation was generated with values of Recommend or Do not Recommend as a class label to determine whether the program of study is a good choice for the student or not. If the student's CGPA is 2.7 or more, we assume the program was of a good choice and would recommend it to students with a similar profile. We used 2.7 as a cut off CGPA after consulting with academic advisors.

Figure 3 shows the condition that was used in RapidMiner, and Figure 4 shows a sample of the generated data.

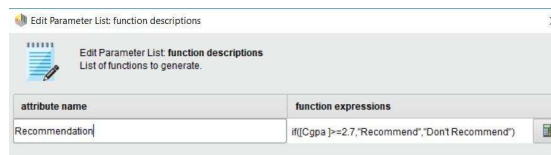


Figure 3: Condition in Rapid Miner

2. We excluded the CGPA from the dataset because we do not want the machine learning algorithm to use it for the recommendation because new students will not have a value for this. Our recommendation is based on pre-college data only.

3. We trained each machine-learning algorithm on one college program at a time by filtering the dataset by programs. Figure 5 shows how to alternate between different programs in RapidMiner. The machine-learning algorithm learns the pattern of historical grades that lead to success in the program from the training set and classify the testing set accordingly.

4. We evaluated the performance of the recommendation for each program and each algorithm.

Figure 6 shows the overall processes in RapidMiner. The three machine-learning classifiers (Decision Tree, Artificial Neural Network and k-Nearest Neighbor) used for recommendation are described in the coming sections.

ExampleSet (1052 examples, 1 special attribute, 5 regular attributes)

Row No.	High School ...	CEPA	CEPA MATH	IELTS Band	Cgpa	Recommendation
1	73	168	150	5	2.800	Recommend
2	71.500	170	154	5	3.370	Recommend
3	78	171	143	5	2.290	Don't Recommend
4	51.400	173	145	5	2.680	Don't Recommend
5	83.700	168	151	5.500	3.180	Recommend
6	68.800	159	156	5	2.380	Don't Recommend

Figure 4: Sample Data

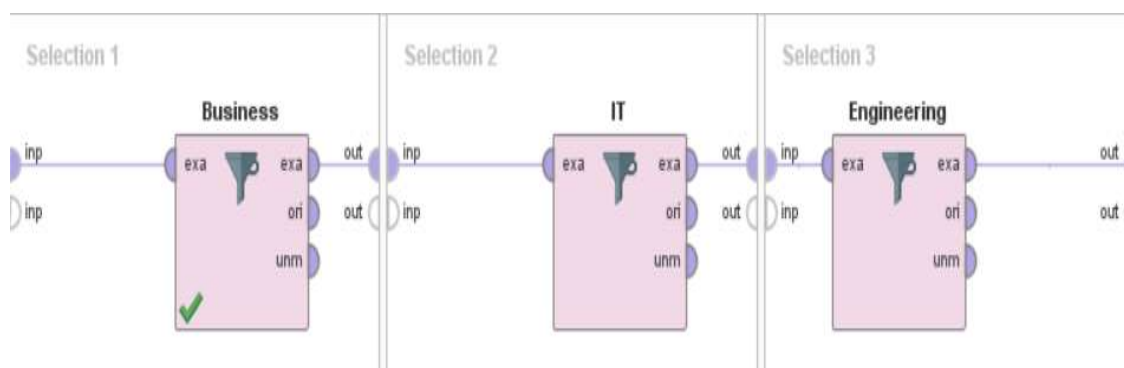


Figure 5: Rapid Miner Select Process

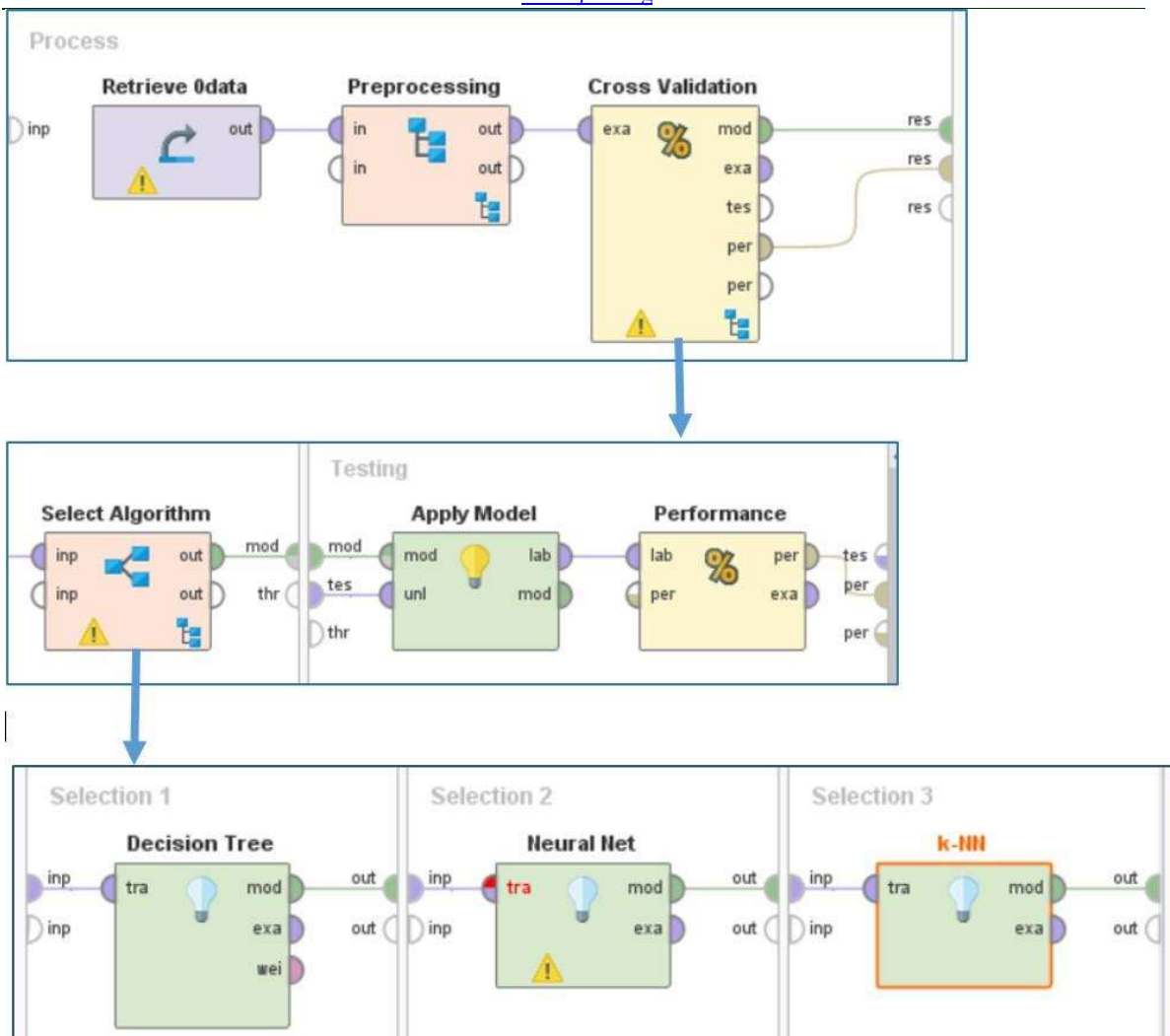


Figure 6: All Rapid Miner Processes

Row No.	High School ...	CEPA	CEPA MATH	IELTS Band	Cgpa	prediction(Cgpa)
1	55.300	154	162	5	3.090	2.494
2	88	207	169	6	3.130	3.396
3	77.800	166	158	5.500	2.900	3.046
4	73.700	191	161	6	3.040	3.049
5	75.300	174	140	5	1.900	2.833
6	93.800	173	152	5.500	3.520	3.369
7	81.100	172	148	5	2.240	2.995
8	74.300	159	133	5	2.830	2.781

Figure 7: Sample Prediction

4.1.1 Decision Tree Algorithm

According to Lior [15], a major drawback of decision trees in recommender systems is the need to create a tree for every item (i.e., program) that is being recommended. Fortunately, in our scenario, this is not a problem since there are only three programs being recommended to students (Business, Engineering, and IT). Furthermore, the decision tree provides an explanation for the recommendation, which is very useful for our scenario.

We built a decision tree for each program, shown in Section 6, based on the training data of that program.

4.1.2 Artificial Neural Network Algorithm

Neural Network algorithm learns a model by using the input data, but it is like a black box, hence does not provide an explanation for the recommendation. However, we decided to use it and compare its performance to other recommendation algorithms. We will discuss the results in Section 6.

4.1.3 k-Nearest Neighbor(k-NN) Algorithm

k-NN is the most common approach to collaborative filtering [16]. We used it to match student record to k similar students using the Euclidean distance similarity measure to recommend a program. We experimented with different values of k and found the best performance to be produced at k=5.

4.2 Linear Regression

Linear regression is a commonly used technique for predictive analytics in machine-learning that is used to predict the outcome of a dependent variable. Simple linear regression uses one independent variable to predict the value of one dependent variable. Multiple linear regression uses several independent variables to predict the value of a dependent variable.

We used multiple linear regression to predict the CGPA that a student might get if he/she chooses a particular program of study. The CGPA (dependent variable) was predicted using four independent variables the high school average, college placement test (CEPA) score, CEPA math and IELTS Band.

We decided to add a prediction of student GPA to give students more insight into their expected performance in the different programs given their precollege data. We believe this will equip the

students to make a better choice using the predicted CGPA as well as their personal preference.

Data was split into 80% training data, and 20% testing data and prediction of GPA was applied to one program at a time. The linear regression algorithm generates a prediction formula using historical data and uses it to predict the GPA for unseen records. Figure 7 shows a sample of the output with the actual and predicted CGPA. The accuracy of the algorithm is discussed in the next section.

5. EVALUATION

In this section we discuss the evaluation method of recommending a program by classification as well as by linear regression.

5.1 Classification (Decision Tree, Neural Network, and k-NN)

Our study investigated recommending college programs to new students based on their high school average, college placement tests and IELTS results. For each program, we applied three classifiers with a total of nine rounds to determine the best machine learning algorithm for our recommender system.

We used three machine learning classifiers (Decision tree, k-NN, and neural network) to generate recommendations by matching similar users and classifying them accordingly to make a recommendation. For the k-NN classifier, we further experimented with different settings for the value of K. The best performance was achieved with k=5.

We used accuracy as a measure to evaluate the performance of the classifiers. Accuracy measures the total number of correct predictions from the total predictions made.

Accuracy may yield misleading results if the dataset is highly imbalanced. However, this is not the case in our scenario, since we used 2.7 CGPA as the cut-off requirement to make a recommendation, so dataset is only slightly imbalanced.

We used cross-validation to iterate through the dataset ten times by partitioning it into a training set and testing set with a ratio of 9:1 in each iteration. This can help generate a more accurate estimate of the prediction performance. We

recorded the accuracy of each classifier for each program: Business, IT and Engineering. The results of the performance are shown in Table 3

Table 3: Classification Performance Results.

	Decision Tree	Neural Network	k-NN
IT (339 records)	64.16%	63.16%	58.64%
Engineering (502 records)	62.47%	65.28%	67.67%
Business (1052 records)	69.49%	67.87%	67.41%

5.2 Linear Regression

We used linear regression algorithm for each program to predict the CGPA a student would achieve if he/she were to select that program, given the high school average, college placement tests, and IELTS band.

To evaluate the performance of linear regression we measured the error rate of the prediction using two evaluation metrics - Mean Absolute Error (MAE) and Root Mean Square Error (RMSE) since these are often used to determine the accuracy of predictions. MAE and RMSE measure the deviation of the recommendation from the actual CGPA. Lower values of MAE and RMSE imply higher prediction accuracy. We captured the prediction error rate for each program as shown in Table 4.

Table 4: Prediction Error Rate Results

	RMSE	MAE
IT (339 records)	0.537	0.425
Engineering (502 records)	0.503	0.412
Business (1052 records)	0.452	0.368

6. DISCUSSION OF RESULTS

In this section we discuss the results of the different approaches used in this study to recommend a program of study.

6.1 Classification (Decision Tree, Neural Network, k-NN)

Overall, the best accuracy of recommendation using was 69.49% achieved using the decision tree algorithm for the Business program as shown in Figure 8. The results show that the performance of the algorithms is proportional to the number of records in the training set, regardless of the classifier used. This emphasizes the importance of large data sets for machine-learning algorithms so they can better learn from that data and produce more accurate classification.

From the three machine-learning classifiers used, the only one that provides an explanation of the recommendation is the Decision Tree. Figure 9, Figure 10, and Figure 11 show the generated decision trees for Business, IT and Engineering programs respectively. The Decision Tree that was generated by the machine-learning algorithm revealed that English level (IELTS band) is the best indicator of student's performance in the Business program -hence appears at the root level of the recommendation tree.

On the other hand, Engineering and IT both have Math placement results (CEPA Math) at the root of the recommendation which can be explained by the scientific nature of these two programs as opposed to the Business program which requires higher English language competency.

6.2 Linear Regression

Figure 12 shows the average error in GPA prediction in all the programs. GPA range is between 0 and 4, and in average, the predicted GPA is around 0.425 far from the actual GPA. Overall, the lowest error rate achieved was 0.368 for predicting the student GPA in the Business program due to the larger data set-. Like accuracy, the results emphasize the importance of large data sets for machine learning algorithms since prediction error rates are inversely proportional to the number of records. The more the records, the more accurate the generated prediction formula is.

7. CONCLUSION AND FUTURE WORK

Choosing a college program after high school is a major decision that affects students for their whole duration of study in addition to their future career. Hence proper guidance is essential to make this decision. In this study, we used four algorithms to better support students and recommend a program of study to them, namely, Decision Trees (DT), Neural Network (NN), K-Nearest neighbor (KNN), and Linear Regression (LR). We discussed and compared the performance of the algorithms which, up to our knowledge, was not done in any prior research. The hope is to reduce the failure rates of students caused by the lack of guidance while selecting programs of study. The algorithms used students' pre-college performance data that includes high school average, CEPA English and CEPA Math placement tests, in addition to standardized IELTS band to recommend a program of study.

We also used linear regression algorithm for each program to predict the CGPA a student would achieve if he/she were to select that program.

We used Accuracy as a measure to evaluate the performance of the classifiers. The accuracy of a classifier is the number of correct predictions divided by the total number of predictions. The result is usually a value between 0% and 100%.

Accuracy=correct predictions/total predictions

Furthermore, we used ten-fold cross-validation to produce a more accurate estimate of the classification performance. Ten-fold cross validation runs for ten times, with each iteration being performed on a training set that comprises 90% of the total training set selected at random, and holding out the remaining 10% to be used for validation.

Overall, the best accuracy achieved was 69.49% using decision tree for the Business program. This shows that the task at hand is a hard one. The result also shows that accuracy is proportional to the number of records in the training set, regardless of the classifier used. Furthermore, decision tree model revealed that the IELTS band is the best indicator of students' performance in business while CEPA math scores are the best indicators for Engineering and IT programs.

To evaluate the performance of the linear regression technique we measured the error rate of the prediction using two evaluation metrics, namely

Mean Absolute Error (MAE), and Root Mean Squared Error (RMSE). Overall, the lowest error rate achieved was 0.368 for predicting the students GPA in the Business program. The predicted GPAs were in average within 0.425 points of the actual GPAs in the worst case.

As our results show, having more records improve the performance of the recommendation. Hence, for our future work, we would like to obtain a larger data set to produce better recommendations. We would also investigate using psychometric tests and/or surveys to improve the recommendations. Other recommendation techniques could be investigated as well, and we can extend recommendations to specializations within the different programs.

A major strength of this work is the use and comparison of multiple approaches to recommend programs. A suggested improvement to the system is to explore using different parameter settings and study the effect of that on the accuracy of the results. Another future improvement would be to integrate the system with the existing college banner so students can use it with ease. It would also be interesting to investigate the actual effect of our recommendation on students' success over the years.

REFERENCES:

- [1] Essam Hussain Al-Lawati, Radhakrishnan Subramaniam, et al. 2017. An empirical study on factors influencing business students choice of specialization with reference to nizwa college of technology, oman. *International Business Research*, 10(9):177.
- [2] Sanjog Ray and Anuj Sharma. 2011. A collaborative filtering based approach for recommending elective courses. In *International Conference on Information Intelligence, Systems, Technology and Management*, pages 330–339. Springer.
- [3] Desi Purwanti Kusumaningrum, Noor Ageng Setiyanto, Erwin Yudi Hidayat, and Khafiizh Hastuti. 2017. Recommendation system for major university determination based on students profile and interest. *Journal of Applied Intelligent System*, 2(1):21–28.
- [4] Sourabh Deorah, Srivatsan Sridharan, and Shivani Goel. 2010. Saes-expert system for advising academic major. In *Advance Computing Conference (IACC), 2010 IEEE 2nd International*, pages 331–336. IEEE.

- [5] Go'khan Engin, Burak Aksoyer, Melike Avdagic, Damla Bozanlı, Umutcan Hanay, Deniz Maden, and Gurdal Ertek. 2014. Rule-based expert systems for supporting university students. *Procedia Computer Science*, 31:22–31.
- [6] Asmaa Elbadrawy and George Karypis. 2016. Domainaware grade prediction and top-n course recommendation. In *Proceedings of the 10th ACM Conference on Recommender Systems*, pages 183–190. ACM.
- [7] Surabhi Dwivedi and VS Kumari Roshni. 2017. Recommender system for big data in education. In *E-Learning & E-Learning Technologies (ELELTECH), 2017 5th National Conference on*, pages 1–4. IEEE.
- [8] Yiu-Kai Ng and Jane Linn. 2017. Crsrecs: A personalized course recommendation system for college students. In *Information, Intelligence, Systems & Applications (IISA), 2017 8th International Conference on*, pages 1–6. IEEE..
- [9] Abdon Carrera Rivera, Mariela Tapia-Leon, and Sergio Lujan-Mora. 2018. Recommendation systems in education: A systematic mapping study. In *International Conference on Information Theoretic Security*, pages 937–947. Springer.
- [10] Mahamudul Hasan, Shibbir Ahmed, Deen Md Abdullah, and Md Shamimur Rahman. 2016. Graduate school recommender system: assisting admission seekers to apply for graduate studies in appropriate graduate schools. In *Informatics, Electronics and Vision (ICIEV), 2016 5th International Conference on*, pages 502–507. IEEE.
- [11] Stein, S.A., M. Weiss, G., Chen, Y. and Leeds, D.D., 2020, September. A College Major Recommendation System. In *Fourteenth ACM Conference on Recommender Systems* (pp. 640-644).
- [12] Park, Y., 2017. A Recommender System for Personalized Exploration of Majors, Minors, and Concentrations. In *RecSys Posters*.
- [13] Alshaikh, K., Bahurmuz, N., Torabah, O., Alzahrani, S., Alshingiti, Z. and Meccawy, M., 2021. Using Recommender Systems for Matching Students with Suitable Specialization: An Exploratory Study at King Abdulaziz University. *International Journal of Emerging Technologies in Learning (iJET)*, 16(3), pp.316-324.
- [14] Lamiaa Mostafa, Giles Oatley, Nermin Khalifa, and Walid Rabie. 2014. A case based reasoning system for academic advising in egyptian educational institutions. In *2nd International Conference on Research in Science, Engineering and Technology (ICRSET2014) March*, pages 21–22..
- [15] Rokach Lior et al. 2014. *Data mining with decision trees: theory and applications*, volume 81. World scientific.
- [16] Xavier Amatriain and Josep M Pujol. 2015. *Data mining methods for recommender systems*. In *Recommender systems handbook*, pages 227–262. Springer.

APPENDIX

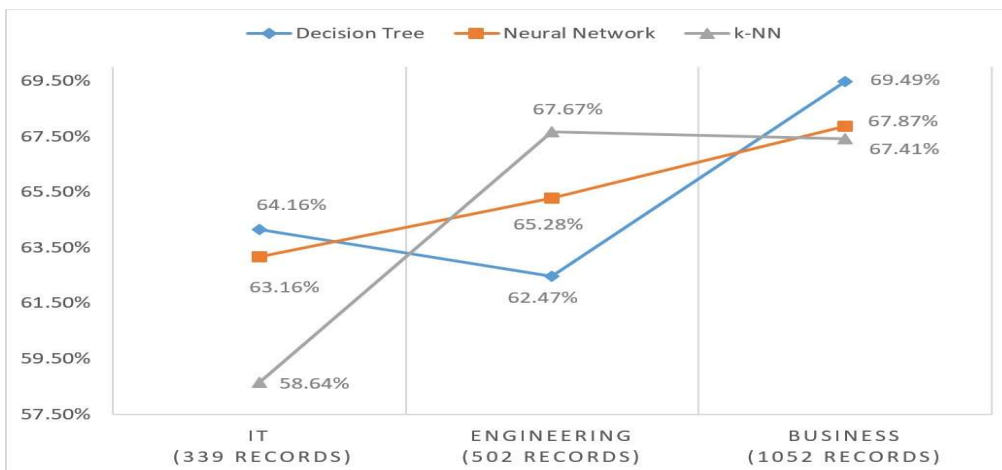


Figure 8: Classification Performance

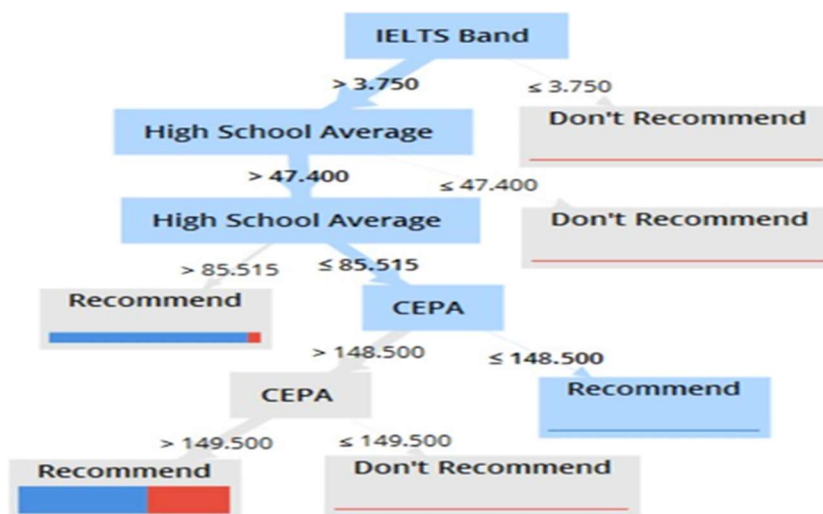


Figure 9: Decision Tree for the Business Program

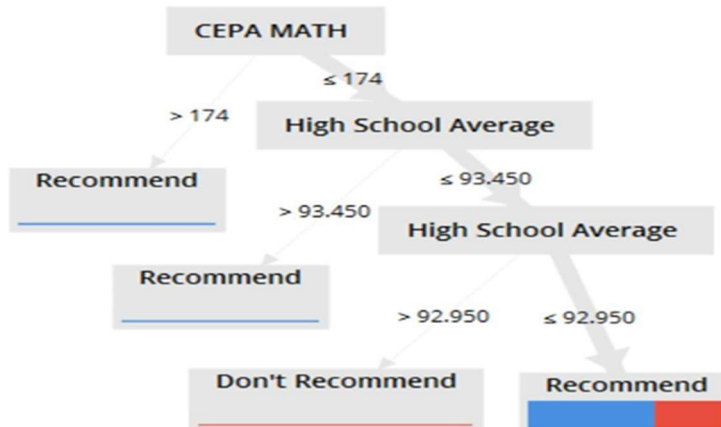


Figure 10: Decision Tree for the IT Program

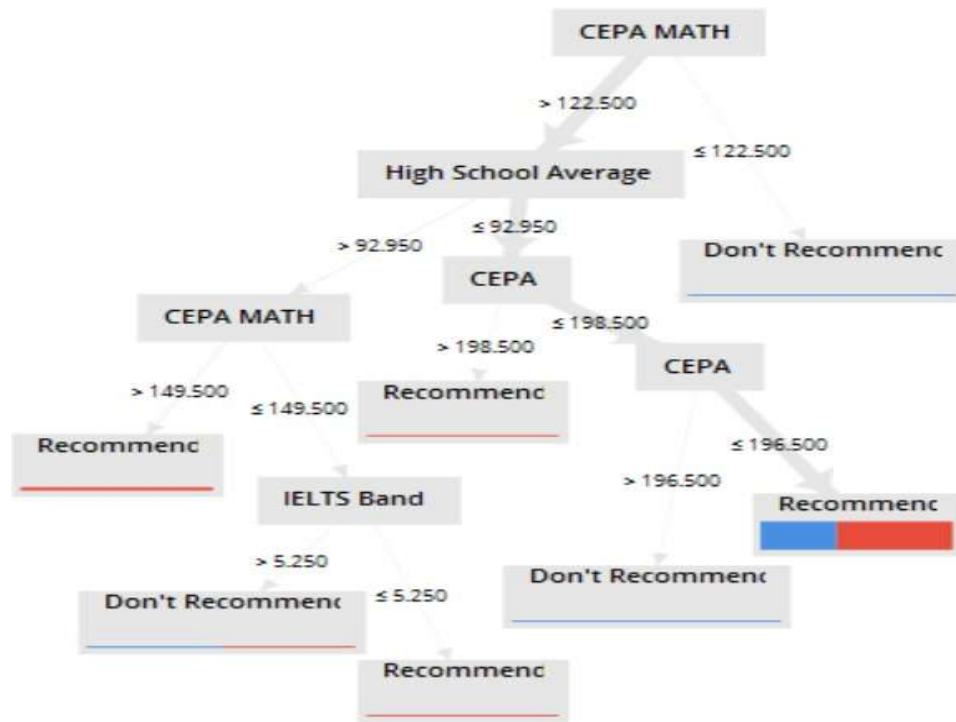


Figure 11: Decision Tree for the Engineering Program

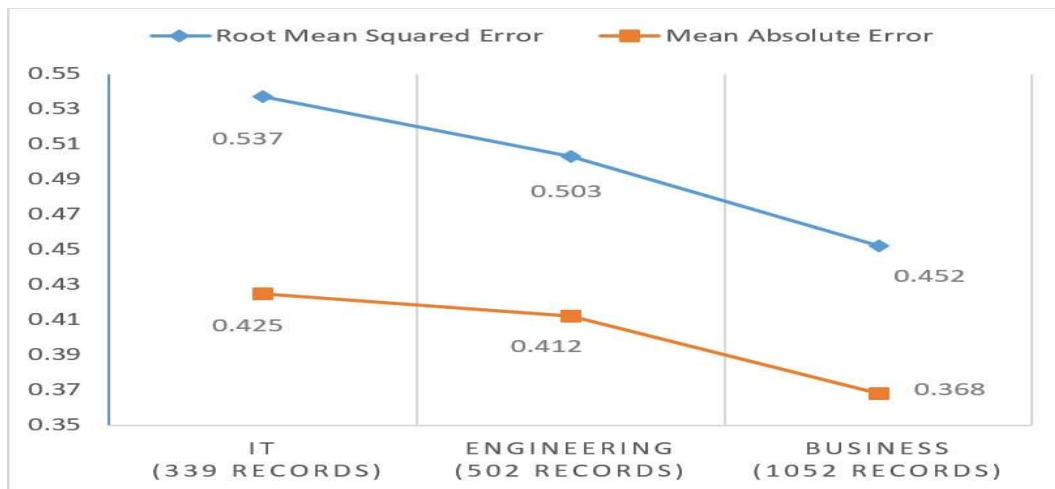


Figure 12: Linear Regression Performance