# Analyzing and Predicting Career Specialization using Classification Techniques

**Jennifer Anne A. Repaso[1], Elenita T. Capariño[2]**
[1]Bulacan State University, Philippines, jarepaso@gmail.com
[2]Bulacan State University, Philippines, elenitacaparino@bulsu.edu.ph

## ABSTRACT

Nowadays, academic institutions conduct studies to attain quality and excellence in student academic performance through Data Mining tools. This paper explores various classification techniques in predicting graduates' career specialization. The data sets used were obtained from Bulacan State University Sarmiento Campus' Information Technology graduates from 2013 to 2016. From these data, a model was created using Naïve Bayes, J48, Random Forest, and Support Vector Machine classification algorithm, with 18 attributes. Among the models built, Naïve Bayes and Random Forest algorithm yielded better accuracy rating, and acceptable ROC and RMSE values. Performance of students per subject area was also determined, and based on this; students perform satisfactorily in both soft and technical skills manifested in the highly satisfactory performance on the Internship course. However, there were few graduates who pursue a career in Networking, and measures must be undertaken to elevate performance in early Programming and Networking courses..

**Key words :** Career Specialization, Classifications Algorithm, Data mining, Predictions.

## 1. INTRODUCTION

The employment of a student after graduation is one of the concerns of an academic institution. Assurance that a student will be employed based on their specialization is a great concern.

The demand for IT Practitioners is high may it be in our country, and abroad. In a survey conducted by ITAA, US firms cite a shortage of qualified IT personnel [1]. In the Philippines, it was noted that the demand for IT specialist is high, but higher institutions cannot supply the right manpower needed by the industry. Although educators are always on the right track in preparing the students technically; however, potential employers were a little bit frustrated that applicants lacked skills coupled with technical skills.

A graduate of an information technology program is expected to land a job appropriate to his specialization. Some of these job titles are : Application Programmer, Software Engineer,Network Administrator, Web Developer, Technical Support [2].

This study will cover the community of Bulacan State University-Sarmiento Campus, particularly graduates of the BS in Information Technology Program. This is significant to the students enrolled in the IT program since they are the direct beneficiaries of the study. This would also help the I.T. department of the campus to be able to guide the students as they pursue their studies, and this will aid the enrolled student in attaining relevant career discipline after they graduate.

This study aims to create a model to be able to predict the I.T. career specialization based on the grades in the general technical skills and soft skills using classification techniques namely Naïve Bayes, J48, Random Forest and Support Vector Machine; compare the different classification techniques by determining the level of accuracy of each model; and analyze the general performance of the students in each subject area.

## 2. RELATED WORKS

Practically every enterprise nowadays is "IT – Dependent." They depend upon the information and communications technology tools and services. Boston Area Advanced Technological Education Connections (BATEC) stated in its Information Technology Workforce Skills Study-3/4 of IT jobs are in IT-Enabled firms, not in IT-Production firms (software and hardware)[1]. IT students do concentrate on their core subjects. But technical skills are not just the skills being searched for by employers. Soft skills or general business skills are equally important [3]. IT/IS educators should emphasize on general technical skills (as opposed to specific technology skills). Among the general-technology capacity levels they listed were: I.T. Project management, testing, security, analysis and design of systems, development of software, integration of systems, networking of systems, maintenance, design of business processes, programming, management of operating systems, architecture and development of internet systems, strategic use of IT and planning, Database design and administration, Web page development, and Human-computer interaction [4]. Similarly, there are categories of different job skills. The category that has the largest number of job skills was programming languages. The other categories were followed by Web development, Database, Operating System and Environments, and Networking [5]. IT/IS educators should emphasize on general business skills as well. Among the top 5 rank of the

general business skills they have identified were: Teamwork, Communications, Interpersonal skills, Organizational problem solving skills and Creativity [4]. On a similar study, another considered category is the soft skills for I.T. jobs as problem variables, and these includes : Communication/Interpersonal, Team Work, Problem Solving/Analytical Thinking, Time Management, Leadership, Creativity, Stress Management and Independent/ Self Motivated [6].

The most important attributes in terms of personal and business capability are the non-technical skills [7]. Accordingly, for IT professional, soft skills are more important specially needed and applied on problem solving, critical thinking and team skills. They also determined the rank and importance of the core knowledge of information systems. Findings of the study proves that technical skills are important in the use of database, programming and web development applications [8].

General technical skills and business skills are essential in the employability of IT graduates. Having been good in these areas during the tertiary years would ensure employability.

Among the top information technology jobs which are in demand were: Application Programmer/Analyst, Business Analyst, Software Engineer, Manager/Project Leader, Sales/Marketing, Web Developer, System Administrator, Database Administrator, Technical Support [2].

In predicting the IT career specialization, the present profession of selected IT graduates will be used, specifically graduates who are: Software Developers/Programmers/Engineers, Network Administrator and Web Developer. Grades incurred by these IT practitioners on the courses which are considered part of the general technical skills and general business skills identified by Peslak et. al., and Koong, et. al. will be used as attributes which form part of the data set, and will be analyzed using data mining technique.

In predicting the I.T. career specialization, data mining techniques were used. Data mining in higher education is a recent research field and this area of research is gaining popularity because of its potentials to educational institutions. Studies were conducted in the academic field to assist the administration of different institution in decision making process.

Data Mining is useful in enhancing our knowledge especially on identifying, extracting and evaluating variables related to the learning process of the students in the educational field [3]. Data Mining refers to pulling out useful information
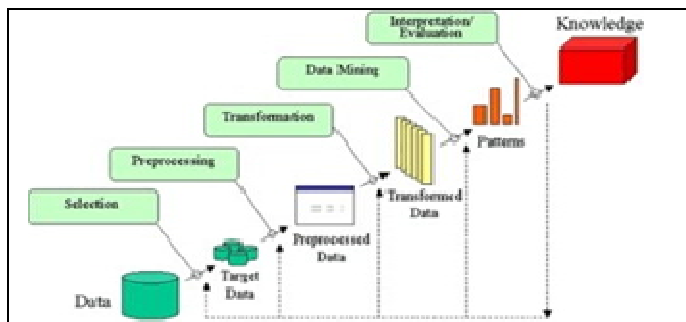


**Figure 1:** Knowledge Discovery Process

from large amounts of data. This is also useful in finding out patterns and relationship on big data that will be helpful in decision making [9]. Classification, Clustering, Regression, Artificial Intelligence, Neural Networks, Association Rules, Decision Trees, Genetic Algorithm, Nearest Neighbor method are just some of the algorithms and techniques that can be used for knowledge discovery from databases. Applying educational data mining (EDM.) is a growing interdisciplinary research field in education. It aims to develop methods for discovering the unique types of data from educational environments. It is helpful for the students to better understand and learn how to expand educational conclusions and to gain understanding from it and clarify educational phenomena [10]. Moreover, in EDM, coming from the academic settings, data which are considered "out of sight" can be mined, analyzed and explored with the applications of newer methods or schemes in consideration of the students' learnings [11].

Many studies have been conducted in relation to mining educational data. Irina Ionita (2015) conducted a study entitled Data Mining for Predicting the Military Career Choice. In this experiment several algorithms have been applied on a sample data of (274 records) such as: logistic regression, J48, JRIP, LMT, REPTree and Simple CART. The results were analyzed and it was observed that Simple CART, followed by J48, LMT and logistic regression has the best rate of classification [12].



**Figure 2:** Framework of the Study

**Table 1:** Attributes and their corresponding skills

| Attributes | General Technology Skills | Soft Skills |
|---|---|---|
| CompProg1 | Programming | Problem solving, creativity |
| CompProg2 | Programming | Problem solving, creativity |
| CompProg3 | Programming | Problem solving, creativity |
| CompProg4 | Programming | Problem solving, creativity |
| DBM1 | Database Management | Problem solving, creativity |
| DBM2 | Database Management | Problem solving, creativity |
| WebDev | Web Development | Problem solving, creativity |
| Network1 | Networking, Maintenance | Problem solving, creativity |
| Network2 | Networking, Maintenance | Problem solving, creativity |
| SAD | Systems Analysis and Design, System Integration | Problem solving, creativity, teamwork, interpersonal skills, strategic utilization of information technology and planning |
| SoftEng | Testing, Security, System Integration | Problem solving, creativity |
| OpSys | Operating System Environments | Problem solving, creativity |
| ComSki1 | | Communication skills |
| ComSki2 | | Communication skills |
| ComSki3 | | Communication skills |
| ProfEth | | Teamwork, interpersonal skills, work ethics |
| CapProj | Project Management | Problem solving, creativity, teamwork, interpersonal skills |
| Internship | | Teamwork, interpersonal skills |

Ivy Tarun, et. al. conducted a study in predicting Licensure Examination Test (LET) performance by implementing PART and JRip classifiers. They were able to predict whether a student will pass or fail the LET Examination after the mock examination using data mining techniques. The study was able to predict that if the mock board rating obtained is lower that 34% of the total points the reviewee will fail the LET exam based on the model obtained using PART and JRip classifiers [13].

In a study conducted by Bhardwaj and Pal (2011) on the performance of 300 students from 5 different degree college offering BCA (Bachelor of Computer Application) course of Dr. R. M. L. Awadh University, Faizabad, India. It was concluded that factors such as senior secondary examination grade, living place, teaching medium, mother's qualification, student behavior, family annual income and family status of the student were highly correlated with academic performance of the students. This was predicted on 17 attribute by way of the Bayesian classification[9].

**Table 2:** Attributes and their values

| Attributes | Course Description | Values |
|---|---|---|
| CompProg1 | Computer Programming 1 | |
| CompProg2 | Computer Programming 2 | |
| CompProg3 | Computer Programming 3 | |
| CompProg4 | Computer Programming 4 | |
| DBM1 | Database Management 1 | |
| DBM2 | Database Management 2 | |
| WebDev | Web Development | |
| Network1 | Networking 1 | |
| Network2 | Networking 2 | |
| SAD | Systems Analysis and Design | |
| SoftEng | Software Engineering | VG, G, F |
| OpSys | Operating Systems | |
| ComSki1 | Communication Skills 1 | |
| ComSki2 | Communication Skills 2 | |
| ComSki3 | Communication Skills 3 | |
| ProfEth | Professional Ethics | |
| CapProj | Capstone Project | |
| Internship | Internship | |
| Year | | Numeric |
| JobClass | | DEV, ITSUPPORT, SYSAD, GA, OTHERS |

Vasile Paul Brefelean (2007) in the study entitled "Analysis and Predictions on Students' Behavior Using Decision Trees in Weka Environment" was an implementation of a J48 algorithm analysis tool on data collected from surveys on different specialization students, with the purpose of differentiating and predicting their choice in continuing their education with post university studies (master degree, Ph.D. studies) through decision trees [14]. This is another application where data mining techniques are applied in the academe.

A decision tree model was applied in the study by Al-Radaideh, et. al. to predict the final grade of students who took the C++ course in Yarmouk University, Jordan in the year 2005. ID3, C4.5, and NaïveBayes are the different classification methods used [15]. Among all methods used Decision Tree model yielded a better prediction.

Galit gave a case study that use students' data to analyze their learning behavior to predict the results and to warn students at risk before their final exams [16].

**Table 3:** Numerical Grade Equivalent of Attributes

| Value | Numerical Grade |
|---|---|
| VG | 1.0, 1.25, 1.5 |
| G | 1.75, 2.0, 2.25 |
| F | 2.5, 2.75, 3.0 |

**Table 4:** Class and Description

| Class | Description |
|---|---|
| DEV | Software Engineer / Developer / Programmer / Web Developer |
| ITSUPPORT | Quality Assurance Personnel / Technical Support / IT Support / Data Analyst |
| SYSAD | System Administrator / Network Architect / Network Administrator |
| GA | Graphic Artist / Multimedia Specialist / Game Developer |
| | |
| OTHERS | Not mentioned from the above |

Recent studies about EDM shows that data mining methods indeed are essential in the decision-making of higher officials for decision-making and in improving process and operations in the academic settings [17] [18] [19]. The framework of the study used by Fayyad, et.al presented in Figure 1. was derived on the knowledge discovery process and knowledge discovery database [20]. The figure. shows how the datasets are being used, processed, and transformed in order to make some forecast or prediction. With these, the authors came up with a framework illustrated in Figure 2. These processes were applied on the study.

## 3. METHODOLOGY

In this study, the IT career specialization is predicted based on the data sets of the graduates from 2013 to 2016. The framework of this study is shown in Figure 2. Student records were retrieved from the Student Information System. The professions of the graduates were also used from the Graduate Tracer. From these, selection was performed, pre-processed and cleaned, and integrated into one database. Attributes were also selected and assigned based on the corresponding skills in general technical and soft skills, as indicated in Table 1.

**Table 5:** Percentage Split (66% Training data)

| Percentage Split | Classifier | Accuracy | ROC | RMSE |
|---|---|---|---|---|
| 66% Training Data | Naïve Bayes | 65% | 0.649 | 0.3303 |
| | Support Vector Machine | 45% | 0.638 | 0.3644 |
| | J48 | 50% | 0.406 | 0.4224 |
| | Random Forest | 65% | 0.704 | 0.3242 |

**Table 6:** Percentage Split (68% Training Data)

| Percentage Split | Classifier | Accuracy | ROC | RMSE |
|---|---|---|---|---|
| 68% Training Data | Naïve Bayes | 68.4211% | 0.720 | 0.3049 |
| | Support Vector Machine | 42.1053% | 0.69 | 0.3616 |
| | J48 | 52.6316% | 0.492 | 0.3885 |
| | Random Forest | 63.1579% | 0.75 | 0.3136 |

The attributes and their corresponding values are also observed in Table 2. Categorical values were derived from the numerical grades of the graduates as shown in Table 3. The classes were also determined during this stage. The class and description are shown on Table 4. The next step was to perform the classification task. The data mining algorithm used and contrasted were Naïve Bayes, J48, Random Forest and Support Vector Machine, using Weka tool. Assessment was also done for each classification task. In addition, charts were created to visualize the results of the grades obtained from the data sets for each course taken by the students when they were taking up the information technology program. For this, a spreadsheet application is used which is MS Excel. A simple analysis was derived from the chart which is discussed in the later section. In addition, the study of [21], [22], [23] and [24] were used as the basis for some of the implementation of the research.

## 4. RESULTS AND DISCUSSION

Weka was used as a tool in the data mining task. There were a total of 60 instances in the data set. A snapshot of the data file done for Weka is shown in Figure 3. Observe that there were some missing values (?) in the set of data. For each classifier, a percentage split from Weka's test options was selected. Initially, a 66% percentage split was chosen. This means that the data set is divided into two, that is, the training set, and the test set. Weka will be the one to split the data and create the model. The training set is 66%, and the remaining 33% is the test set. As observed in Table 5, accuracy, ROC, and root mean square error (RMSE) were indicated. For this, it is Naïve Bayes and Random Forest classifier which yielded an accuracy rate of 65%, with an acceptable ROC and RMSE value. Support Vector Machine has a low accuracy rate; however, it obtained a ROC value of 0.638 and RMSE of 0.3644 which may be considered acceptable as well. For Table 6, a percentage split of 68% was selected for the training data set. Obviously, it is Naïve Bayes and Random Forest algorithm that obtained an accuracy rate of 68.4211% and 63.1579% respectively; with an acceptable ROC and RMSE value. On the other hand, notice that J48 yielded an accuracy rate of 52.6316% which is higher than that of Support Vector Machine; however, has lower values of ROC and RMSE.

```
@relation career_alldata

@attribute Year NUMERIC
@attribute CompProg1 {VG, G, F}
@attribute CompSki1 {VG, G, F}
@attribute CompProg2 {VG, G, F}
@attribute CompSki2 {VG, G, F}
@attribute CompProg3 {VG, G, F}
@attribute CompProg4 {VG, G, F}
@attribute DBM1 {VG, G, F}
@attribute ComSki3 {VG, G, F}
@attribute DBM2 {VG, G, F}
@attribute Network1 {VG, G, F}
@attribute SAD {VG, G, F}
@attribute OpSys {VG, G, F}
@attribute Network2 {VG, G, F}
@attribute SoftEng {VG, G, F}
@attribute ProfEth {VG, G, F}
@attribute WebDev {VG, G, F}
@attribute CapProj {VG, G, F}
@attribute Internship {VG, G, F}
@attribute JobClass {DEV, ITSUPPORT, SYSAD, GA, OTHERS}

@data
2016,?,G,G,G,F,G,F,G,G,G,G,G,G,G,G,G,VG,ITSUPPORT
2016,F,F,F,F,F,F,F,G,G,G,F,G,G,VG,DEV
2016,VG,?,VG,G,G,VG,G,VG,G,G,?,VG,G,F,F,G,G,VG,OTHERS
2015,G,G,G,G,F,?,G,G,G,F,G,G,G,G,F,G,VG,ITSUPPORT
2013,F,VG,F,G,F,VG,G,F,G,F,G,G,G,G,G,ITSUPPORT
2013,F,F,G,F,G,VG,VG,?,VG,F,G,VG,G,F,G,VG,VG,G,DEV
2013,G,G,G,G,VG,VG,VG,?,VG,F,VG,G,G,G,VG,G,VG,ITSUPPORT
2013,G,G,G,G,VG,VG,G,G,F,VG,G,G,F,G,VG,VG,ITSUPPORT
2013,VG,?,VG,F,VG,VG,G,?,VG,G,G,VG,G,G,F,G,G,VG,DEV
2013,F,G,F,G,G,VG,VG,G,VG,F,G,G,G,G,G,G,G,SYSAD
```
**Figure 3:** Arff File for WEKA

To sum it up, for the data set used in this study, it is better to use Naïve Bayes or Random Forest as a classification algorithm since they show acceptable accuracy, ROC and RMSE values.

In analyzing the general performance of students, charts were created using MS Excel to visualize the results as shown in Figure 4 to Figure 8. As observed, students consistently performed satisfactorily in Communication Skills and Ethics as illustrated in Figure 4. However, students' performances were noticeably levelled up on the succeeding courses.

As observed, there were only a few graduates who pursued a career related to Networking despite being satisfactory in Networking 2 course as shown in Figure 7. Furthermore, based on Figure 8, efforts must be exerted to achieve a better performance in Software Engineering aside from the other subject areas like Programming where, noticeably in Figure 5, most students were fairly performing in the first programming course, but has improved on their higher programming courses. This may be due to the adjustment period being experienced by freshmen students during the first quarters of their tertiary level in the Information Technology program. Obviously, as seen in Figure 8, students performed very satisfactorily in their internship.

## 5. CONCLUSION

The results shown in this paper revealed that Naïve Bayes and Random Forest classification algorithm may be appropriate in predicting the career specialization of a graduate. The result was illustrated in figures 4-8. Also, more data for training and testing may be used in order to increase the level of accuracy. Also, freshmen students show fair performance in their first Programming and Networking course. Students perform satisfactorily in their professional courses especially in Systems Analysis and Design, Capstone Project, and

Internship. In this regard, the information technology department of the campus and the faculty members may implement programs to improve the performance level of students especially in the first year level specifically in Programming, and on some courses in the 2nd and 3rd year level which is Networking and Software Engineering. Moreover, good practices of the department should be maintained in sustaining the students' performance particularly on their communication, ethics, and professional courses.
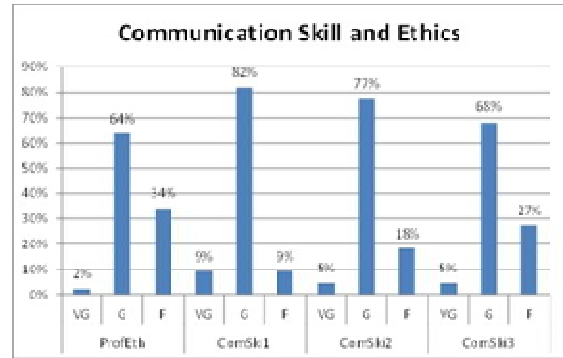


**Figure 4:** Performance on Communication and Ethics course
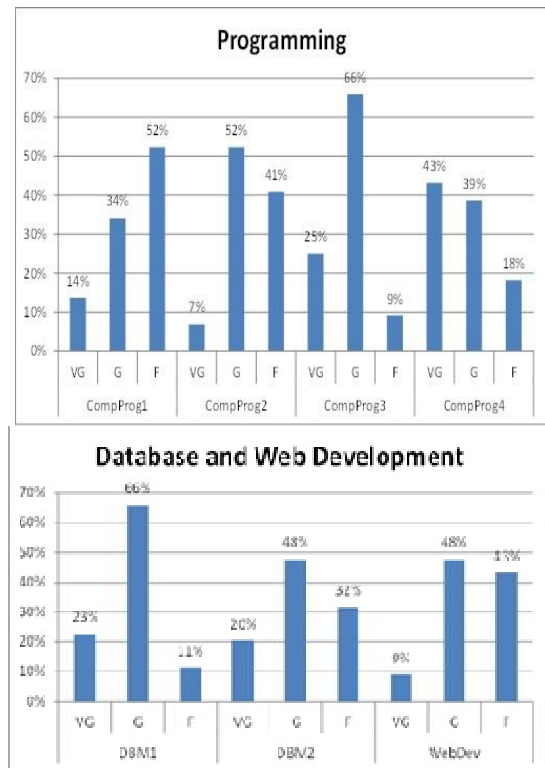
**Figure 5**: Performance on Programming courses





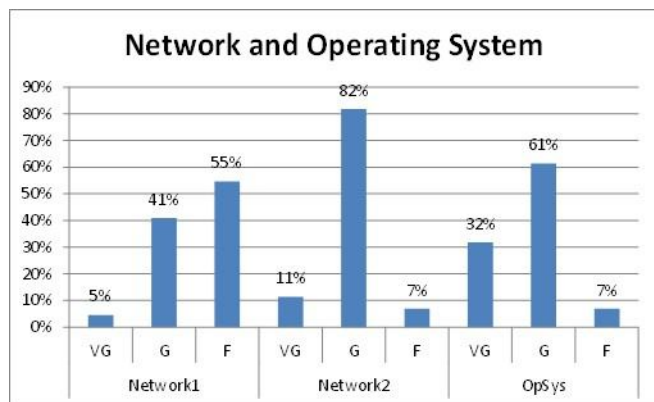**Figure 6:** Performance on Database and Web Development course

**Figure 7:** Performance on Network and Operating System course
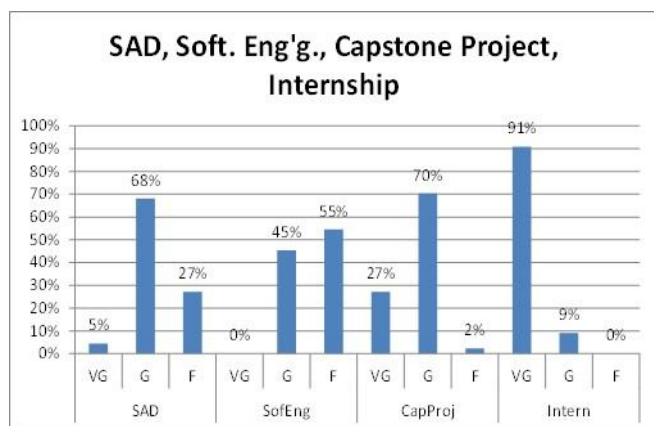


**Figure 8:** Performance on other professional courses

## REFERENCES

1. U. of M. Boston, "Information Technology Workforce Skills," *BATEC Inf. Technol. Work. Ski. Study*, 2007.
2. S. D. Galup, R. Dattero, and J. J. Quan. **The Demand for Information Technology Knowledge and Skills : An Exploratory Investigation**, *Journal of International Technology and Information Management*, vol. 13, no. 4, 2004.
3. A. El-Halees, **Mining Students Data To Analyze Learning Behavior : a Case Study Educational Systems**, *Dep. Comput. Sci. Islam. Univ. Gaza P.O.Box 108 Gaza, Palest.*, no. February, 2008.
4. A. R. Peslak and G. A. Davis, **An Empirical Study Of The Relative Importance Of Specific Technology Skills , General Business Skills , And General Technology Skills**, vol. X, no. 2, pp. 430–438, 2009.
5. K. S. Koong and L. C. Liu, **A Study of the Demand for Information Technology Professionals in Selected Internet Job Portals**, *Journal of Information Systems Education, Vol 13(1)*, May, 2015.
6. A. A. Bakar and Choo-Yee Ting, **Soft skills recommendation systems for IT jobs: A Bayesian network approach**, *Conf. Data Min. Optim.*, pp. 82–87, 2011.
7. M. E. McMurtrey, J. P. Downey, S. M. Zeltmann, and W. H. Friedman, **Critical Skill Sets of Entry-Level IT Professionals: An Empirical Examination of Perceptions from Field Personnel**, *Journl of Information Technology Education*, Vol. 7, pp. 101–120, 2008.
   https://doi.org/10.28945/181
8. S. Hawk *et al.*, **A Typology of Requisite Skills for Information Technology Professionals A Typology of Requisite Skills for Information Technology Professionals**, *International Conference on System Sciences*, 2011.
9. B. K. Bhardwaj, S. Pal, **Mining Educational Data to Analyze Student' Performance**, *International Journal of Advanced Computer Science and Applications*, Vol. 2, No. 6, pp. 63-69, 2012.
10. C. Romero and S. Ventura, **Data mining in education**, vol. 3, no. February, pp. 12–27, 2013.
11. B. K. Bhardwaj, D. J. Finch, L. K. Hamilton, and R. Unveils, **A Critically Review of Data Mining Segment : A New Perspective**, *International Journal of Advanced Trends in Computer Science and Engineering,* Vol. 8 No. 6, pp. 2984–2987, 2019.
    https://doi.org/10.30534/ijatcse/2019/50862019
12. I. Ionita, "**Data Mining For Predicting The Military Career Choice**, vol. 3, no. 3, 2015.
13. I. M. Tarun, B. D. Gerardo, and B. T. Tanguilig III, **Generating Licensure Examination Performance Models Using PART and JRip Classifiers: A Data Mining Application in Education**, *International Journal of Computer and Communication Engineering*, Vol. 3, No. 3, pp. 202–207, 2014.
14. V. P. Breúfelean, **Analysis and Predictions on Students' Behavior Using Decision Trees in Weka Environment**, *Int. Conf. on Information Technology Interfaces,* pp. 51–56, 2007.
15. M. I. Al-Radaideh, Qasem A. Emad Al-Shawakfa, Emad M. and Al-Najjar, **Mining Student Data Using Decision Trees**, *2006 International Arab Conference on Information Technology*, Vol. 40, No. 2, pp. 1–18, 2015.
16. G. Ben-Zadok, A. Hershkovitz, R. Mintz, and R. Nachmias, **Examining online learning processes based on log files analysis: A case study**, *Res. Reflections Innov. Integr. ICT Educ.*, pp. 55–59, 2015.
17. A. O. Gamao and B. D. Gerardo, **Prediction-Based Model for Student Dropouts using Modified Mutated Firefly Algorithm**, *International Journal of Advanced Trends in Computer Science and Engineering*, Vol. 8 No.6, pp. 3461-3469, 2019.
    https://doi.org/10.30534/ijatcse/2019/122862019
18. J. S. Gil, I. Journal, J. S. Gil, A. J. P. Delima, and R. N. Vilchez, **Predicting Students ' Dropout Indicators in Public School using Data Mining Approaches**, *International Journal of Advanced Trends in Computer Science and Engineering Available*, Vol. 9, No. 1, pp. 5–9, 2020.
    https://doi.org/10.30534/ijatcse/2020/110912020
19. P. G. L. Denila, A. J. P. Delima, and R. N. Vilchez, **Analysis of IT Graduates Employment Alignment Using C4 . 5 and Naïve Bayes Algorithm**, *International Journal of Advanced Trends in Computer Science and Engineering Available*, Vol. 9 No.1, pp 745-752, 2020.

https://doi.org/10.30534/ijatcse/2020/106912020

20. U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, **From Data Mining to Knowledge Discovery in Databases**, *AI Mag.*, vol. 17, no. 3, p. 37, 1996.

21. G. Alcober, T. Revano, & M. Garcia, **E-Safety in the Use of Social Networking Application**, *International Journal of Advanced Trends in Computer Science and Engineering*, Vol. 9, No. 1.2 2020, https://doi.org/10.30534/ijatcse/2020/1291.22020

22. R. Dellosa, **An Efficient Position Estimation of Indoor Positioning System Based on Dynamic Time Warping**, *International Journal of Advanced Trends in Computer Science and Engineering*, Vol. 9, No. 1.2 2020, https://doi.org/10.30534/ijatcse/2020/0491.22020.

23. J. Victoriano & L. Lacatan, **A Geospatial Analysis and Kernel Density Estimation of River Quality Parameter in Bulacan, Philippines**, *International Journal of Advanced Trends in Computer Science and Engineering*, Vol. 9, No. 1.2 2020, https://doi.org/10.30534/ijatcse/2020/1191.22020.

24. Elenita T. Capariño, Ariel M. Sison, Ruji P. Medina, **Application of the Modified Imputation Method to Missing Data to Increase Classification Performance**, *2019 IEEE 4th International Conference on Computer and Communication Systems (ICCCS),* 2019, DOI: 10.1109/CCOMS.2019.8821632.

25. Miriam P. Pariñas, Ariel M. Sison, Ruji P. Medina, **A Modified Apriori Algorithm to Mine Association Rules using Relative Multiple Supports**, *International Journal of Advanced Trends in Computer Science and Engineering*, Vol. 9, No. 1.2 2020, https://doi.org/ 10.30534/ijatcse/2020/1091.12020.