

# Contextual Data Mining for Higher Educational Institutions

Subhashini Sailesh Bhaskaran, Mansoor Al Aali, Kevin Lu

**Abstract:** Context-awareness research has been carried out by many researchers. However, little has been done in building a context-aware data mining methodology which supports decision-making in HEIs. The usefulness of Knowledge discovery data mining process (KDDM process) in HEIs were investigated to discover hidden knowledge that is contextualized, resident in student datasets and use them in decision making. It was experimented and found that not any of the KDDM processes include a contextual factor mining stage that is essential to take out hidden knowledge from datasets described by contextual factors. Therefore a new process was introduced in KDDM process that uncovered contextual data to be used to support business goal and produced a dataset at the preparation stage which generated data mining model that was contextual leading to the unearthing of course taking patterns that are contextualized. This discovery has enabled forecasting of optimum CGPA and time-to-degree.

**Index Terms:** HEIs, Data Mining, KDDM, Time to Degree, Student Performance, Context-Awareness.

## I. INTRODUCTION

In the delivery of education in HEIs decision making processes assume significance. Those decision making processes are concerned with a number of activities in HEIs including teaching and student learning experience. For instance decisions could impact such aspects as curriculum design, categorization of students, timetabling improvements and student assessment that can be related to teaching and student learning experience (BIS, 2014). However accuracy and applicability of those decisions in HEIs are areas of major concern, because common decisions that are taken in HEIs at the institutional level do not always get implemented at the academic level (BIS, 2014). One reason for this could be that decisions usually made in HEIs use data that contain observable patterns or phenomena not entirely adequate to make accurate decisions. Observable phenomena include data and information about student gender etc. However literature points out that there is a need to consider unobservable phenomena also that could have serious implications to the decision making process. Unobservable phenomena could include data and knowledge associated to student education (for instance, student potential ,course difficulty). A major phenomenon that is found to be unobservable is the contextual factor that could contain knowledge helpful in HEIs decision making process. For instance when decisions are being made in HEIs regarding

optimizing the student time to degree, knowledge related to course difficulty or student potential could help HEIs in determining the categorization of students and courses as well as providing additional support to students who need support. While literature review shows that contextual factors have the potential to enhance the quality of decision making, extracting such contextual factors from known resources and discover useful knowledge about those contextual factors has been a challenge (Vert et al., 2010). One resource that is very promising that could be used to extract contextual factors and knowledge about those contextual factors is student dataset. Recently researchers have started to focus on contextual factors that are important to HEIs that could be used in their decision making process although more needs to be done. One such area that is clearly providing evidence of this kind of interest and provides a large scope for investigation is Knowledge discovery. However serious questions are raised on the usefulness of existing KDDM processes to enable decision making involving contextual factors related to a particular business environment like HEI. For instance Vert et.al (2010) argues that current KDDM processes need to be enhanced for their usefulness in supporting accurate decision making using contextual factors, an argument that implies that current KDDM processes are not capable of addressing unobservable knowledge that is contextualized. Some argue that further research is needed in this area to support users of KDDM processes (Singh, 2003; Vajirkar, 2003). Thus on the one hand there is a need to use KDDM processes to discover knowledge and on the other those KDDM processes must be enhanced to deal with contextual factors leading to the discovery of knowledge that is contextualized and useful to make more accurate decisions. In the context of HEIs the current level of understanding of the decision makers is limited to visible and tangible factors available in the student data set (For e.g. Gender, age, nationality, grades). Most HEIs are yet to deploy any type of KDDM process to discover hidden knowledge in their decision making process. For instance, a question could be raised regarding predictability of a pattern of courses a student could register in a semester or semesters that would lead the student to achieve a lower time to degree (the time taken by a student to graduate). Such a question is unlikely to be answered straightaway using the data readily available in the dataset like student factors namely semester based number of courses enrolled or Grade Point Average (GPA) scores as the knowledge related to pattern of courses that can be used to predict the time to degree is not readily available in the dataset. Here KDDM processes could be very useful although HEIs do not appear to have deployed such processes yet in their decision making process.

**Manuscript published on 30 June 2019.**

\* Correspondence Author (s)

**Dr. Subhashini Sailesh**, MIS Department, Ahlia University/ College of Business/ Ahlia University, Bahrain.

**Prof. Mansoor Al Aali**, President , Ahlia University, Bahrain.

**Dr. Kevin Lu**, Brunel University London.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](http://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

The reason for this situation are that manual methods or commonly used computing methods(e.g. SQL queries) are not adequate enough to discover patterns from hidden patterns from large datasets and HEIs lack knowledge to extract those patterns from student dataset using KDDM process. Moreover the use of extracted patterns from student dataset using KDDMs in decision making can suffer in terms of accuracy or achievement of business goals as the extracted patterns lack any knowledge related to contextual factors. The main reason for this is that KDDM processes that are currently used in the business domains or discussed in the literature, that have the potential to be used in HEIs, do not have steps to integrate contextual knowledge concealed in the dataset, needed for making better decisions to attain the business goals. This is a major gap in the literature. Thus the objective of this research is to address this gap to some extent by demonstrating a modified KDDM process that extracts hidden knowledge from the student dataset characterized by contextual factors that could be used in making better decisions in HEIs to achieve specific business goals. Towards this, in this research two contextual factors namely course difficulty and semester have been taken as examples to

produce course taking patterns from student dataset to forecast optimum time to degree and Cumulative Grade Point Average (CGPA). This knowledge about course taking patterns could be useful in taking many decisions in HEIs that have bearing on teaching and learning, for instance course categorization, student categorization, forecast of optimal time to degree of students and achievement of optimum CGPA.

II. RELATED LITERATURE

A. Review Stage

This section covers two aspects. One related to current knowledge on contextual factors and their utility in KDDM processes and the other related to KDDM processes. Literature review shows that contextual factors have been of interest researchers in various fields. Table1 provides a glimpse of current efforts of researchers taken in different fields that utilize contextual factors in KDDM processes to extract hidden knowledge from datasets that are contextualized.

Table 1, Current Efforts

Meaning of Context	Author and Paper	Remarks
Context, the cumulative history that is derived from data observations about entities (people, places, and things), is a critical component of analytic decision process. Without context, business conclusions might be flawed.	ContextBased Analytics in a Big Data World: Better Decisions, An IBM® Redbooks® Point-of-View publication  Sokol, 2013	By using context analytics with big data, organizations can derive trends, patterns, and relationships from unstructured data and related structured data. These insights can help an organization to make fact based decisions to anticipate and shape business outcomes. Entities are defined as people, places, things, locations, organizations, and events. Entities are an important focus of big data analytics. Context is defined as a better understanding of how entities relate. Cumulative context is the memory of how entities relate over time.
A general approach for contextaware adaptive mining of data streams that aims to dynamically and autonomously adjust data stream mining parameters according to changes in context and situations	ContextAware Adaptive Data Stream Mining,  Haghighia et. al, 2009	The researchers proposed an overall method for contextaware adaptive data mining that includes contextawareness into universal data stream mining and allows real time examination of data on board mobile devices in a clever and costeffective manner. They achieved Contextawareness through Fuzzy Situation Inference (FSI) that assimilates fuzzy logic in the CS model, an official context modeling and cognitive approach for assisting pervasive computing environments.
Context– background information  Contextual item set mining extracts frequent associations among items considering background	ContextualItem set Mining in DBpedia,  Rabatel et al, 2014	The authors exhibit the capacity of contextual item set mining. Contextual item set mining excerpts frequent associations between items bearing in mind the background

Context as a concept has been identified as important in the field of KDDM recently, due to the potential it has in adding value to the discovered knowledge. Lack of context in the discovered knowledge has been argued to be a limitation in various KDDM processes by researchers (Schilit et al., 1994; Dey, 2001; Bolchini et al., 2007) due to the important role context could play in decision making. (Vert et al., 2010)claim that a dataset characterized by contextual attributes if fed to the KDDM processes at a definite stage, then it is possible that the knowledge discovered through the mining processes will not only contain the normal attributes but also the contextual attributes. These arguments could be

extended to the HEIs also. Thus based on these arguments and as explained earlier there is a necessity to investigate the usefulness of KDDM processes in HEIs to discover hidden knowledge that is contextualised, resident in student datasets and use them in decision making. Notably if a method to extract contextual factors is identified and employed in KDDM processes then it is probable to dig out hidden knowledge that is contextualised for better decision making.



Furthermore, amongst the various factors that have been discussed by researchers in supporting decision making process in HEIs, prediction of optimum time to degree and CGPA by means of students' course taking patterns in association with contextual factors namely course difficulty and semester do not find place. While there are arguments supporting a linkage between course taking patterns, CGPA and time to degree, there is hardly any research outcome that has substantiated this claim using data mining techniques particularly considering semester and course difficulty as contextual factors. This research aims to address this aspect.

### A. KDDM Process

A number of KDDM processes have been used in many domains namely medicine, manufacturing, education and other businesses. Widely used KDDM processes are tabulated in table 2 which indicates the steps involved in them. While literature shows the usefulness of KDDM processes in extracting hidden knowledge from datasets for decision making, those processes are also found to have limitations. Table 3 provides a good idea about the limitations affecting different KDDM processes.

**Table 2, Evaluation of KDDM (adapted from Kurgan and Musilek, 2006)**

Model	Fayyad <i>et al.</i>	Cabena <i>et al.</i>	Anand & Buchner	CRISP-DM	Cios <i>et al.</i>	Generic model
Area	Academic	Industrial	Academic	Industrial	Academic	N/A
No of steps	9	5	8	6	6	6
Refs	(Fayyad <i>et al.</i> , 1996d)	(Cabena <i>et al.</i> , 1998)	(Anand & Buchner, 1998)	(Shearer, 2000)	(Cios <i>et al.</i> , 2000)	N/A
Steps	1 Developing and Understanding of the Application Domain	1 Business Objectives Determination	1 Human Resource Identification	1 Business Understanding	1 Understanding the Problem Domain	1 Application Domain Understanding
	2 Creating a Target Data Set	2 Data Preparation	2 Problem Specification	2 Data Understanding	2 Understanding the Data	2 Data Understanding
	3 Data Cleaning and Preprocessing		3 Data Prospecting	3 Data Preparation	3 Preparation of the Data	3 Data Preparation and Identification of DM Technology
	4 Data Reduction and Projection		4 Domain Knowledge Elicitation			
	5 Choosing the DM Task		5 Methodology Identification			
	6 Choosing the DM Algorithm		6 Data Preprocessing			
	7 DM	3 DM	7 Pattern Discovery	4 Modeling	4 DM	4 DM
	8 Interpreting Mined Patterns	4 Domain Knowledge Elicitation	8 Knowledge Post-processing	5 Evaluation	5 Evaluation of the Discovered Knowledge	5 Evaluation
	9 Consolidating Discovered Knowledge	5 Assimilation of Knowledge		6 Deployment	6 Using the Discovered Knowledge	6 Knowledge Consolidation and Deployment

**Table 3, Limitations of KDDM**

Process Model	Limitation	Limitation Identified By	Modified step or stage	New Concept
KDD(Fayyad <i>et al.</i> (1996)) 5 step)	Lack of data collection step which is vital for the KDD techniques in some real applications such as information security and medical treatment	(Ruan, 2007)	Data Collection step using previous mining results. It was added before data selections shown in Fig 2	Inclusion of Data Collection step in data mining process to filter irrelevant data leading to better decision making. (Ruan, 2007).
	Lack of domain knowledge leading to decision making that maybe useful if such a knowledge is not part of the mined data.	(Redpath & Srinivasan, 2004)	Proposed an architecture based on domain knowledge as shown in Fig 3	Introduction of Domain knowledge before the data selection (Redpath & Srinivasan, 2004)
	Unlike in other process models loop back to the second, third or fourth step are necessary due to prepared data which is not suitable for the mining process.	(Kurgan & Musilek, 2006)	Not addressed in the literature	Not addressed in the literature
	Lack of contextual information in the data	(Vert <i>et al.</i> , 2010)	Not addressed in the literature	Not addressed in the literature
CRISP - DM	Insufficiency to handle multidimensional temporal data resulting in knowledge that cannot support decision making which is dependent on temporal issues	(Catley <i>et al.</i> , 2009)	phases 1 (business understanding), 2 (data understanding), 4 (data modelling), and 6 (deployment) were enhanced to suit temporal data as shown in Fig 4	Introduction of Intelligent Data Analysis architecture in the mining process for ensuring the mined knowledge to support decision making that need temporal aspects(Catley <i>et al.</i> , 2009).
	The current CRISP -DM model is limited to address data that is free of human intervention	(Li <i>et al.</i> , 2009)	Proposed model had two backbones of the model, namely data mining and applied intelligent system, three participation elements namely On-Line Analytical Processing (OLAP), six sigma, and domain knowledge as shown in Fig 5	On-Line Analytical Processing (OLAP), six sigma, and domain knowledge (Li <i>et al.</i> , 2009).
	Lack of integrated process model.	(Sharma <i>et al.</i> , 2012)	Identification of task-task dependencies (between tasks of the same phase and different phases) is the first step	Building an integrated process model(Sharma <i>et al.</i> , 2012)

The previous information reveals that not any of the KDDM processes include a contextual factor mining stage which is essential to unearth unseen knowledge from data identified by contextual factors. Even though there are few papers that have talked about context knowledge and context driven data mining, such knowledge do not propose how to link data mining to contextual factor to help in decision making in businesses including HEIs that lead to accomplishment of business goals. When such addition could be attained, the resultant KDDM model can have a higher foretelling power which is essential for making better decisions. This has been addressed in this research in a limited manner by adding the CRISP-DM process with a contextual factor mining stage. In this context, this research relies upon KDDM process developed by Chapman. CRISP-DM model is given in Figure 1.

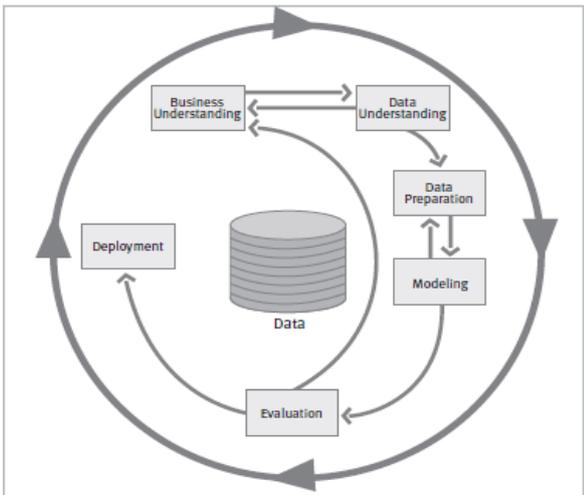


Figure 1, Current CRISP-DM model described in the literature (Adapted from Chapman et al., 2000)

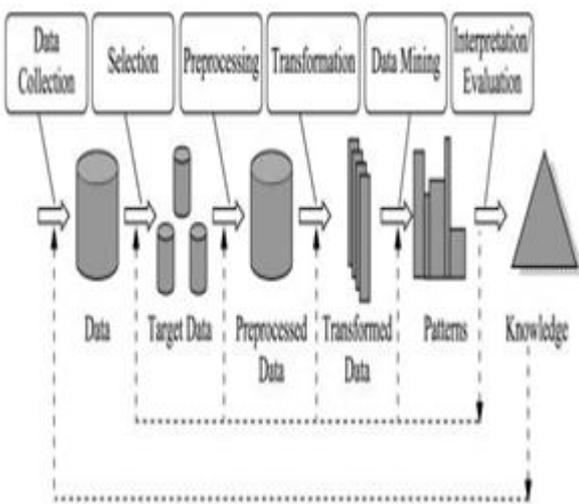


Figure 2, KDDM process model adapted from (Ruan, 2007)

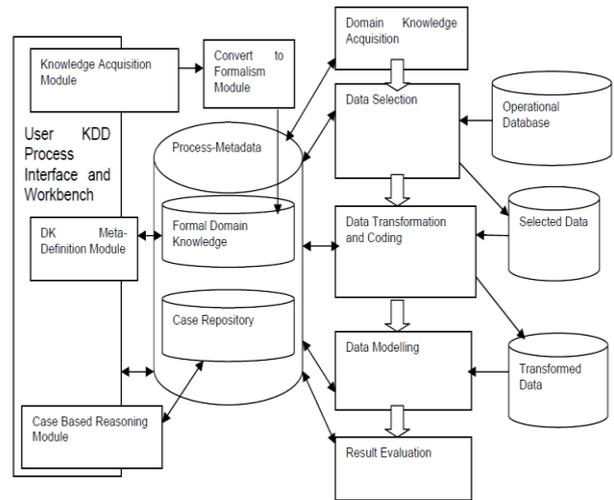


Figure 3, KDDM process model adapted from (Redpath and Srinivasan, 2004)



Figure 4, KDDM process (Catley et al., 2009)

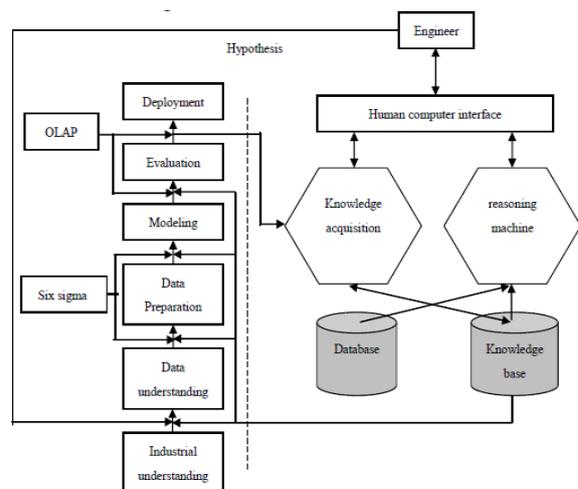


Figure 5, KDDM process model adapted from (Li et al., 2009)

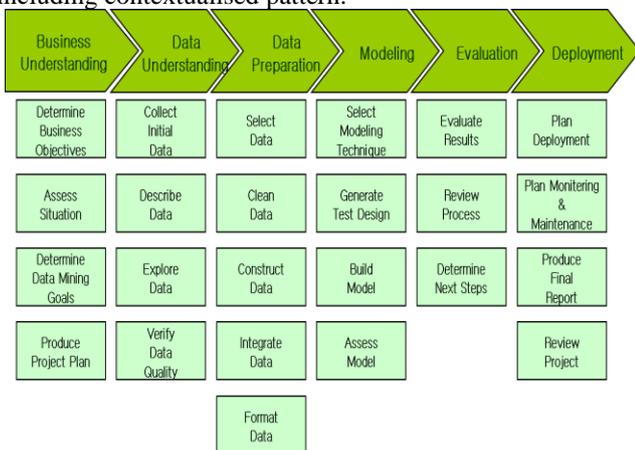


**I. THREE STAGE PROCESS TO MODIFY KDDM**

This section covers three stages. The first one covers a discussion on CRISP-DM process including its limitation to handle contextual data, modification needed to enable into handle contextual data and how a dataset with contextual factor can be used to discover knowledge through the process. The dataset thus identified is used for the next stage of discovering hidden knowledge. Second stage informs how the CRISP-DM method is actually modified to handle contextual dataset. The third stage provides information on the experiments conducted on a contextualised dataset by means of the tailored CRISP-DM process which include a method to detect the presence or absence of a particular contextual factor within a dataset using an algorithm developed for the purpose. In addition the third stage also depict how the modelling stage is linked to the contextual dataset that enable the CRISP-DM process to extract hidden knowledge from that student dataset described by contextual factor that could be useful to make decisions. The results of the experiment have been provided. In order to discuss the above mentioned aspects it is important to note that the researcher used the dataset of a HEI in Bahrain pertaining to BS programme in Accounting and finance (BSAF). The student dataset was extracted from the student registration system and has limited data pertaining to time to degree and CGPA. The method used in this research to determine the presence of context factor (course difficulty and semester) was the one suggested by Vert et al. (2010) which uses pseudo code. Each one of the stages mentioned above is explained below.

**II. TESTING AND EVALUATION**

Use Refer figure 1. A broad summary of the different stages of CRISP-DM process has been provided in table 5 because the detailed discussion on the CRISP-DM process is beyond the scope of this paper. Further due to paucity of space each stage and its components have not been discussed in detail. Instead important components that are necessary to support the experiments conducted in this research are given below (for a detailed description of various components of CRISP-DM refer to Chapman et.al (2000). In addition this section does not include discussion on any step related to contextual factor as the aim of this section is to fetch to see whether it can produce hidden pattern from the mined data including contextualised pattern.



**Figure 6, Current CRISP-DM model steps described in the literature (Adapted from Chapman et al., 2000)**  
The essential steps of CRISP-DM are given below.

The detailed specification of the stages and components can be had from the authors.(i) Business Understanding: To forecast the optimal time to degree using course taking patterns and CGPA and to categorize courses and students with unearthed knowledge. For this purpose data from student registration system of a university in Bahrain was used. Weka was used as the data mining tool.

(ii) Data Understanding: SQL queries was used to retrieve the data (see table 3). The outcome was a dataset that can be mined to attend the business problem. An investigation of the dataset discovered probable association between semester wise student course registration data and CGPA. The dataset was evaluated for quality problems like extreme values, incomplete values.

**Table 3, Data set without contextual factors extracted**

Student_ID	gpa	len	sem_gradepoints	sem_pass_credit	Course_Cr
2E+08	3.78	3.5	3.668		15 ACCT 101
2E+08	3.78	3.5	3.668		15 ARAB 101
2E+08	3.84	3.5	4		15 ARAB 101
2E+08	3.84	3.5	3.6		15 ACCT 101
2E+08	3.84	3.5	3.6		15 ARAB 101
2E+08	3.84	3.5	3.776667		18 ACCT 101
2E+08	2.18	4	1.553333		9 ARAB 101
2E+08	2.18	4	1.75		12 ACCT 101
2E+08	2.62	4	2.934		15 ARAB 101
2E+08	2.62	4	2.61		18 ACCT 101
2E+08	2.62	4	2.666667		9 ACCT 101
2E+08	2.62	4	2.666667		9 ARAB 101

(iii) Data Preparation: This stage concerned the choice of records of students from the dataset, clean-up, building, combining and arranges data. The resulting dataset enclosed partial information associated to enrolled student courses, semester information, CGPA and time to degree. The dataset contained 50 entries.

(iv) Modelling: This stage involved the use of modelling algorithm that enabled selection of modelling technique, generation of test design, building a model and assessment of the model.

The modeling technique chosen was classification and the chosen algorithm is Genetic Algorithm which is a heuristic search algorithm that uses cross over, mutation and fitness function to produce the course taking pattern. As part of the test design training, test and validation data sets were generated. The model was created using the following parameter setting: CGPA ≤ 4, time to degree ≥ 3.5 and ≤ 6, programme under consideration – bachelors in accounting, number of courses registered in ≥ 4 and ≤ 6 and semester = 3. The final model is a table with student records and fields (CGPA, semester and time to degree) (see Table 5) having columns student ID (student\_id), Cumulative CGPA (GPA), time to degree(len), and course taking patterns(course).

The time to degree were assessed using the course taking patterns and CGPA see table 5 which showed that time to degree can be associated courses in terms of registered number of courses and CGPA of students only.

(v) Evaluation: The results in Table 5 implied that the time to degree could be interpreted only in terms of the registered number of courses of students (without an observable pattern) for a particular CGPA, considered the highest amongst students belonging to the same programme (BSAF).



However in the absence of an observable pattern the relationship between time to degree and the registered number of courses can only be interpreted using the number of courses registered in to achieve lower time to degree as in a semester as the attribute .However it can be seen that higher the number of courses, lower will be the time to degree which is perhaps obvious regardless of the CGPA which does not require data mining. Such a result will not be very useful to the HEIs to accurately classify the courses or students to allow them to attain optimal time to degree and CGPA because of lack of additional information that could be used to link course pattern to time to degree significantly. Besides it lacks predictive power because of the obviousness associated with the outcome. Instead if a contextual factor was included in the dataset, for instance course difficulty, then it could act as missing information that could determine a enrolled course pattern of students. In this case the level of course difficulty associated with every course in the set of registered courses of students could be used to determine the optimum time to degree for a given CGPA. For instance student (stud21) has achieved the maximum CGPA of 3.88 possibly because all the six courses could be having a lower difficulty rating. However whether this actually is the case or not, cannot be confirmed due to lack of knowledge about course difficulty in the outcome produced by CRISP-DM process. This point towards the usefulness of contextual factors in discovering knowledge pertaining to course taking pattern. How this can be achieved is discussed in Section V.

(vi) Deployment: The model could not be deployed as evaluation of the outcome of CRISP-DM (see Table 5) showed that the discovered knowledge is not useful to accurately predict the course taking pattern and optimal CGPA and time to degree and classify either the students or courses which is the business goal.

At this point it can be summarized that if the business goal has to be achieved there is a need to discover hidden patterns in the dataset characterized by context which is only possible when the CRISP-DM process has steps to include contextualized dataset and mine it. Thus this research has developed an altered CRISP-DM process that could achieve the above as described next.

student_id	CGPA	CGPA	course
stud1	3.5	3.5	101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000,1001,1002,1003,1004,1005,1006,1007,1008,1009,1010,1011,1012,1013,1014,1015,1016,1017,1018,1019,1020,1021,1022,1023,1024,1025,1026,1027,1028,1029,1030,1031,1032,1033,1034,1035,1036,1037,1038,1039,1040,1041,1042,1043,1044,1045,1046,1047,1048,1049,1050,1051,1052,1053,1054,1055,1056,1057,1058,1059,1060,1061,1062,1063,1064,1065,1066,1067,1068,1069,1070,1071,1072,1073,1074,1075,1076,1077,1078,1079,1080,1081,1082,1083,1084,1085,1086,1087,1088,1089,1090,1091,1092,1093,1094,1095,1096,1097,1098,1099,1100,1101,1102,1103,1104,1105,1106,1107,1108,1109,1110,1111,1112,1113,1114,1115,1116,1117,1118,1119,1120,1121,1122,1123,1124,1125,1126,1127,1128,1129,1130,1131,1132,1133,1134,1135,1136,1137,1138,1139,1140,1141,1142,1143,1144,1145,1146,1147,1148,1149,1150,1151,1152,1153,1154,1155,1156,1157,1158,1159,1160,1161,1162,1163,1164,1165,1166,1167,1168,1169,1170,1171,1172,1173,1174,1175,1176,1177,1178,1179,1180,1181,1182,1183,1184,1185,1186,1187,1188,1189,1190,1191,1192,1193,1194,1195,1196,1197,1198,1199,1200,1201,1202,1203,1204,1205,1206,1207,1208,1209,1210,1211,1212,1213,1214,1215,1216,1217,1218,1219,1220,1221,1222,1223,1224,1225,1226,1227,1228,1229,1230,1231,1232,1233,1234,1235,1236,1237,1238,1239,1240,1241,1242,1243,1244,1245,1246,1247,1248,1249,1250,1251,1252,1253,1254,1255,1256,1257,1258,1259,1260,1261,1262,1263,1264,1265,1266,1267,1268,1269,1270,1271,1272,1273,1274,1275,1276,1277,1278,1279,1280,1281,1282,1283,1284,1285,1286,1287,1288,1289,1290,1291,1292,1293,1294,1295,1296,1297,1298,1299,1300,1301,1302,1303,1304,1305,1306,1307,1308,1309,1310,1311,1312,1313,1314,1315,1316,1317,1318,1319,1320,1321,1322,1323,1324,1325,1326,1327,1328,1329,1330,1331,1332,1333,1334,1335,1336,1337,1338,1339,1340,1341,1342,1343,1344,1345,1346,1347,1348,1349,1350,1351,1352,1353,1354,1355,1356,1357,1358,1359,1360,1361,1362,1363,1364,1365,1366,1367,1368,1369,1370,1371,1372,1373,1374,1375,1376,1377,1378,1379,1380,1381,1382,1383,1384,1385,1386,1387,1388,1389,1390,1391,1392,1393,1394,1395,1396,1397,1398,1399,1400,1401,1402,1403,1404,1405,1406,1407,1408,1409,1410,1411,1412,1413,1414,1415,1416,1417,1418,1419,1420,1421,1422,1423,1424,1425,1426,1427,1428,1429,1430,1431,1432,1433,1434,1435,1436,1437,1438,1439,1440,1441,1442,1443,1444,1445,1446,1447,1448,1449,1450,1451,1452,1453,1454,1455,1456,1457,1458,1459,1460,1461,1462,1463,1464,1465,1466,1467,1468,1469,1470,1471,1472,1473,1474,1475,1476,1477,1478,1479,1480,1481,1482,1483,1484,1485,1486,1487,1488,1489,1490,1491,1492,1493,1494,1495,1496,1497,1498,1499,1500,1501,1502,1503,1504,1505,1506,1507,1508,1509,1510,1511,1512,1513,1514,1515,1516,1517,1518,1519,1520,1521,1522,1523,1524,1525,1526,1527,1528,1529,1530,1531,1532,1533,1534,1535,1536,1537,1538,1539,1540,1541,1542,1543,1544,1545,1546,1547,1548,1549,1550,1551,1552,1553,1554,1555,1556,1557,1558,1559,1560,1561,1562,1563,1564,1565,1566,1567,1568,1569,1570,1571,1572,1573,1574,1575,1576,1577,1578,1579,1580,1581,1582,1583,1584,1585,1586,1587,1588,1589,1590,1591,1592,1593,1594,1595,1596,1597,1598,1599,1600,1601,1602,1603,1604,1605,1606,1607,1608,1609,1610,1611,1612,1613,1614,1615,1616,1617,1618,1619,1620,1621,1622,1623,1624,1625,1626,1627,1628,1629,1630,1631,1632,1633,1634,1635,1636,1637,1638,1639,1640,1641,1642,1643,1644,1645,1646,1647,1648,1649,1650,1651,1652,1653,1654,1655,1656,1657,1658,1659,1660,1661,1662,1663,1664,1665,1666,1667,1668,1669,1670,1671,1672,1673,1674,1675,1676,1677,1678,1679,1680,1681,1682,1683,1684,1685,1686,1687,1688,1689,1690,1691,1692,1693,1694,1695,1696,1697,1698,1699,1700,1701,1702,1703,1704,1705,1706,1707,1708,1709,1710,1711,1712,1713,1714,1715,1716,1717,1718,1719,1720,1721,1722,1723,1724,1725,1726,1727,1728,1729,1730,1731,1732,1733,1734,1735,1736,1737,1738,1739,1740,1741,1742,1743,1744,1745,1746,1747,1748,1749,1750,1751,1752,1753,1754,1755,1756,1757,1758,1759,1760,1761,1762,1763,1764,1765,1766,1767,1768,1769,1770,1771,1772,1773,1774,1775,1776,1777,1778,1779,1780,1781,1782,1783,1784,1785,1786,1787,1788,1789,1790,1791,1792,1793,1794,1795,1796,1797,1798,1799,1800,1801,1802,1803,1804,1805,1806,1807,1808,1809,1810,1811,1812,1813,1814,1815,1816,1817,1818,1819,1820,1821,1822,1823,1824,1825,1826,1827,1828,1829,1830,1831,1832,1833,1834,1835,1836,1837,1838,1839,1840,1841,1842,1843,1844,1845,1846,1847,1848,1849,1850,1851,1852,1853,1854,1855,1856,1857,1858,1859,1860,1861,1862,1863,1864,1865,1866,1867,1868,1869,1870,1871,1872,1873,1874,1875,1876,1877,1878,1879,1880,1881,1882,1883,1884,1885,1886,1887,1888,1889,1890,1891,1892,1893,1894,1895,1896,1897,1898,1899,1900,1901,1902,1903,1904,1905,1906,1907,1908,1909,1910,1911,1912,1913,1914,1915,1916,1917,1918,1919,1920,1921,1922,1923,1924,1925,1926,1927,1928,1929,1930,1931,1932,1933,1934,1935,1936,1937,1938,1939,1940,1941,1942,1943,1944,1945,1946,1947,1948,1949,1950,1951,1952,1953,1954,1955,1956,1957,1958,1959,1960,1961,1962,1963,1964,1965,1966,1967,1968,1969,1970,1971,1972,1973,1974,1975,1976,1977,1978,1979,1980,1981,1982,1983,1984,1985,1986,1987,1988,1989,1990,1991,1992,1993,1994,1995,1996,1997,1998,1999,2000,2001,2002,2003,2004,2005,2006,2007,2008,2009,2010,2011,2012,2013,2014,2015,2016,2017,2018,2019,2020,2021,2022,2023,2024,2025,2026,2027,2028,2029,2030,2031,2032,2033,2034,2035,2036,2037,2038,2039,2040,2041,2042,2043,2044,2045,2046,2047,2048,2049,2050,2051,2052,2053,2054,2055,2056,2057,2058,2059,2060,2061,2062,2063,2064,2065,2066,2067,2068,2069,2070,2071,2072,2073,2074,2075,2076,2077,2078,2079,2080,2081,2082,2083,2084,2085,2086,2087,2088,2089,2090,2091,2092,2093,2094,2095,2096,2097,2098,2099,2100,2101,2102,2103,2104,2105,2106,2107,2108,2109,2110,2111,2112,2113,2114,2115,2116,2117,2118,2119,2120,2121,2122,2123,2124,2125,2126,2127,2128,2129,2130,2131,2132,2133,2134,2135,2136,2137,2138,2139,2140,2141,2142,2143,2144,2145,2146,2147,2148,2149,2150,2151,2152,2153,2154,2155,2156,2157,2158,2159,2160,2161,2162,2163,2164,2165,2166,2167,2168,2169,2170,2171,2172,2173,2174,2175,2176,2177,2178,2179,2180,2181,2182,2183,2184,2185,2186,2187,2188,2189,2190,2191,2192,2193,2194,2195,2196,2197,2198,2199,2200,2201,2202,2203,2204,2205,2206,2207,2208,2209,2210,2211,2212,2213,2214,2215,2216,2217,2218,2219,2220,2221,2222,2223,2224,2225,2226,2227,2228,2229,2230,2231,2232,2233,2234,2235,2236,2237,2238,2239,2240,2241,2242,2243,2244,2245,2246,2247,2248,2249,2250,2251,2252,2253,2254,2255,2256,2257,2258,2259,2260,2261,2262,2263,2264,2265,2266,2267,2268,2269,2270,2271,2272,2273,2274,2275,2276,2277,2278,2279,2280,2281,2282,2283,2284,2285,2286,2287,2288,2289,2290,2291,2292,2293,2294,2295,2296,2297,2298,2299,2300,2301,2302,2303,2304,2305,2306,2307,2308,2309,2310,2311,2312,2313,2314,2315,2316,2317,2318,2319,2320,2321,2322,2323,2324,2325,2326,2327,2328,2329,2330,2331,2332,2333,2334,2335,2336,2337,2338,2339,2340,2341,2342,2343,2344,2345,2346,2347,2348,2349,2350,2351,2352,2353,2354,2355,2356,2357,2358,2359,2360,2361,2362,2363,2364,2365,2366,2367,2368,2369,2370,2371,2372,2373,2374,2375,2376,2377,2378,2379,2380,2381,2382,2383,2384,2385,2386,2387,2388,2389,2390,2391,2392,2393,2394,2395,2396,2397,2398,2399,2400,2401,2402,2403,2404,2405,2406,2407,2408,2409,2410,2411,2412,2413,2414,2415,2416,2417,2418,2419,2420,2421,2422,2423,2424,2425,2426,2427,2428,2429,2430,2431,2432,2433,2434,2435,2436,2437,2438,2439,2440,2441,2442,2443,2444,2445,2446,2447,2448,2449,2450,2451,2452,2453,2454,2455,2456,2457,2458,2459,2460,2461,2462,2463,2464,2465,2466,2467,2468,2469,2470,2471,2472,2473,2474,2475,2476,2477,2478,2479,2480,2481,2482,2483,2484,2485,2486,2487,2488,2489,2490,2491,2492,2493,2494,2495,2496,2497,2498,2499,2500,2501,2502,2503,2504,2505,2506,2507,2508,2509,2510,2511,2512,2513,2514,2515,2516,2517,2518,2519,2520,2521,2522,2523,2524,2525,2526,2527,2528,2529,2530,2531,2532,2533,2534,2535,2536,2537,2538,2539,2540,2541,2542,2543,2544,2545,2546,2547,2548,2549,2550,2551,2552,2553,2554,2555,2556,2557,2558,2559,2560,2561,2562,2563,2564,2565,2566,2567,2568,2569,2570,2571,2572,2573,2574,2575,2576,2577,2578,2579,2580,2581,2582,2583,2584,2585,2586,2587,2588,2589,2590,2591,2592,2593,2594,2595,2596,2597,2598,2599,2600,2601,2602,2603,2604,2605,2606,2607,2608,2609,2610,2611,2612,2613,2614,2615,2616,2617,2618,2619,2620,2621,2

- c) Data Preparation (Step C): Like in Step B this step also have been divided into 2 subsections which step C1 and C2 with respect to subsection C of Section IV. While step C1 deals with contextual data namely course difficulty, in this step data preparation is conducted in similar line as described in subsection (iii) of section IV. Again step C2 is implemented in the same way as described in sub section (iii) of section IV. It can be seen in Figure 7 that the step B1 is linked to step C1 and step B2 is linked to step C2 in a reversible fashion which is another modification of the original CRISP-DM process.
- d) Added Data Preparation (Step D): The outcomes of this step are a dataset that has extra information about course difficulty. This step is not originally part of the CRISP-DM process developed Chapman et. al, 2000. In this step the outcomes of steps C1 and C2 are merged into a single dataset that will be mined further. This is a major modification of the original CRISP-DM process.
- e) Contextualized Modelling (Step E): This step although very similar to subsection (iv) of section IV differs from that step significantly as it has been modified to handle contextual data using the additional link provided to connect with step C1 and C2 in the reverse direction as well as using a modified Genetic Algorithm that uses course difficulty data. In addition the model was created The model was created using the following parameter setting:  $CGPA \leq 4$ , time to degree  $\geq 3$  and  $\leq 6$ , programme under consideration – bachelors in accounting and finance (BSAF), number of courses registered in  $\geq 4$  but  $\leq 6$ , semester = 3 and course difficulty which is measured using a 5 point scale (very difficult, difficult, average, easy and very easy). The final model expected to be produced is a table with student records and fields CGPA and time to degree) (see Table 6) which has columns student identification (student\_id), time to degree(len), Cumulative CGPA (GPA), course difficulty(difficulty) and course taking patterns. The time to degree were assessed using the course taking patterns and CGPA see table 5 which showed that time to degree can be associated courses in terms of registered number of courses and CGPA of students only.
- f) Evaluation (Step F): Time to degree could be interpreted in terms of the registered number of courses of students, with an observable pattern that could be linked to course difficulty for a particular CGPA, considered the highest amongst students belonging to the same programme (BSAF). The use of course difficulty lies in enabling the interpretation of the most optimum course taking pattern

of students that could lead to achieve the optimum time to degree and not by mere number of courses. For instance a student could register in six courses all of which could fall under the category of difficult and yet achieve a time to degree of 3 and the highest CGPA. On the other hand a student could register in six courses that could fall under the category of course difficulty “easy” yet achieve a time to degree higher than 3 and lower CGPA. In this case it is possible to interpret that despite registering in the maximum number of courses in a semester, with the difficulty levels measured as “difficult”, a student can achieve a lower time to degree and highest CGPA. Such an interpretation could help in identifying and categorizing students who are capable of achieving lower time to degree regardless of the fact that they have registered in courses whose difficulty level could be higher and score high CGPA. The same argument could be used to categorize other students depending on the course taking pattern and the difficulty levels of those courses they have registered in but have taken longer time to degree and scored lower CGPA. This process could be continued to categorize students and the courses they have registered in to determine the level academic support that needs to be provided to enable them achieve optimum time to degree and high CGPA. Thus it can be seen that the business goals of categorizing courses and students to predict optimum time to degree and highest CGPA could be achieved accurately using the contextual factors and a modified process which is demonstrated practically in Section VI.

g) Deployment (Step G): The deployment of the modified CRISP-DM process involves an analysis and understanding of the results produced by the process and at this stage a decision to deploy cannot be taken. Hence this section has been addressed in Section VI where results of the analysis can be used to determine the deployment.

In summary it can be seen that the newly designed modified CRISP-DM process has the potential to achieve the business goal of predicting reasonably accurately the optimum time to degree and CGPA in terms of course taking pattern of students and course difficulty leading to decision making about categorizing courses and students and provide better support to students. To verify this aspect the next section provides the details of the experiments conducted and analysis of the results.

Table 6, Course-taking patterns extracted from Dataset with Contextual Attributes mined from Student Data

Student ID	Time to degree	CGPA	SGPA	difficulty	course
stud1	4.5	3.33	2.982	Difficult,Difficult,Difficult,Average,Difficult	ACCT 301,FINC 310,FINC 321,FREN 101,STAT 202
stud2	4.5	3.06	2.4	Difficult,Difficult,Difficult,Average,Difficult	ACCT 301,FINC 310,FINC 321,FREN 101,STAT 202
stud3	5	2.44	2.666	Average,Difficult,Average,Difficult,Difficult	ACCT 312,BANK 220,CULT 101,ENGL 202,ITMA 201
stud4	6	2.28	2.4	Average,Difficult,Average,Difficult,Difficult	ACCT 312,BANK 220,CULT 101,ENGL 202,ITMA 201
stud5	4.5	3.1	3.334	Difficult,Average,Difficult,Difficult,Difficult	ACCT 311,ACCT 321,ENGL 202,FINC 321,STAT 202
stud6	4.5	3.49	3.468	Difficult,Average,Difficult,Difficult,Difficult	ACCT 311,ACCT 321,ENGL 202,FINC 321,STAT 202
stud7	4.5	3.75	3.6	Difficult,Difficult,Easy,Difficult,Difficult	ACCT 301,ACCT 311,ARAB 102,FINC 421,ITMA 201
stud8	4.5	3.41	2.8	Difficult,Difficult,Easy,Difficult,Difficult	ACCT 301,ACCT 311,ARAB 102,FINC 421,ITMA 201
stud9	4	3.4	3.668333	Difficult,Difficult,Average,Easy	ACCT 321,ARAB 201,BANK 220,FINC 431,ITCS 121,PHOT 101
stud10	4	3.46	3.723333	Difficult,Difficult,Average,Easy	ACCT 321,ARAB 201,BANK 220,FINC 431,ITCS 121,PHOT 101
stud11	4.5	2.55	2.168333	Difficult,Average,Difficult,Easy,Difficult,Easy	ACCT 321,ACCT 402,BANK 302,ENGL 201,FINC 421,PHOT 101
stud12	4.5	2.33	2.333333	Difficult,Average,Difficult,Easy,Difficult,Easy	ACCT 321,ACCT 402,BANK 302,ENGL 201,FINC 421,PHOT 101
stud13	5	2.5	1	Average,Difficult,Easy,Difficult,Difficult	ACCT 403,BANK 302,ENGL 201,FINC 320,FINC 421
stud14	5.5	2.36	0.8	Average,Difficult,Easy,Difficult,Difficult	ACCT 403,BANK 302,ENGL 201,FINC 320,FINC 421
stud15	4	3.57	3.732	Difficult,Average,Difficult,Difficult,Average	ACCT 301,CULT 102,ENGL 202,FINC 320,ITCS 121
stud16	4	3.84	3.666	Difficult,Average,Difficult,Difficult,Average	ACCT 301,CULT 102,ENGL 202,FINC 320,ITCS 121
stud17	4	3.6	3.984	Average,Average,Difficult,Easy,Easy	ACCT 402,ACCT 403,BANK 302,ENGL 201,VIDEO 101
stud18	4	3.22	3.202	Average,Average,Difficult,Easy,Easy	ACCT 402,ACCT 403,BANK 302,ENGL 201,VIDEO 101
stud19	4	3.42	3.338	Average,Average,Easy,Difficult,Difficult	ACCT 312,ACCT 320,ACCT 341,BANK 302,FINC 310
stud20	4	3.23	3.2	Average,Average,Easy,Difficult,Difficult	ACCT 312,ACCT 320,ACCT 341,BANK 302,FINC 310
stud21	3	3.88	3.778333	Difficult,Difficult,Difficult,Difficult,Difficult,Difficult	ACCT 404,BANK 302,ECON 301,FINC 320,FINC 421,STAT 202
stud22	3	3.53	3.556667	Difficult,Difficult,Difficult,Difficult,Difficult,Difficult	ACCT 404,BANK 302,ECON 301,FINC 320,FINC 421,STAT 202

The input to the KDDM process is the dataset characterised by contextual factor course difficulty. The business goals to be achieved are determining the optimum time to degree and CGPA, categorisation of courses that could predict optimum time to degree and CGPA using course registration patterns and course difficulty. Categorisation of students based on optimum time to degree, CGPA, course patterns and course difficulty. Data preparation and data understanding stage covers the data quality factors of data. For detailed steps, readers have to refer to Chapman et al. (2000). Genetic algorithm was used with knowledge of contextual factors course difficulty. Model evaluation stage comprises of pattern extracted through the process with and without contextual detection stage. The overall CRISP-DM process has not been discussed in detail and the discussion here is limited to demonstrating the difference between the knowledge discovered using patterns extracted through CRISP-DM process with and without the detection of course difficulty and semester as contextual factor. As far as the introduction of contextualised dataset was concerned initially the characteristics of course difficulty pertaining to every course under consideration in a particular semester 3 were utilised which included 'course GPA', 'weighted average of the course GPA', 'current course', 'course equivalent to the current one', 'Set of equivalence courses for course', 'Total number of students in the course'. Using these characteristics pertaining to course difficulty a scale was developed to signify the level of difficulty of each course. The course difficulty was measured over 5 point scale very easy, easy, average, difficult and very difficult.

### I. EXPERIMENTS ON MODIFIED CRISP-DM PROCESS AND ANALYSIS OF RESULTS(STAGE 3)

This section describes experiments conducted on modified CRISP-DM process given in figure 7. The subsection (a1) to

(g1) in this section has one to one correspondence to the subsection A to G in section V. Thus each subsection provides the experimental and analytical details conducted according to the design specification provided in the subsection A to G in section V.

a1) Business understanding (Step A): predicting reasonably accurately the optimum time to degree and CGPA in terms of course taking pattern of students and course difficulty leading to decision making about categorising courses and students and provide better support to students.

b1) Data Understanding (Step B): This step is divided into two: one addressing contextual data and the other addressing general data.

b1.1 Finding the existence or nonexistence of course difficulty using a specially designed algorithm that was based on the formula devised by Zainuddin, 2012 to compute the course difficulty level (refer Zainuddin, 2012) and a pseudo code written based on the guideline given by Vert et al., 2010.

The outputs derived from the execution of the pseudocode is shown in the screenshot below.

The screenshot shows a data table with columns: Student\_ID, gpa, len, sem\_gradepts, sem\_pass\_cred, Course\_Code, and Semester. The data is filtered for Student\_ID 200300087. A dialog box is overlaid on the table, stating "Data 1 has 1 contextual variable" with an "OK" button.

Student_ID	gpa	len	sem_gradepts	sem_pass_cred	Course_Code	Semester
200300087	3.78	3.5	3.668	15	ACCT 101	1
200300087	3.78	3.5	3.668	15	ARAB 101	1
200300087	3.78	3.5	3.668	15	ECON 101	1
200300087	3.78	3.5	3.668	15	ENGL 050	1
200300087	3.78	3.5	3.668	15		
200300087	3.78	3.5	3.668	15		
200300087	3.78	3.5	3.668	15		
200300087	3.78	3.5	3.668	15		
200300087	3.78	3.5	3.734	15		
200300087	3.78	3.5	3.734	15		
200300087	3.78	3.5	3.734	15	FREN 101	2
200300087	3.78	3.5	3.734	15	MAGT 221	2
200300087	3.78	3.5	3.734	15	STAT 101	2

student_id	len	prevgpa	gpa	semester	sem_grade	course
200300056	3.5	95.3	3.48	1	3.2475	BANK 2
200300056	3.5	95.3	3.48	1	3.2475	ECON 1
200300056	3.5	95.3	3.48	1	3.2475	ITMA 20
200300056	3.5	95.3	3.48	1	3.2475	MAGT 2
200300056	3.5	95.3	3.48	2	3.75	BANK 3
200300056	3.5	95.3	3.48	2	3.75	CULT 11
200300056	3.5	95.3	3.48	2	3.75	ECON 3
200300056	3.5	95.3	3.48	2	3.75	FINC 32
200300056	3.5	95.3	3.48	3	3.8325	ACCT 41
200300056	3.5	95.3	3.48	3	3.8325	ACCT 41
200300056	3.5	95.3	3.48	3	3.8325	CULT 11
200300056	3.5	95.3	3.48	3	3.8325	ECON 4
200300056	3.5	95.3	3.48	4	2.75	ACCT 11

student_id	len	prevgpa	gpa	semester	sem_grade	course_code	difficulty
2E+08	3.5	95.3	3.48	4	2.75	ACCT 101	Very Difficult
2E+08	3.5	95.3	3.48	4	2.75	ARAB 101	Difficult
2E+08	4	97.6	3.87	1	3.868	ACCT 101	Very Difficult
2E+08	4	97.6	3.87	2	3.466	ACCT 101	Very Difficult
2E+08	4	97.6	3.87	2	3.466	ARAB 101	Difficult
2E+08	4	97.6	3.87	5	3.868	ARAB 101	Difficult
2E+08	3.5	95.4	3.78	1	3.668	ACCT 101	Very Difficult
2E+08	3.5	95.4	3.78	1	3.668	ARAB 101	Difficult
2E+08	3.5	96.2	3.84	1	4	ARAB 101	Difficult
2E+08	3.5	96.2	3.84	2	3.6	ACCT 101	Very Difficult
2E+08	3.5	96.2	3.84	2	3.6	ARAB 101	Difficult
2E+08	3.5	96.2	3.84	5	3.776667	ACCT 101	Very Difficult
2E+08	4	77.3	2.18	1	1.553333	ARAB 101	Difficult

b1.2 General Data

This section is same as previous sections as it deals with general data of students. The data extracted using the SQL queries is provided in Table 3.

c1) Data Preparation (Step C): Both the outputs of b1.1 and b1.2 were subjected to the same steps mentioned subsection iii of section IV.

d1) Added Data Preparation (Step D): The data prepared in step C1 and C2 see figure 7 are merged and the resulting dataset is shown in Table 4. This dataset is further fed into the modeling stage.

Table 4, Data set with contextual factors extracted

Table 7, Course taking patterns extracted from Dataset with Contextual Factors Extracted from Student Data (KJ Lu, 2016).

Student ID	Time to degree	CGPA	SGPA	difficulty	course
stud1	4.5	3.33	2.982	Difficult,Difficult,Difficult,Average,Difficult	ACCT 301,FINC 310,FINC 321,FREN 101,STAT 202
stud2	4.5	3.06	2.4	Difficult,Difficult,Difficult,Average,Difficult	ACCT 301,FINC 310,FINC 321,FREN 101,STAT 202
stud3	5	2.44	2.666	Average,Difficult,Average,Difficult,Difficult	ACCT 312,BANK 220,CULT 101,ENGL 202,ITMA 201
stud4	6	2.28	2.4	Average,Difficult,Average,Difficult,Difficult	ACCT 312,BANK 220,CULT 101,ENGL 202,ITMA 201
stud5	4.5	3.1	3.334	Difficult,Average,Difficult,Difficult,Difficult	ACCT 311,ACCT 321,ENGL 202,FINC 321,STAT 202
stud6	4.5	3.49	3.468	Difficult,Average,Difficult,Difficult,Difficult	ACCT 311,ACCT 321,ENGL 202,FINC 321,STAT 202
stud7	4.5	3.75	3.6	Difficult,Difficult,Easy,Difficult,Difficult	ACCT 301,ACCT 311,ARAB 102,FINC 421,ITMA 201
stud8	4.5	3.41	2.8	Difficult,Difficult,Easy,Difficult,Difficult	ACCT 301,ACCT 311,ARAB 102,FINC 421,ITMA 201
stud9	4	3.4	3.668333	Difficult,Difficult,Average,Easy	ACCT 321,ARAB 201,BANK 220,FINC 431,ITCS 121,PHOT 101
stud10	4	3.46	3.723333	Difficult,Difficult,Average,Easy	ACCT 321,ARAB 201,BANK 220,FINC 431,ITCS 121,PHOT 101
stud11	4.5	2.55	2.168333	Difficult,Average,Difficult,Easy,Difficult,Easy	ACCT 321,ACCT 402,BANK 302,ENGL 201,FINC 421,PHOT 101
stud12	4.5	2.33	2.333333	Difficult,Average,Difficult,Easy,Difficult,Easy	ACCT 321,ACCT 402,BANK 302,ENGL 201,FINC 421,PHOT 101
stud13	5	2.5	1	Average,Difficult,Easy,Difficult,Difficult	ACCT 403,BANK 302,ENGL 201,FINC 320,FINC 421
stud14	5.5	2.36	0.8	Average,Difficult,Easy,Difficult,Difficult	ACCT 403,BANK 302,ENGL 201,FINC 320,FINC 421
stud15	4	3.57	3.732	Difficult,Average,Difficult,Difficult,Average	ACCT 301,CULT 102,ENGL 202,FINC 320,ITCS 121
stud16	4	3.84	3.666	Difficult,Average,Difficult,Difficult,Average	ACCT 301,CULT 102,ENGL 202,FINC 320,ITCS 121
stud17	4	3.6	3.984	Average,Average,Difficult,Easy,Easy	ACCT 402,ACCT 403,BANK 302,ENGL 201,VIDEO 101
stud18	4	3.22	3.202	Average,Average,Difficult,Easy,Easy	ACCT 402,ACCT 403,BANK 302,ENGL 201,VIDEO 101
stud19	4	3.42	3.398	Average,Average,Easy,Difficult,Difficult	ACCT 312,ACCT 320,ACCT 341,BANK 302,FINC 310
stud20	4	3.23	3.2	Average,Average,Easy,Difficult,Difficult	ACCT 312,ACCT 320,ACCT 341,BANK 302,FINC 310
stud21	3	3.88	3.778333	Difficult,Difficult,Difficult,Difficult,Difficult,Difficult	ACCT 404,BANK 302,ECON 301,FINC 320,FINC 421,STAT 202
stud22	3	3.53	3.556667	Difficult,Difficult,Difficult,Difficult,Difficult,Difficult	ACCT 404,BANK 302,ECON 301,FINC 320,FINC 421,STAT 202

f1) Evaluation (Step F): In line with the details of the design given in sub-section (f) of Section V the results obtained were evaluated. Evaluation included exploring the course taking pattern by allocating the course difficulty level for each enrolled course and links the course pattern to time to degree and assess the scored CGPA (Table 7). The course difficulty level was calculated using a five point scale (see Section V). The conjecture here is that more the number of enrolled courses of a student having difficulty level is greater

e1) Modelling (Step E): This step provides the model generated at step E see of the modified CRISP-DM process using the experimental details provided in subsection e in section V. The model is a report is shown in Table 7 (KJ Lu, 2016).

than or equal to 'difficult' then more would be the time to degree and lesser will be the CGPA. For instance, take a student e.g. stud 6 in Table 7 who is shown to have registered in five courses of these the course difficulty is 'Difficult' with the left over 1 course as 'Average'.



The student CGPA is 3.49. It is sensible to anticipate that this student will take more time to graduate as he was able to finish the entire number of 44 courses in 4.5 years on an average of five courses in a semester and a CGPA around the same figure of 3.49. The course difficulty stage has possibly not allowed the student to enrol in more courses in a semester, say 6. Had 'stud6' enrolled in 6 courses/semester, then the student might have graduated in 4 years. Thus course taking pattern and the corresponding course difficulty measure of a course can be believed to have an effect on the time to degree. On the other hand if 'stud6' had registered in 6 courses whose course difficulty scores were a combination of 'Easy', 'Average' and 'Difficult' (for instance, a course taking pattern of 6 courses may be ARAB 102, ACCT 320, ENGL 201, BANK 302, STAT 202 with course difficulty measured as 'Easy', 'Average', 'Easy', 'Difficult', 'Difficult' then there is a probability that 'stud6' might have scored the same CGPA of 3.49 or more and had taken lesser time to degree (4 years). Therefore it can be witnessed that course taking pattern if linked with the course difficulty then it is probable to forecast the optimal time to degree at an optimal CGPA. That is to say, it is possible to anticipate that when students enrol in courses where difficulty level is "difficult" then lesser would be the number of such courses a student could enrol owing to harder effort to complete with a higher CGPA. From this situation, it is probable to deduce that when students enrol in a set of courses in which difficulty levels are "difficult", then the no. of courses in that set is liable to be lesser. This might be because, more the course difficulty level, difficult is the attempt needed to finish the course with higher CGPA. Therefore lower number of courses in the set provides greater opportunity to work hard leading to better chance of scoring higher CGPA at the cost of time to degree. That is:

f11) Course difficulty negatively influences the number of courses registered in a semester.

f12) Course taking pattern characterized by high number of difficult courses negatively influences CGPA and positively influences time to degree.

f13) Course taking pattern characterized by low number of difficult courses positively influences CGPA and negatively influences time to degree.

However this inference was not so obvious to predict and generalize when one considers the hidden knowledge in the dataset and results provided Table 7. Entirely opposite results have been obtained in the case two students namely 'stud21' and 'stud22'. From Table 7 it can be seen that the statements in (f11), (f12) and (f13) are completely falsified. The evaluation of the results of the students 'stud21' and 'stud22' clearly shows:

f14) Course difficulty does not influence the number of courses registered in a semester. The students 'stud21' and 'stud22' have registered in a maximum of six courses all of them rated as 'difficult'.

f15) Course taking pattern characterized by high number of difficult courses does not influence CGPA and negatively influences time to degree. It can be seen that the pattern of courses registered in by the student 'stud21' although comprises six courses whose level of difficulty is rated as 'difficult' has the scored the maximum CGPA and achieved the lowest time to degree of 3 years. Similarly for the student 'stud22' the course taking pattern although comprising of six courses whose level of difficulty is rated as 'difficult' has not

scored the maximum CGPA but achieved the lowest time to degree of 3 years. This shows that the course taking pattern characterized by high number of difficult courses does not influence CGPA but negatively influences time to degree.

f16) Course taking pattern characterized by low number of difficult courses positively influences CGPA and negatively influences time to degree is falsified because high number of course rated as difficult can also positively influence CGPA and negatively influence time to degree as mentioned in (f15). The previous findings and discussions evidently point out that the optimal time to degree and the maximum CGPA can be fixed by a set of six courses (highest authorized by a HEI in a semester, refer section III), with every course fixed as 'Difficult'. When this finding is treated as optimal, in that case the following decision could be facilitated in the HEIs.

1. Categorizing students akin to 'stud21' and 'stud22' to allow them to graduate in 3 years.
2. Categorizing students akin to 'stud6' who have enrolled in just 5 courses and investigate their performance in terms of their ability to score the same CGPA or higher to enable them to achieve a performance similar to 'stud21' and 'stud22'.
3. Categorizing students like 'stud17' and 'stud18' who have taken 4 courses and analyse their performance in terms of their ability to score the same CGPA or higher to enable them to achieve a performance similar to 'stud21' and 'stud22'.
4. Assert the performance of students who have enrolled in 4 or 5 courses in a semester at an early stage in their academic career in the university using modified CRISP-DM to forecast their optimal time to degree with higher CGPA based on past student results.

g1) Deployment (Step G): The deployment of the modified CRISP-DM process involves an analysis and understanding of the results produced by the process. From the discussions provided in above section f1 it is clear that the modified CRISP-DM process can enable HEIs to achieve business goals using the discovered knowledge in terms of course taking patterns and the context of course difficulty to determine optimum time to degree and CGPA. The results obtained are consistent and hence the tailored CRISP-DM process can be deployed through teaching and learning and achieve the business goal of supporting students to achieve to optimize the time to degree and CGPA using course taking patterns and course difficulty as parameters. The tests deployed in modifying CRISP-DM process are very similar to original CRISP-DM process.

## I. CONCLUSION, LIMITATION AND FUTURE RESEARCH

The results obtained in this research as explained in Sections V and IV clearly point out to that the main objective of demonstrating a modified KDDM process that extracts hidden knowledge from the student dataset characterized by contextual factors that could be useful in making better decisions in HEIs to achieve specific business goals has been achieved.



The modified process uses CRISP-DM as an example which has not been deployed in the context HEIs. Thus the outcome of this research also provides a demonstration of the use of CRISP-DM process in HEIs. Finally the modified CRISP-DM process provides a DM process that could uncover contextual data to be used to support business goal and produce a dataset at the preparation stage which is contextualized. This provides a new way of dealing with the model development at the DM stage in the CRISP-DM process. At this stage a modified genetic algorithm has been specifically developed to enable the CRISP-DM process to generate a model that is contextual leading to the discovery of course taking patterns that are contextualized. This discovery has enables prediction of optimum time to degree and CGPA, an aspect not covered in the extant literature. As far as limitations of the research are concerned, it must be pointed out that this research has used only one contextual factor namely course difficulty of the registered courses of just 25 students in only one semester. There is a need to cross-check this with more number of semesters and additional contextual factors for greater number of students. Future research could investigate into more number of semesters and use additional contextual factors including student potential, course weight age and course complexity which promise to enable decision makers to make more accurate decisions.

## REFERENCES

- Adelman, C., 1999. Answers in the toolbox: Academic Intensity, attendance patterns, and Bachelor's Degree attainment. Washington DC, Department of Education, Office of Educational Research and Improvement.
- Astin, A., 1971. Predicting academic performance in college: Selectivity data for 2300 American colleges.. New York: The Free Press.
- Athiyaman, A., 1997. Linking student satisfaction and service quality perceptions: The case of university education. *European Journal of Marketing*, 31(7), pp. 528-540.
- Bahr, P. R., 2010. The bird's eye view of community colleges: A behavioral typology of first-time students based on cluster analytic classification.. *Research in Higher Education*, 51(8), pp. 724-749.
- Bolchini, C. et al., 2007. A dataoriented survey of context models. *SIGMOD Rec.(ACM)* doi:10.1145/1361348.1361353. ISSN, 36(4), pp. 19-26.
- Brachman, R. & Anand, T., 1996. The process of knowledge discovery in databases: A human-centered approach. In U. Fayyad, G. Piatetsky-Shapiro, P. Smith, & R. Uthuruswamy (Eds.), *Advances in knowledge discovery and data mining*. AAAI Press, pp. 36-57.
- Catley, C., Smith, K., McGregor, C. & Tracy, M., 2009. Extending CRISP-DM to Incorporate Temporal Data Mining of Multi dimensional Medical Data Streams: A Neonatal Intensive Care Unit Case Study. s.l., *Computer-Based Medical Systems*, 2009. CBMS 2009. 22nd IEEE International Symposium.
- Chapman, P., Clinton, J., Kerber, R. & Khabaza, T., 2000. CRISPDM 1.0 step-by-step data mining guide. Technical report, CRISP-DM , s.l.: CRISP-DM.
- Daempfle, P. A., 2003. An Analysis of the High Attrition Rates Among First Year College Science, Math, and Engineering Majors. *Journal of College Student Retention*, 5(1), pp. 37-52.
- Davenport, T. . H., 2010. The New World of "Business Analytics", s.l.: International Institute for Analytics.
- DeShields, O. W., Kara, A. & Kaynak, E., 2005. Determinants of business student satisfaction and retention in higher education: Applying Herzberg's two-factor theory.. *International Journal of Educational Management*, 19(2), pp. 128-139.
- Dey, A. K., 2001. Understanding and Using Context. *Personal Ubiquitous Comput.*, 5(1), pp. 4-7.
- Elliott, K. M. & Healy, M. A., 2001. Key factors influencing student satisfaction related to recruitment retention. *Journal of Marketing for Higher Education*, 10(4), pp. 1-11.
- Fayyad, U., Piatetsky-Shapiro, G., Smyth, P. & Uthurusamy, R., 1996. *Advances in knowledge discovery and data mining*.. MIT Press..
- Fettke, P., Vella, A. L. & Loos, P., 2012. From Measuring the Quality of Labels in Process Models to a Discourse on Process Model Quality: A Case Study. Maui, HI , IEEE.
- Helgesen, O. & Nettet, E., 2007. What accounts for students' loyalty? Some field study evidence.. *International Journal of Educational Management*, 21(2), pp. 126-143.
- Hevner, A. R., March, S. T., Park, J., & Ram, S. (2004). Design science in information systems research. *MIS Quarterly*, 28(1), 75-105.
- John, G. H., 1997. Enhancements to the data mining process. s.l.: PhD thesis, Stanford University.
- Knight, W. E., 1994. Why the five-year (or longer) bachelors degree ? : An exploratory study of time to degree attainment. New Orleans, LA, Association for Institutional Research forum.
- Kovacic, Z. J., 2010. Early prediction of student success: Mining student enrollment data. s.l., *Proceedings of Informing Science & IT Education Conference* .
- Kurgan, L. & Musilek, P., 2006. A survey of knowledge discovery and data mining process models.. *Knowledge Engineering Review*, 21(1), pp. 1-24.
- Li, J., Yang, B. & Song, W., 2009. A New Data Mining Process Model for Aluminum Electrolysis. Qingdao, P. R. China, *Proceedings of the International Symposium on Intelligent Information Systems and Applications (IISA'09)*.
- Lotkowski, V. A., Robbins, S. B. & Noeth, R. J., 2004. The Role of Academic and Non-Academic Factors in Improving College Retention, Iowa City, IA: ACT Policy Report.
- Marbán, Ó., Mariscal, G. & Segovia, J., 2009. A Data Mining & Knowledge Discovery Process Model. I-Tech, Vienna, Austria: *Data Mining and Knowledge Discovery in Real Life Applications*, Julio Ponce and Adem Karahoca (Ed.), ISBN: 978-3-902613-53-0.
- Minaei-Bigdoli, B., Kashy, D. A., Kortemeyer, G. & Punch, W. F., 2003. 33rd ASEE/IEEE Frontiers in Education Conference. Boulder, CO, IEEE.
- Pheng, L. S. & Arain, F. M., 2006. A KNOWLEDGE-BASED SYSTEM AS A DECISION MAKING TOOL FOR EFFECTIVE MANAGEMENT OF VARIATIONS AND DESIGN IMPROVEMENT: LEVERAGING ON INFORMATION TECHNOLOGY APPLICATIONS. *ITcon*, Volume 11.
- Rao, K., Govardhan, A. & Rao, K., 2012. SPATIOTEMPORAL DATA MINING: ISSUES, TASKS AND APPLICATIONS. *International Journal of Computer Science & Engineering Survey (IJCSES)*, 3(1).
- Redpath, R. & Srinivasan, B., 2004. A Model for Domain Centered Knowledge Discovery in Databases. Budapest, Hungary, *Proceedings of the IEEE 4th International Conference On Intelligent Systems Design and Application August (ISDA 2004)*, ISBN: 9637154302.
- Rennolls, K. & Al-Shawabkeh, A., 2008. Formal structures for data mining, knowledge discovery and communication in a knowledge management environment. *Intelligent Data Analysis*, 12(2), p. 147-163.
- Ronco, S. L., 1996. How Enrollment Ends: Analyzing the Correlates of Student Graduation, Transfer and Dropout with a Competing Risks Model, Tallahassee, Fla.: AIR Professional File, No. 61 Association for Institutional Research.
- Ruan, T. . L. D., 2007. An extended process model of knowledge discovery in database. *Journal of Enterprise Information Management*, 20(2), pp. 169 - 177.
- SAS, I., 2008. SEMMA data mining methodology. [Online] Available at: <http://www.sas.com>
- Schilit, B., Adams, N. & Want, R., 1994. Context-Aware Computing Applications. s.l., *First International Workshop on Mobile Computing Systems and Applications*.
- Sharma, S., Osei-Bryson, K.-M. & Kasper, G. M., 2012. Evaluation of an integrated Knowledge Discovery and Data Mining process model. *Expert Systems with Applications*, Volume 39, p. 11335-11348.
- Tinto, V., 1975. Dropouts from higher education: A theoretical synthesis of recent literature. *A Review of Educational Research*, Volume 45, pp. 89-125.

36. Vert, G., Chennamaneni, A. & Iyengar, S. S., 2010. Potential Application of Contextual Information Processing To Data Mining.. Las Vegas Nevada, USA., Proceedings of the 2010 International Conference on Information & Knowledge Engineering, IKE 2010, July 12-15, 2010.
37. Vialardi, . C., Chue, J., Peche, J. P. & Alvarado, G., 2011. A data mining approach to guide students through the enrollment process based on academic performance. User Modeling and User - Adaptation Interaction, Volume 21, pp. 217-248.
38. Zhang, C., Yu, P. . S. & Bell, D., 2010. Introduction to the Domain-Driven Data Mining Special Section. IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, 22(6).
39. BIS (2014) Improving the Student Learning Experience – a national assessment (BIS Research Paper no 169). London: Department for Business, Innovation and Skills.
40. G. Eason, B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," Phil. Trans. Roy. Soc. London, vol. A247, pp. 529–551, April 1955.
41. Abeer Badr El Din Ahmed and Ibrahim Sayed Elaraby, "Data Mining: A prediction  
42. for Student's Performance Using Classification Method", World Journal of Computer Application and Technology 2(2): 43-47, 2014
43. Lu, K.J. and Sailesh, S.B., 2016. Context driven data mining to classify students of higher educational institutions.
44. R. J. Vidmar. (1992, August). On the use of atmospheric plasmas as electromagnetic reflectors. IEEE Trans. Plasma Sci. [Online]. 21(3). pp. 876—880. Available: <http://www.halcyon.com/pub/journals/21ps03-vidmar>

### AUTHORS PROFILE

**Dr. Subhashini Bhaskaran Sailesh**, Assistant Professor, Ahlia University, Bahrain Research interests include Data Mining, Data Science Business Analytics, Data Analytics

**Prof. Mansoor Al Aali**, President, Ahlia university, Bahrain Research interests include Data Mining, Data Science, Artificial Intelligence

**Dr. Kevin Lu Brunel University London**