



# Student performance analysis and prediction in classroom learning: A review of educational data mining studies

Anupam Khan<sup>1</sup>  · Soumya K. Ghosh<sup>1</sup>

Received: 4 February 2020 / Accepted: 17 May 2020 / Published online: 01 July 2020  
© Springer Science+Business Media, LLC, part of Springer Nature 2020

## Abstract

Student performance modelling is one of the challenging and popular research topics in educational data mining (EDM). Multiple factors influence the performance in non-linear ways; thus making this field more attractive to the researchers. The widespread availability of educational datasets further catalyse this interestingness, especially in online learning. Although several EDM surveys are available in the literature, we could find only a few specific surveys on student performance analysis and prediction. These specific surveys are limited in nature and primarily focus on studies that try to identify possible predictor or model student performance. However, the previous works do not address the temporal aspect of prediction. Moreover, we could not find any such specific survey which focuses only on classroom-based education. In this paper, we present a systematic review of EDM studies on student performance in classroom learning. It focuses on identifying the predictors, methods used for such identification, time and aim of prediction. It is significantly the first systematic survey of EDM studies that consider only classroom learning and focuses on the temporal aspect as well. This paper presents a review of 140 studies in this area. The meta-analysis indicates that the researchers achieve significant prediction efficiency during the tenure of the course. However, performance prediction before course commencement needs special attention.

**Keywords** Student performance · Classroom learning · Performance prediction · Literature review · Educational data mining

---

✉ Anupam Khan  
[anupamkh@iitkgp.ac.in](mailto:anupamkh@iitkgp.ac.in)

Soumya K. Ghosh  
[skg@cse.iitkgp.ac.in](mailto:skg@cse.iitkgp.ac.in)

<sup>1</sup> Department of Computer Science and Engineering, Indian Institute of Technology Kharagpur, Kharagpur, West Bengal 721302, India

# 1 Introduction

Student performance analysis and prediction are the two widely explored research topics in education system literature. Although their objectives are different, the outcomes of performance analysis significantly influence the prediction studies. The statistical methods may not always be sufficient for establishing the association of various factors with performance (Natek and Zwilling 2014). The use of sophisticated algorithms may yield impressive knowledge that could help the educators as well as students (Baker and Yacef 2009; Kumar et al. 2018; Liu et al. 2018b; Romero and Ventura 2007). The advancement of data mining technologies has influenced many researchers to investigate more profound insights into the knowledge dissemination process. Some of them have applied data mining approaches for this purpose long ago (Fausett and Elwasif 1994; Gedeon and Turner 1993). However, the number of such studies are very less at that time. The growing availability of digital data captured by several academic information management systems and educational software in recent years catalyses this process to improve the quality of education (Aghabozorgi et al. 2014; Asif et al. 2017; Baker 2014; Khanna et al. 2016; Livieris et al. 2018; Loh and Sheng 2015; Ogor 2007; Wook et al. 2016). Interestingly, Baker and Inventado (2014) estimated that all educational research would involve analytics and data mining by 2022.

The researchers have earlier conducted multiple surveys on EDM studies. Some of them consider a broader scope to outline multiple aspects of educational processes. A few of them specifically focuses on student performance analysis and prediction as well. Let us briefly introduce these surveys here. The first one was published more than a decade ago in which Romero and Ventura (2007) consolidate the studies published between 1995 and 2005 and outline the probable objectives of EDM. In 2010, they published another exhaustive review (Romero and Ventura 2010) which considers nearly 300 papers published between 1993 and 2009. It divides the research area of EDM into eleven broad sub-areas, and student performance prediction is one of them. The area of learning analytics has also gained momentum by that time. It is essential to mention that two research communities, the international educational data mining society and the society for learning analytics research, initially drive the analytical studies in education. Siemens and Baker (2012) have mentioned the overlap of research interest between these two communities. In the same year, Ferguson (2012) outlines the drivers and challenges behind learning analytics studies. In 2013, Chrysafiadi and Virvou (2013) carried out a specialised review on student modelling approaches. It may be the first of this kind that focuses on a specific area of EDM. Peña-Ayala (2014) has published another exhaustive survey in 2014, which considers 240 EDM studies between 2010 and the first quarter of 2013. According to this study, the number of works on student performance modelling is increasing significantly since 2010. The researchers have additionally carried out specific EDM survey on the psychology of learning and application of clustering approaches as well (Koedinger et al. 2015; Dutt et al. 2017). Bakhshinategh et al. (2017) have probably conducted the most recent review of EDM studies in which they mention that student performance prediction can help the educators and administrators significantly.

In addition to the broad survey on EDM studies, a few specific surveys on student performance analysis and prediction are also available in the literature. Shahiri et al. (2015) have published first such survey in 2015 to the best of our knowledge. It primarily focuses on the predictor attributes and the prediction methods used in 39 studies published between 2002 and 2015. Later on, Khanna et al. (2016) outline the parameters used for predicting students performance in 25 research papers. Furthermore, another review of 40 papers in 2017 provides an overview of prediction accuracy in addition to the factors and methods used (Hu et al. 2017). Kumar et al. (2017) conducted a similar survey with 16 papers in the same year with the primary focus on used predictors, methods and accuracy. In another two short reviews (Shingari and Kumar 2018; Kumar et al. 2018), the authors outline the data mining techniques used for predicting performance in 17 and 36 papers, respectively.

Although the researchers have applied EDM techniques in both traditional and computer-based online education, the application in conventional education is comparatively less than the other alternatives available. The authors of two separate EDM review papers (Peña-Ayala 2014; Romero and Ventura 2010) therefore found only a few studies on the traditional education system. Moreover, all specific EDM surveys on student performance consider the studies on both online and traditional classroom-based learning. Importantly, we could not find any EDM survey paper which focuses explicitly on classroom-based education. The previous surveys on student performance primarily focus on the predictors, methods and prediction efficiency. However, none of them considers the temporal aspects. The prediction before and after course commencement may serve a completely different purpose. There is a pressing demand for next-term performance prediction nowadays to provide timely and effective support (Backenköhler and Wolf 2017; Polyzou and Karypis 2019; Shingari and Kumar 2018; Sweeney et al. 2016). Prior knowledge of performance influencing factors can help in implementing such advanced prediction (Helal et al. 2018). Therefore, this survey primarily focuses on the following research questions:

- What are the factors influencing student performance in classroom learning?
- Which methods are used for finding these factors?
- Is it possible to predict before course commencement?
- Can we predict the actual grade or score?

In this paper, we present a systematic survey of student performance related EDM studies with a specific focus on classroom learning. Various technologies like online discussion forum, assessment system, assignment and learning tools are used in classroom learning nowadays. This survey does not restrict to the studies only on those educational systems that follow the traditional instruction delivery mechanism. Instead, it considers all works which are related to the students of the traditional education system. In this paper, we present a review of 140 relevant studies published between 2000 and 2018. In order to carry out this survey, we have created a taxonomy of research directions and categorise the relevant works accordingly. The meta-analysis finally helps in identifying future research directions.

The anatomy of the paper is as follows: Section 2 elaborates the methodology adopted for conducting this survey. Sections 3 and 4 present the meta-analysis of relevant works. In Section 3, we first outline the existing studies of student performance

analysis and Section 4 later presents the studies on performance prediction before and after course commencement separately. Section 5 discusses the observations and recommendations for future work. Finally, Section 6 concludes the survey.

## 2 Methodology

This study has followed the recommendation provided by Cooper (1988) for determining the focus, goal, perspective, coverage, organization and audience of this survey (refer Table 1). The goal here is to identify the gaps or central issues with a particular focus on student performance analysis and prediction related literature. This neutral survey is not an exhaustive one; rather, it represents the relevant literature pivotal to the research questions framed. In this paper, Fig. 1 presents the methodology adopted for conducting this survey. The following parts of this section mention the detailed steps.

In a literature survey, a well-planned and well-executed search strategy is essential for finding every relevant piece of work (Shahiri et al. 2015). Therefore, we have first created a taxonomy of research directions as a part of the survey strategy. Figure 2 presents the defined taxonomy here. The shaded part of Fig. 2 is the primary focus of this survey. We have organised the paper in these research directions. Once the taxonomy is defined, we have adopted a hybrid approach for searching the relevant literature. The Google Scholar helps us in configuring an alert with (“educational data mining” and “student performance”) string. Due to this configuration, Google Scholar regularly sends a list of recently published relevant papers to our email since August 2015. We have additionally adopted the conventional approach for searching relevant literature as well. In this approach, we search the Google Scholar with the following search strings: (i) “educational data mining” and “student performance”, (ii) “learning analytics” and “student performance”, (iii) “student performance” and “teaching quality”, (iv) “student performance” and “domain knowledge”. This approach facilitates us to get the studies published before August 2015 as well. Importantly, Google Scholar is selected here for better coverage of the search space (Walters 2007).

The above hybrid search approach retrieves multiple articles which are filtered systematically to obtain the list of selected articles. The recommendation of the

**Table 1** Scope of the survey

Characteristics	Scope
Focus	Research Outcome
Goal	Identification of central issues
Perspective	Neutral representation
Coverage	Representative, pivotal
Organisation	Conceptual
Audience	Specialised scholars

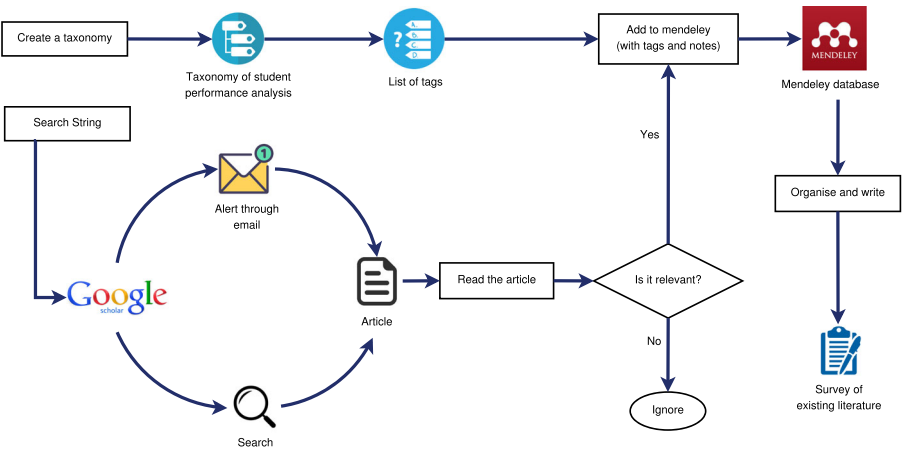


Fig. 1 Methodology adapted for this survey

Prisma Statement (Moher et al. 2009) helps in the selection process here. As an initial step, we have checked the title to identify studies related to student performance analysis or prediction. This process helps us in identifying a total of 362 distinct studies as a potential candidate for this survey. We read the abstract and keyword

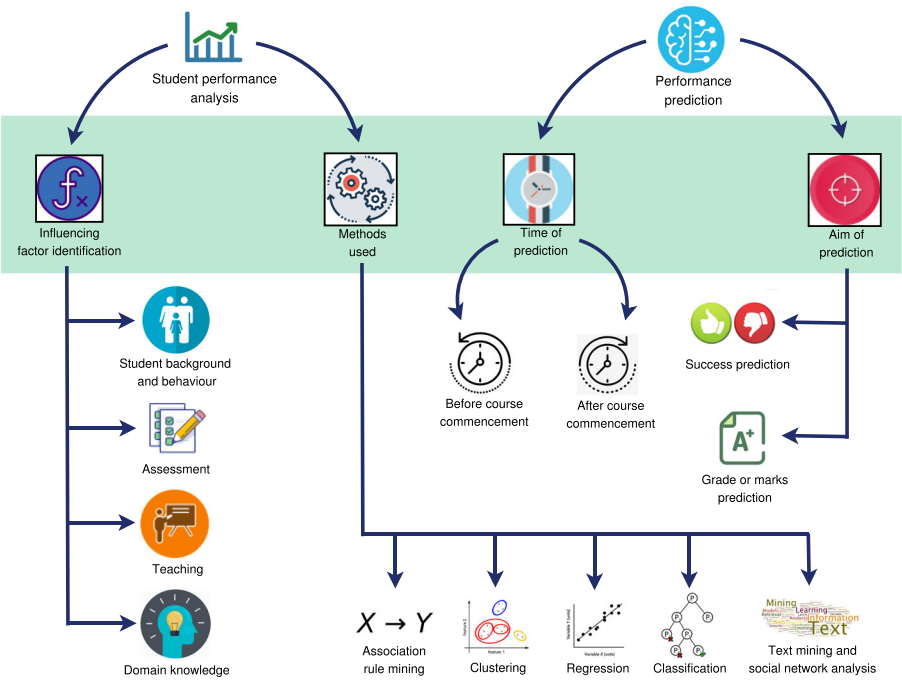


Fig. 2 Taxonomy of research directions

mentioned in each one of them for further screening, and it filters out 65 irrelevant studies. In the next step, we have read the full-text and tried to identify whether the contribution of the study is pivotal to the research questions of this survey. This process helps us in excluding the studies which are not related to the students of traditional classroom-based education as well. We also discard those studies which are not well-presented, contains unclear methodology, dataset or contribution. Due to these criteria, 157 studies do not pass the eligibility test, and 140 relevant articles finally get selected for this review. In the next step of the qualitative synthesis, we add the shortlisted articles in Mendeley (Zaugg et al. 2011), a reference management tool provided by Elsevier. During inclusion, we also classify each article with multiple tags. Each item in the defined taxonomy except those on the shaded region of Fig. 2 are the tags used here. It helps us in grouping the articles further. Moreover, we use the general note section of Mendeley to keep a track on the key findings, significant contributions and limitations of each article which eventually facilitates us to complete the meta-analysis.

### 3 Student performance analysis

The choice of predictors plays a vital role in any efficient prediction task. The factors affecting knowledge dissemination process, therefore, should be identified first before predicting student performance (Helal et al. 2018). Furthermore, the decrease in students' success rate and frequent drop-outs are the primary concern of many educational institutes nowadays. Prior analysis of student performance could help the institutes in this regard (Mat et al. 2013). This section presents a meta-analysis of existing studies which tries to identify and measure the influencing factors of student performance. It first introduces the factors found in the literature and later on discusses the methods used for establishing the association of these factors with student performance.

#### 3.1 Influencing factor identification

Various factors impact the performance in non-linear ways; thus, their identification for subsequent use as predictors of success is a complex problem (Zollanvari et al. 2017). The social and academic background of the student, their behaviour, previous assessments, quality of teaching, previous knowledge on the course topics are a few broad areas explored by the researchers. Here we present a meta-analysis of existing literature. Besides this, Tables 2, 3, 4 and 5 summarise them in tabular format for providing a comprehensive view to the readers.

##### 3.1.1 Student background and behaviour

The exploration of student demographics, behaviour, academic and social background related factors are widely available in student performance analysis literature. The influence of gender is one of the widely debated topics among them. Ahmed et al. (2014) show that the student is more likely to be a male student in case of poor

**Table 2** Student background and behaviour related studies

Study	Explored predictor	Method
Ahmed et al. (2014)	Gender, residence type	ARM
Al-Obeidat et al. (2017)	Parent education and job	CLA
Bayer et al. (2012)	Communication skill, social behaviour	SNA, CLA
Bendikson et al. (2011)	Socio-economic status	REG
Campagni et al. (2015)	Order of examination taken	CTR
Carter et al. (2017)	Programming and social behaviours	REG
Daud et al. (2017)	Family expenditure and personal information	CLA
Dvorak and Jia (2016)	Study habits	REG
Felisoni and Godoi (2018)	Cell phone usage behaviour	REG
Fernandes et al. (2018)	Age, Social and academic background	CLA
García et al. (2007)	Learning management system usage behaviour	ARM
Gasevic et al. (2017)	Learning strategies	CTR
Gowda et al. (2013)	Learning style	REG
Gray et al. (2014)	Psychometric factors	CLA
Guruler et al. (2010)	Registration category and family income	CLA
Hart et al. (2017)	Attitude, cognitive skills, and engagement	REG
Hattie and Clinton (2012)	Physical activity	REG
Hassan and Rasiah (2011)	Poverty level	REG
Helal et al. (2018)	Social background	CTR
Ivančević et al. (2010)	Spatial position in laboratory	ARM, CTR
Mishra et al. (2014)	Social, academic background, emotional skills	CLA
Natek and Zwilling (2014)	Past academic and social background	CLA
Pal and Chaurasia (2017)	Alcohol consumption	CLA
Papamitsiou et al. (2014)	Temporal behaviour of answering questions	REG
Quadri and Kalyankar (2010)	Scholarship status, first child	CLA
Ramesh et al. (2013)	Parent occupation	REG
Romero et al. (2008)	Activity in moodle	CLA
Romero et al. (2013)	Participation in online discussion forum	SNA, CTR
Saarela and Kärkkäinen (2015)	Learning capabilities	REG, CTR
Saxena and Govil (2009)	Parent occupation, student demographics	CTR
Van Inwegen et al. (2015)	Assignment tool usage pattern	REG
Wang et al. (2016)	Assignment tool usage pattern	REG
Zacharis (2015)	Activity in learning management system	REG
Zimmermann et al. (2015)	Performance in undergraduate level	REG

ARM: Association rule mining, CLA: Classification, CTR: Clustering, REG: Regression, SNA: Social network analysis

performance. In contrast, Quadri and Kalyankar (2010) did not find any influence of gender on student performance. They instead find a significant impact of low family income on performance and drop-out. A few other studies have also reported a similar observation (Daud et al. 2017; Guruler et al. 2010; Hassan and Rasiah 2011). Inter-

**Table 3** Assessment related studies

Study	Explored predictor	Method
Ahmed et al. (2014)	Internal and external assessment, attendance	ARM
Al-Barrak and Al-Razgan (2016)	Performance in other courses	CLA
Beemer et al. (2018)	Internal assessment	REG
Bucos and Druagulescu (2018)	Attendance, average score and credits, student activity	CLA
Buldu and Üçgün (2010)	Previous failure	ARM
Buniyamin et al. (2015)	CGPA, past academic performance	CLA
Chaturvedi and Ezeife (2013)	Assignment marks	ARM
Christian and Ayub (2014)	CGPA	CLA
Fernandes et al. (2017)	Attendance	CLA
Golding and Donaldson (2006)	Assessment in first year courses	REG
Guarin et al. (2015)	Admission test score	CLA
Helal et al. (2018)	Internal assessment	CTR
Huang and Fang (2013)	CGPA, internal and external assessment	REG
Jishan et al. (2015)	CGPA, attendance and internal assessment	CLA
Kamley et al. (2016)	Performance in school	ARM
Kaviyarasi and Balasubramanian (2018)	Internal test scores	CLA
Li et al. (2013)	Score in first year mathematics	REG
O'Connell et al. (2018)	Performance in prior courses	REG
Parack et al. (2012)	Internal assessment	ARM
Romero et al. (2008)	Internal assessment data in moodle	CLA
Superby et al. (2006)	Attendance, external assessment	CLA
Tair and El-Halees (2012)	Performance in secondary education	ARM, CLA
Xiong et al. (2014)	Prerequisite skills	REG

ARM: Association rule mining, CLA: Classification, CTR: Clustering, REG: Regression

estingly, the fact of whether the student is the first child of the parent also found to be influential on performance (Quadri and Kalyankar 2010). Some studies have even reported the possible impact of parent occupation (Al-Obeidat et al. 2017; Ramesh et al. 2013), registration category (Guruler et al. 2010; Helal et al. 2018), past academic and socio-economic background of the student (Bendikson et al. 2011; Helal et al. 2018; Natek and Zwilling 2014), neighbourhood, previous school and age (Fernandes et al. 2017, 2018) as well. In contrast, Hattie and Clinton (2012) show that physical activity does not influence student performance at all. However, Gray et al. (2014) exhibit that psychometric factors like learner ability, personality, motivation, and learning strategies are useful to identify the college student who is at risk of failure. In a similar study, Mishra et al. (2014) mention that the social and academic integration and various emotional skills of the student are also capable of building a performance prediction model.



**Table 4** Teaching related studies

Study	Explored predictor	Method
Abrami et al. (2007)	SETe	REG
Balam and Shannon (2010)	SETe, gender	REG
Brocato et al. (2015)	SETe, gender of teacher and course level	REG
Centra (2003)	Difficulty level, workload	REG
Figlio and Lucas (2004)	Grading pattern of the teacher	REG
Galbraith et al. (2012)	SETe, class size	CLA
Goos and Salomons (2016)	SETe	REG
Khan and Ghosh (2016)	SETe	ARM
Khan and Ghosh (2018)	SETe, number of student evaluation	ARM
Lin et al. (2018)	Size of the class taught by the teacher	REG
Macfadyen et al. (2015)	SETe, gender and age of teacher	REG
Nikolic et al. (2015)	Quality of notes, facilities and equipments in laboratory	REG
Pandey and Pal (2011)	Language used for teaching	ARM
Pong-Inwong and Rungworawut (2012)	Open-end questions of SETe	TM, CLA
Price et al. (2017)	SETe, gender of the teacher and student	REG
Rani and Kumar (2017)	Students' comment on teaching	TM
Stronge et al. (2007)	Instruction type, SETe, personal qualities	REG
Üstünlüoğlu (2016)	Nationality and gender of teacher	REG
Uttl et al. (2017)	SETe	REG
Yin et al. (2016)	SETe in theoretical courses	REG

ARM: Association rule mining, CLA: Classification, REG: Regression, TM: Text mining

The researchers have analysed the impact of students' behaviour on their performance as well. In a study, Papamitsiou et al. (2014) have found the influence of students' temporal behaviour while answering the questions during the examination. Even the spatial position in the laboratory influences student performance (Ivančević et al. 2010). Campagni et al. (2015) have analysed the performance relating to the order of examination taken by the student. The result shows that superior students took most exams according to the order planned by the curriculum. Dvorak and Jia (2016) find that students can earn higher grades if they start working on assignments earlier. In another interesting study, Pal and Chaurasia (2017) observe the negative effect of alcohol consumption on student performance. Some studies have even reported the possible impact of cell phone usage (Felisoni and Godoi 2018), social behaviour (Bayer et al. 2012; Carter et al. 2017), and active participation in the online discussion forum (Romero et al. 2013). Several universities nowadays use learning management tools for enhancing the knowledge dissemination process. The data extracted from these tools have helped some studies (Romero et al. 2008; Van Inwegen et al. 2015; Wang et al. 2016) to analyse student performance as well. In a similar work, Zacharis (2015) finds the frequency of reading and posting messages,

**Table 5** Domain knowledge related studies

Study	Explored predictor	Method
Adjei et al. (2016)	Prerequisite skills	REG
Bydžovská (2016)	Performance in similar courses	CTR
Chen et al. (2017)	Performance in prerequisite	CTR
Damaševičius (2010)	Difficulty level of the course topic	ARM
She et al. (2012)	Previous performance in related course	REG

ARM: Association rule mining, CTR: Clustering, REG: Regression

viewing files, efforts on content creation, solving quiz may be considered as the predictor of students' grade. Besides, the learning strategy is also capable of indicating student performance (Gasevic et al. 2017).

### 3.1.2 Assessments

EDM researchers have widely explored the cumulative grade point average (CGPA), internal and external assessments for predicting student performance. The internal assessment in literature refers to the internally evaluated scores in assignments, quiz, lab work, class test or other similar components. In contrast, the marks obtained in other courses are well-known as external assessments among researchers (Shahiri et al. 2015). Chaturvedi and Ezeife (2013) have analysed the association of assignment marks with final grade and found a positive relationship between these two. In fact, multiple studies have explored various types of internal assessment, like performance in the quiz, laboratory, mid-term examination etc. Most of them found a positive influence on final performance (Ahmed et al. 2014; Helal et al. 2018; Jishan et al. 2015; Kaviyarasi and Balasubramanian 2018; Romero et al. 2008). Besides this, evidence of the positive impact of attendance (Ahmed et al. 2014; Bucos and Druagulescu 2018; Fernandes et al. 2017; Jishan et al. 2015) and teamwork score (Parack et al. 2012) also exist in the literature.

In addition to internal assessments, the researchers have also explored the external assessments for analysing student performance. In one such study, Buldu and Üçgün (2010) find that the students who could not succeed in a numerical course can fail again in further attempts. They even observe a high correlation of failure in mathematics, physics and chemistry. The performance in first-year mathematics or computer science can act as a performance indicator of other courses later (Golding and Donaldson 2006; Li et al. 2013). Al-Barrak and Al-Razgan (2016) propose a mechanism for identifying the most critical mandatory courses that can help in predicting final grade. The achievements in previous courses are indeed a generic indicator of student performance (O'Connell et al. 2018). Importantly, the CGPA considers all prior performances and consolidates them to a single value. It is therefore quite evident that many studies could establish CGPA as a significant predictor of student performance (Buniyamin et al. 2015; Christian and Ayub 2014; Huang and Fang 2013; Jishan et al. 2015). Some studies have even found the association with

school-level performance (Kamley et al. 2016; Tair and El-Halees 2012), admission test score (Guarin et al. 2015), previous academic experience (Superby et al. 2006) as well.

### 3.1.3 Teaching

Analysing the effect of teaching is a long-debated topic in literature. A majority of existing studies analyse the impact of teaching on indirect factors such as students' motivation, satisfaction etc. Only a few studies exist which explores the direct influence on student performance. Here we present a meta-analysis and a summary (refer Table 4) of these studies which analyse various aspects of teaching and their direct or indirect impact on student performance.

In one of the related study, Stronge et al. (2007) mention that teaching excellence with efficient knowledge dissemination, student assessment, classroom management, and superior personal qualities helps the student in achieving their desired goal. Another study (Pong-Inwong and Rungworawut 2012) shows that prior knowledge of the class sentiment is helpful for the teacher to improve their instruction process which could, in turn, impact student performance. The value addition is a measure which can be useful in determining the learning achievement (Stronge et al. 2007). Some studies (Yin et al. 2016; Khan and Ghosh 2016, 2018) indeed observe that effective teaching motivates the student to perform better whereas performance degrades with poor teaching. A regression analysis shows that student satisfaction is related to the quality of notes, facilities and equipment in the laboratory as well (Nikolic et al. 2015). In another study, Pandey and Pal (2011) observe a significant impact of the instruction delivery language on class attendance which influences the performance as well. Higher grading standard by teachers can also be beneficial to the students. The magnitude of such benefit depends on student quality (Figlio and Lucas 2004). The class size also influences the teaching quality which in turn affects the final grade of students (Lin et al. 2018).

Abrami et al. (2007) have mentioned that the relationship between teaching quality and learner achievement is an exciting research area, and it needs further detail investigation. However, determining the proper measure of teaching quality is still a research challenge (Khan and Ghosh 2018). Importantly, student evaluation of teaching excellence (SETE) is a powerful instrument which can help to overcome this challenge. Üstünlüoğlu (2016) indeed put more emphasis in favour of SETE related research. However, the author reminds the necessity of student, employer and teacher-centric SETE. In addition to formal measures, sentiment analysis of student feedback is an indirect assessment which facilitates the teachers to assess the students' interest in class (Rani and Kumar 2017). However, the community is sceptical about the validity of such evaluation. For example, Uttl et al. (2017) do not find a strong correlation between learning achievement and teacher rating. Galbraith et al. (2012) observe the ability of SETE in predicting student learning achievement though raises a concern about possible bias of class size. Goos and Salomons (2016) also reminded about the inaccuracy of SETE due to varying response rate across courses. Not only that, the bias between student gender, age, specialisation area, final grade

and SETE is also evident in the literature (Macfadyen et al. 2015; Price et al. 2017; Zabaleta 2007).

In contrast, some researchers consider SETE as a reliable, stable and relatively valid indicator of effective teaching (Balam and Shannon 2010; Moore and Kuol 2005; Khan and Ghosh 2016, 2018). Marsh (2007) thinks that SETE may be considered as a useful instrument for teachers to gather feedback on teaching, student to select courses, and administration to make decisions. However, a systematic feedback collection process is necessary, and proper analysis of it can help in many ways (Hattie and Timperley 2007; Jara and Mellar 2010). Although possible gender bias is evident in the literature, Brocato et al. (2015) report that it is not very significant in a traditional classroom setting. Centra (2003) even rejects the myth of possible bias between final grade and SETE. Importantly, Marsh (1984) has explained the reason behind this myth long ago. According to the author, the co-occurrence of higher ratings and higher grades might be a result of better learning due to more effective teaching.

### 3.1.4 Domain knowledge

According to Blooms' taxonomy (Bloom et al. 1956), knowledge involves the recall or recognition of terms, ideas, processes, and theories. In the learning context, domain knowledge is the cognition on a similar topic learned earlier. However, quantifying the domain knowledge is a challenge on its own (Lorenzetti et al. 2016). The researchers have used the domain knowledge to analyse various areas such as self-regulated learning (Moos and Azevedo 2008), web-based problem-solving pattern (She et al. 2012), reading comprehension (McCarthy and Goldman 2017), internet searching for learning objects (Willoughby et al. 2009) etc. These studies primarily measure the domain knowledge based on previous experience, specialisation, or through a set of questionnaires. The syllabus is available in digital format nowadays, which contains valuable information about the course content. The course content indeed represents the domain knowledge repositories and cognitive models of the knowledge components to be learned (Chung and Kim 2016; Peña-Ayala 2014). However, only a limited number of studies measuring the direct impact of domain knowledge on student performance exists in the literature. Here we present a few closely related studies and summarises them in Table 5.

In one of the few studies available, Adjei et al. (2016) show that availability of prerequisite skills affect student performance. Chen et al. (2017) have used a prerequisite dataset to measure domain knowledge which facilitates the prediction of student response in assessments. The difficulty level of course topic also has a significant effect on performance. Prior knowledge could help in such cases (Damaševičius 2010). Bydžovská (2016) have clustered the courses based on cosine similarity measures and later on analyses the student performance. They found that these cluster of courses can help in performance prediction. In another study, She et al. (2012) use past performance in related courses to quantify the domain knowledge. Even they choose similar courses manually but able to establish an association of domain knowledge with student performance.

### 3.2 Methods used for identifying the factors

Researchers have applied various techniques such as clustering, classification, association rule mining etc. for extracting valuable knowledge from the educational dataset. This subsection briefly discusses the EDM approaches applied in the existing literature.

#### 3.2.1 Association rule mining

Agrawal et al. (1993) introduced the concept of association rule mining for market basket analysis. Later on, it was widely used in various other domains as well. The association rule mining is one of the well-known and popular data mining techniques used extensively for educational purposes (Kamley et al. 2016). It is beneficial indeed for understanding the pedagogical aspects of learning which in turn help the academic administrator to frame policies (Angeli et al. 2017; Damaševičius 2010; García et al. 2007). There is a reasonable number of studies exist which use the association rule mining for analysing student performance (Ahmed et al. 2014; Buldu and Üçgün 2010; Chaturvedi and Ezeife 2013; Damaševičius 2010; García et al. 2007; Kamley et al. 2016; Pandey and Pal 2011; Parack et al. 2012; Tair and El-Halees 2012).

In order to obtain a set of important rules, it is essential to pre-determine the minimal support and confidence. However, it is difficult for an educator to decide these two input parameters in advance. Furthermore, the number of obtained rules may be too high in some cases, and most of them are non-interesting and with low comprehensibility. Instead of filtering, the ranking of important rules may be one possible solution to address this problem (García et al. 2007). A combined measure of cumulative interestingness may also be effective in this context (Damaševičius 2010). Above all, a non-parametric technique of data mining can be more useful for the end-users in educational settings (Zorilla et al. 2010).

#### 3.2.2 Regression

The regression is a prevalent technique for statistical analysis, and it is quite popular in the data mining context as well. It is probably the most straightforward method for investigating the functional relationship between variables (Kotsiantis and Pintelas 2005). An equation between the dependent variable and one or multiple independent variables defines the relationship in regression analysis. The equation is the most critical product here. In prediction settings, it is also used for evaluating the importance of individual predictors and understand the correlation among variable (Chatterjee and Hadi 2015; Lipovetsky and Conklin 2015).

The widespread application of regression is evident in educational data mining literature. It is indeed frequently used in EDM studies for establishing or disproving the impact of teaching quality. Many research works have explored the regression analysis for student performance analysis (Abrami et al. 2007; Adjei et al. 2016; Centra 2003; Goos and Salomons 2016; Li et al. 2013; Macfadyen et al. 2015; Papamitsiou et al. 2014; Stronge et al. 2007; Üstünlüoğlu 2016; Yin et al. 2016; Zabaleta 2007).

### 3.2.3 Classification

It is one of the primary data mining technique widely used for predicting group membership of data instances. Researchers have widely applied this method of classifying each item in a dataset into one of the predefined class or group (Kesavaraj and Sukumaran 2013). It is possibly the most widely used data mining techniques in the educational context as well. To support this, Peña-Ayala (2014) mentions that 42.15 % of EDM research has applied the classification technique between 2010 and the first quarter of 2013. There are several classification approaches available in literature such as decision tree methods, rule-based classification, memory-based learning, neural networks, bayesian network, and support vector machines. The decision tree provides a powerful formalism for representing comprehensible and accurate classifier (Quinlan 1990). The decision tree classification is the most popular method used in educational research (Shahiri et al. 2015).

The predictive ability of the classification technique makes it a superior choice for student performance analysis. The primary objective of classification in EDM is to identify the critical factors that contribute to the final grade of the student (Al-Barrak and Al-Razgan 2016). Moreover, student classification can help teachers to understand the specific requirements and determine an appropriate approach as well (Hidayah et al. 2013). There are various studies available in EDM literature that uses classification approaches for performance predictor identification (Al-Barrak and Al-Razgan 2016; Buniyami et al. 2015; Christian and Ayub 2014; Gray et al. 2014; Jishan et al. 2015; Guarin et al. 2015; Mishra et al. 2014; Natek and Zwilling 2014; Pong-Inwong and Rungworawut 2012; Quadri and Kalyankar 2010; Ramesh et al. 2013; Romero et al. 2008, 2013).

### 3.2.4 Clustering

It is a concept of dividing data into groups of similar objects where each group, or cluster, consists of objects that are similar to one another and dissimilar to the object of other groups (Berkhin 2006). It is a data modelling technique which has a historical milestone in the field of mathematics and statistics. However, in machine learning, clusters correspond to hidden patterns, the search for a cluster is unsupervised learning, and the resulting system represents the knowledge (Kaufman and Rousseeuw 2009). Partitioning and hierarchical methods are the two broad categories of clustering, whereas k-means and agglomerative hierarchical techniques are very popular among them.

Although clustering is popular in data mining literature, the researchers have reported some limitations of clustering approaches. These are quite relevant for educational datasets as well. Traditional clustering approaches are not suitable when both quantitative and qualitative attributes coexist (Chen et al. 2016). However, most of the educational dataset contains both of these types. Most of the clustering approach additionally require an input parameter which determines the number of resulting cluster. It is difficult for an educator to determine this parameter (Bogarin et al. 2014). Nevertheless, EDM researchers have considerably used this unsupervised technique

in student performance analysis literature (Bydžovská 2016; Campagni et al. 2015; Helal et al. 2018; Ivančević et al. 2010; Romero et al. 2013; Saxena and Govil 2009).

### 3.2.5 Text mining and social network analysis

Text mining is the full or partially automated process of extracting hidden information from a large amount of unstructured data. Although text mining is a part of the general data mining, it differs from other approaches. It extracts the patterns from natural language text rather than from structured data (Delen and Crossland 2008; Feldman and Sanger 2007; Romero and Ventura 2007). Moreover, the social network is a web-based service commonly used for information dissemination, personal activities posting, product reviews, online pictures sharing, professional profiling, advertisements, sentiment expression etc. The social network is getting tremendously popular in the last decade (Zuber 2014). Nowadays, people are relying more on the social network for interaction with other users. It, in turn, generates massive data characterised by three computational issues, namely, size, noise and dynamism. Data mining provides various techniques for extracting useful knowledge from these massive datasets (Kagdi et al. 2007).

Researchers have identified some topics on text mining and social network analysis, which are growing popularity nowadays. Community detection, semantics, opinion and sentiment analysis are a few examples among them (Zuber 2014). Interestingly these are very relevant in education as well. Detecting student communities and analysing community performance can be a thought-provoking study. Similarly, analysis of student or teacher sentiment may open a new direction of understanding the knowledge dissemination process. The answer script, textual feedback, communication between student and teacher can turn out to be valuable inputs for such analysis (Peña-Ayala 2014). A few interesting EDM studies are available on text mining and social network analysis (Akçapinar 2015; Bayer et al. 2012; Chung and Kim 2016; Foley and Allan 2016; Montuschi et al. 2015; Pong-Inwong and Rungworawut 2012; Rani and Kumar 2017; Rekha et al. 2012; Romero et al. 2013). However, these are the initial steps only, and there is ample scope of exploring the educational data with the help of these technologies further (Yim and Warschauer 2017). The application of semantic approaches has a huge potential in education as it manages the knowledge objects through unstructured data like syllabus, textbook, question paper etc. (Peña-Ayala 2014). Although some semantic studies are available in EDM literature (Elouazizi et al. 2017; Hsiao et al. 2016; Hsiao and Lin 2017; Liu et al. 2018a; Montuschi et al. 2015; Nakayama 2016), it is incipient in the field of student performance analysis.

## 4 Student performance prediction

The final grade is supposed to summarise how well a student understands and apply the knowledge conveyed in a course. However, predicting the actual grade is a challenging task. This fact has energised more EDM researchers to discover an efficient model of student performance (Bresfelean et al. 2008; Guo et al. 2015; Meier



et al. 2016; Márquez-Vera et al. 2013). In addition to efficiency, time of prediction is another crucial aspect which facilitates the decision-makers and teachers in taking a right and timely measures (Quille and Bergin 2018). The performance modelling studies use various attributes as the predictor. However, only a few of these are available before course commencement which creates more challenges in early prediction.

## 4.1 Time of prediction

In this part, we present the existing performance modelling literature classified as per their temporal nature of prediction. It separately presents the meta-analysis of studies that tries to predict performance after or before course commencement. Besides this, Tables 6, 7, 8 and 9 summarises them in tabular format.

### 4.1.1 During the tenure of the course

This survey observes a significant number of studies which predict the performance during the tenure of the course efficiently. Tables 6 and 7 provides a list of these studies with the details on predictor and method used, aim and observed efficiency. It is quite apparent that efficient predictors are required to yield a better result. Fortunately, many predictor data are readily available during the tenure of the course. According to Huang and Fang (2013), the CGPA is the sole attribute capable of predicting the overall academic performance of a course and the grades in prerequisites, mid-term examination score and the CGPA are the most suitable predictors for the individual outcome. A majority of studies have indeed used these as predictors. A few studies have used only internal assessments for predicting student performance (Hasheminejad and Sarvmili 2018; Meier et al. 2016; Sivakumar and Selvaraj 2018). Some studies use external assessment like grades in other courses and prerequisites along with internal assessment records for a better result (Guo et al. 2015; Jishan et al. 2015; VeeraManickam et al. 2018; Yu et al. 2018). Some studies use data on student demographics and socio-economical factors along with internal assessment as well (Kotsiantis and Pintelas 2005; Koutina and Kermanidis 2011; Márquez-Vera et al. 2013; Natek and Zwilling 2014; Pandey and Taruna 2016; Santana et al. 2017; Uddin and Lee 2017; Xu et al. 2017). Researchers have also used student behavioural data captured from various online tools in this regard. For example, Hong et al. (2017), Lu et al. (2018) and Yang et al. (2018) have used video viewing behaviour for estimating student performance. Student behaviour in an online forum (Mueen et al. 2016; Ornelas and Ordóñez 2017; Romero et al. 2008; Widyahastuti and Tjhin 2017; Yoo and Kim 2014; Yu et al. 2018), learning management system (Conijn et al. 2017; Kim et al. 2018; Ostrow et al. 2015; Sandoval et al. 2018; Xing et al. 2015), their movement pattern (Zhang et al. 2018), and activity during web browsing (Chaturvedi and Ezeife 2017) also help in predicting student performance. Thai-Nghe et al. (2009) and Zollanvari et al. (2017) have used teaching quality and psychological factors of students for classifying student performance.

In addition to predictors, researchers have explored various methods for predicting student performance as well. Some of them (Hämäläinen and Vinni 2006; Hong



**Table 6** Classification studies during the tenure of the course

Study	Predictor used	Aim	Efficiency
Ahmed and Sadiq (2018)	Gender, internal assessment, attendance	Pass-fail	Acc: 83.56 %
Bresfelean et al. (2008)	Student profile and questionnaire	Pass-fail	Acc: 76.00 %
Cen et al. (2016)	Student interaction	Pass-fail	Acc: 80.00 %
Chaturvedi and Ezeife (2017)	Web usage pattern, internal assessment	Pass-Fail	Acc: 96.00 %
Chen et al. (2018)	Student interaction data	Final grade	Acc: 70.00 %
Grivokostopoulou et al. (2014)	Internal assessment, gender, year of study	Pass-fail	Acc: 97.80 %
Guo et al. (2015)	Background, internal and external assessment	Final grade	Acc: 77.20 %
Hämäläinen and Vinni (2006)	Performance in exercise	Pass-fail	Acc: 80.00 %
Hasan et al. (2018)	Student academic and activity data	Category	Acc: 100 %
Hasheminejad and Sarvmili (2018)	Internal assessment	Final score	Acc: 92.00 %
Hong et al. (2017)	Video viewing behaviour	Good-poor	AUC: 0.8100
Jishan et al. (2015)	Internal and external assessment	Final grade	Acc: 75.00 %
Kim et al. (2018)	Learning activity data	Pass-fail	AUC: 96.00 %
Koutina and Kermanidis (2011)	Student demographics, internal assessment	Category	Acc: 85.70 %
Meier et al. (2016)	Internal assessment	Good-poor	Acc: 76.00 %
Mueen et al. (2016)	Internal assessment and forum participation	Pass-Fail	Acc: 86.00 %
Natek and Zwilling (2014)	Internal assessment, demographics, extra-curricular activities	Final grade	Acc: 90.00 %
Ornelas and Ordonez (2017)	Student interaction data	Pass-fail	Acc: 90.00 %
Pandey and Taruna (2016)	Student demographics and assessment data	Pass-fail	Acc: 87.84 %
Quille and Bergin (2018)	Student background and psychological data	Pass-fail	Acc: 71.00 %
Romero et al. (2008)	Time spend in quiz, assignment	Final grade	Acc: 67.02 %
Santana et al. (2017)	Background, assessments	Pass-fail	FM: 0.8300
Sivakumar and Selvaraj (2018)	Internal assessment	Category	Acc: 98.56 %
Thai-Nghe et al. (2009)	Teaching quality	Pass-fail	FM: 0.6150
Uddin and Lee (2017)	Background, external assessment, social network interaction	Pass-fail	Acc: 70.00 %
Xu et al. (2017)	Background, internal assessment	Rank	Acc: 25.00 %
Yoo and Kim (2014)	Participation in online forum	Final grade	Acc: 91.40 %
Yu et al. (2018)	Usage pattern of online forum, internal and external assessment, attendance	Pass-fail	Acc: 78.08 %

**Table 6** (continued)

Study	Predictor used	Aim	Efficiency
Xing et al. (2015)	Log data from geometry learning tool	Pass-fail	Acc: 80.20 %
Zhang et al. (2018)	Student movement pattern	Good-poor	Acc: 95.00 %
Zollanvari et al. (2017)	Psychological factors	High-low	Acc: 82.00 %

Acc: Accuracy, AUC: Area under curve, FM: F-measure

et al. 2017; Huang and Fang 2013) observe that the applied method does not influence the prediction efficiency. In contrast, many studies have compared the efficiency of various algorithms and reported the superiority of a few among them. Quite a few studies have observed better result with decision tree classifier (Cen et al. 2016; Grivokostopoulou et al. 2014; Natek and Zwilling 2014; Chaturvedi and Ezeife 2017; Romero et al. 2008). The superiority of naïve bayes (Jishan et al. 2015; Ornelas and Ordóñez 2017), artificial neural network (Jishan et al. 2015; Widyahastuti and Tjhin

**Table 7** Regression studies during the tenure of the course

Study	Predictor used	Aim	Efficiency
Cen et al. (2016)	Student interaction	Final grade	MAE: 0.0800
Conijn et al. (2017)	Internal assessment and usage pattern of moodle	Final grade	$R^2$ : 0.4300
Huang and Fang (2013)	CGPA, internal and external assessment	Final score	APA: 90.10 %
Kotsiantis and Pintelas (2005)	Student demographics, internal assessment	Final-score	MAE: 1.21
Lu et al. (2018)	Internal assessment, video viewing behaviour	Final score	$R^2$ : 0.4000
Márquez-Vera et al. (2013)	Background, external assessments, socio-economic factors	Pass-fail	Acc: 98.70 %
Ostrow et al. (2015)	ASSISTment usage pattern	Final score	$R^2$ : 0.9610
Sandoval et al. (2018)	Background, academic records, behaviour in learning management system	Final score	MAE: 0.0615
VeeraManickam et al. (2018)	Internal and external assessment	Final score	RMSE: 4.665
Widyahastuti and Tjhin (2017)	Attendance, usage pattern of online forum	Final grade	MAE: 0.1185
Yang et al. (2018)	Internal assessment, video viewing behaviour	Final score	Acc: 80.00 %

Acc: Accuracy, APA: Average prediction accuracy, MAE: Mean absolute error, RMSE: Root mean square error

**Table 8** Classification studies before course commencement

Study	Predictor used	Aim	Efficiency
Chanlekha and Niramitranon (2018)	Student demographics, performance in school and admission	At-risk	Acc: 62.50 %
García and Mora (2011)	Student demographics, social life, and result in diagnostic test	Category	Acc: 58.64 %
Hidayah et al. (2013)	Intelligence level, interests, talents, motivation	Category	RMSE: 0.2561
Kabakchieva (2012)	Social and academic background, previous failures	Category	Acc: 73.59 %
Kabra and Bichkar (2011)	High school performance, family background	Pass-fail	Acc: 69.94 %
Martínez (2001)	High school performance, student background	Pass-fail	Acc: 62.60 %
Mimis et al. (2018)	Gender, motivation, previous performance	Final grade	Acc: 52.32 %
Mishra et al. (2014)	Social, academic background, emotional skills	Pass-fail	Acc: 94.42 %
Osmanbegović and Suljić (2012)	Social and academic background, attitudes towards study	Pass-fail	Acc: 76.65 %
Ramesh et al. (2013)	Gender, social and academic background	Final grade	Acc: 72.00 %
Sullivan et al. (2017)	Socio-economic factors	Below-above	Acc: 86.00 %

Acc: Accuracy, RMSE: Root mean square error

2017; Yu et al. 2018), support vector machine (Santana et al. 2017), random forest (Ahmed and Sadiq 2018; Chen et al. 2018; Hasan et al. 2018; Sandoval et al. 2018), deep learning (Guo et al. 2015; Kim et al. 2018) and linear regression (Lu et al. 2018; Yang et al. 2018) are also evident in the literature. Apart from the application of traditional data mining approaches, researchers have also proposed some specialised algorithms to predict student performance (Hasheminejad and Sarvmili 2018; Márquez-Vera et al. 2013; Meier et al. 2016; Uddin and Lee 2017; Xu et al. 2017; Zollanvari et al. 2017). Some of them have evolved from state-of-the-art algorithms after customisation, whereas a few are entirely new approach.

Moreover, these studies have reported significant efficiency. Some of them even achieved more than 90 % accuracy in classification settings (Chaturvedi and Ezeife 2017; Grivokostopoulou et al. 2014; Hasan et al. 2018; Hasheminejad and Sarvmili 2018; Márquez-Vera et al. 2013; Natek and Zwilling 2014; Ornelas and Ordóñez 2017; Sivakumar and Selvaraj 2018; Yoo and Kim 2014; Zhang et al. 2018). Some of these studies report high efficiency, but that may be due to the imbalanced dataset used. Unfortunately, many of the studies available in the literature ignore the class

**Table 9** Regression studies before course commencement

Study	Predictor used	Aim	Efficiency
Backenköhler and Wolf (2017)	CGPA, age, semester, number of attempt	Final grade	MSE: 0.4106
Bahritidinov and Sánchez (2017)	Domain knowledge, workload, availability	Final grade	RMSE: 0.7255
Cakmak (2017)	Grades in similar courses	Final grade	MAE: 0.2600
Elbadrawy et al. (2016)	Student demographics, content and grades of previous courses	Final score	RMSE: 0.7443
Gray et al. (2014)	Psychometric factors	Pass-fail	Acc: 73.30 %
Ibrahim and Rusli (2007)	Previous knowledge, performance, background	Final score	RASE: 0.1714
Pardos et al. (2010)	Skills identified through coding test	Final score	MAD: 4.210
Polyzou and Karypis (2016)	Grades in last semester	Final grade	RMSE: 0.6710
Sweeney et al. (2016)	Grades in previous semester	Final grade	RMSE: 0.7758

Acc: Accuracy, MAE: Mean absolute error, MSE: Mean square error, RASE: Root average squared error, RMSE: Root mean square error, MAD: Mean average deviance

imbalance issue. Although Hasan et al. (2018) have predicted all instances correctly, it is essential to mention here that they have used records of a few students only.

#### 4.1.2 Before course commencement

Although the EDM literature on performance prediction during the tenure of the course is mature enough, prediction before course commencement is still incipient. Here we present a meta-analysis of these studies, where Tables 8 and 9 provide a list of them with relevant details.

The different characteristics of predictors make the early prediction studies different from the former type of studies. In the present scenario, the predictors should be available before course commencement. Collecting the relevant data and their measurement is the major challenge here. As an example, Hidayah et al. (2013) tried to predict student performance based on their level of intelligence, interests, talents, and motivation. It collects the data using a questionnaire and quiz at the beginning of the semester. Gray et al. (2014) have used psychometric indicators like personality, motivation and learning strategies in a similar study. It is not feasible always to collect data through a questionnaire as it requires students' involvement. Therefore, EDM researchers have widely used the student demographics data and past performance in many early prediction studies (Backenköhler and Wolf 2017; Cakmak 2017; Chanlekha and Niramitranon 2018; Elbadrawy et al. 2016; Garcia and Mora 2011; Ibrahim and Rusli 2007; Kabakchieva 2012; Kabra and Bichkar 2011; Martinez 2001; Mimis et al. 2018; Polyzou and Karypis 2016; Sweeney et al. 2016).

Some studies have even utilised the domain knowledge for this purpose (Bahritidinov and Sánchez 2017; Ibrahim and Rusli 2007). Besides this, researchers have used socio-economic factors (Sullivan et al. 2017), prior skills (Pardos et al. 2010), and performance in the diagnostic test (Garcia and Mora 2011) as well.

Tables 8 and 9 present the multiple prediction studies which have applied classification and regression approach respectively. In one of these studies, Osmanbegović and Suljić (2012) observe that naïve bayes, artificial neural network, and decision tree classifier perform almost equally well for classification. In contrast, some of them have reported superiority of artificial neural network (Kabakchieva 2012; Mimis et al. 2018; Ramesh et al. 2013), support vector machine (Gray et al. 2014) and random forest (Mishra et al. 2014) algorithms. Besides this, the usage of advanced techniques like fuzzy neural network (Hidayah et al. 2013), matrix factorization (Sweeney et al. 2016), collaborative filtering (Cakmak 2017), and decision tree with entropy splitting (Sullivan et al. 2017) are also evident in the literature. Importantly, the efficiency of these studies is not as high as observed in prediction studies during the tenure of the course. In the case of classification, only Mishra et al. (2014) and Sullivan et al. (2017) have reported more than 80 % accuracy. However, Mishra et al. (2014) have used a dataset which contains only 215 records of single specialisation for analysis. Sullivan et al. (2017) have used socio-economic factors like free and reduced lunch receiving percentage, jobless rate, crime rate etc. However, the data for some of these factors may not be available in all circumstances.

## 4.2 Aim of prediction

In this section, we outline the aim or broad objective of the prediction studies discussed earlier. These studies have tried to achieve various objectives which can be grouped broadly in two categories. Some of them try to predict the student performance either in binary classes like pass-fail, success-failure etc., whereas others have attempted to predict the performance in terms of the final grade or actual score. Prediction in binary terms seems to be more popular amongst these two categories.

### 4.2.1 Success prediction

Modelling the actual grade or marks is a challenging task as it depends on diverse factors such as demographics, educational background, personal, psychological, academic progress and other environmental variables (Guo et al. 2015). Relationship between many of these variables are yet to understand and, therefore, performance classification turns out to be more popular among researchers. It is quite reasonable that a prediction is more straightforward with fewer target classes (Hu et al. 2017). Many researchers have tried to classify the student performance in binary terms such as pass-fail, below and above a reference level, good-poor etc. (Ahmed and Sadiq 2018; Cen et al. 2016; Chaturvedi and Ezeife 2017; Gray et al. 2014; Hong et al. 2017; Meier et al. 2016; Sullivan et al. 2017; Zollanvari et al. 2017). However, some studies try to classify in more generic terms with more than two class values (Hasan et al. 2018; Koutina and Kermanidis 2011; Sivakumar and Selvaraj 2018;

Garcia and Mora 2011). Importantly, Gardner and Brooks (2018) mention that non-parametric tree-based algorithms can be more effective in success prediction rather than techniques which require hyper-parameter tuning. Even the ensemble methods can provide a promising result as their predictive performance is generally high (Miguéis et al. 2018).

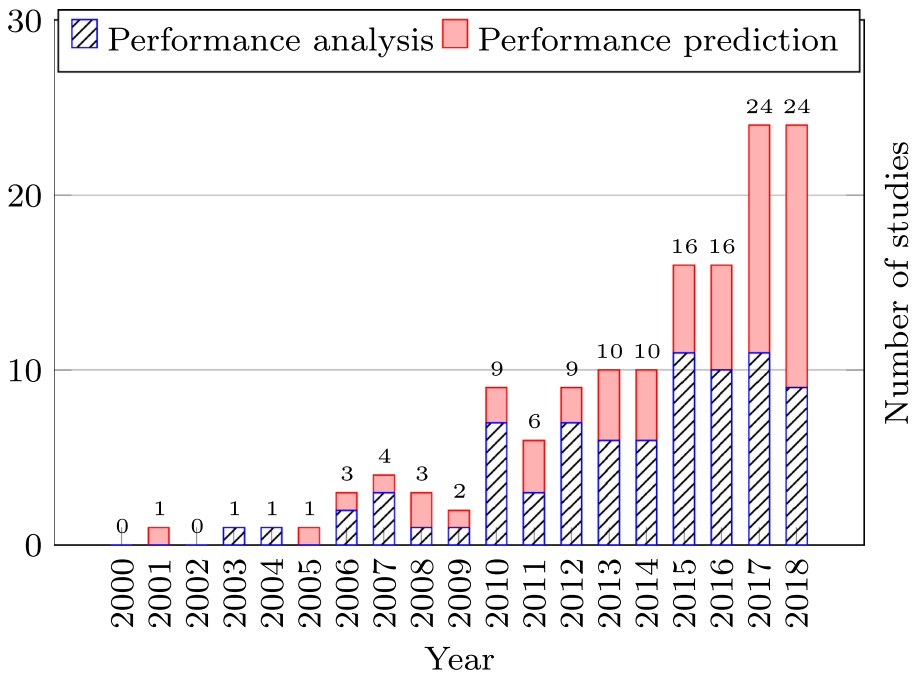
#### 4.2.2 Grade or marks prediction

As mentioned earlier, some researchers have even tried to predict actual marks (Guo et al. 2015; Huang and Fang 2013; Lu et al. 2018; Romero et al. 2008; Sandoval et al. 2018; Yang et al. 2018). The existing studies have used both regression and classification techniques for this purpose. A majority of them use the internal assessment as input data to predict grade or marks during the tenure of the course. A few prediction studies before course commencement tried to estimate the final grade or marks as well (Backenköhler and Wolf 2017; Bahritidinov and Sánchez 2017; Cakmak 2017; Elbadrawy et al. 2016; Mimis et al. 2018; Polyzou and Karypis 2016; Ramesh et al. 2013; Sweeney et al. 2016). However, they are not as efficient as the previous one. Moreover, Hu et al. (2017) observe that the number of studies which attempt to predict actual grade or score is comparatively less than the success prediction cases.

### 5 Discussion

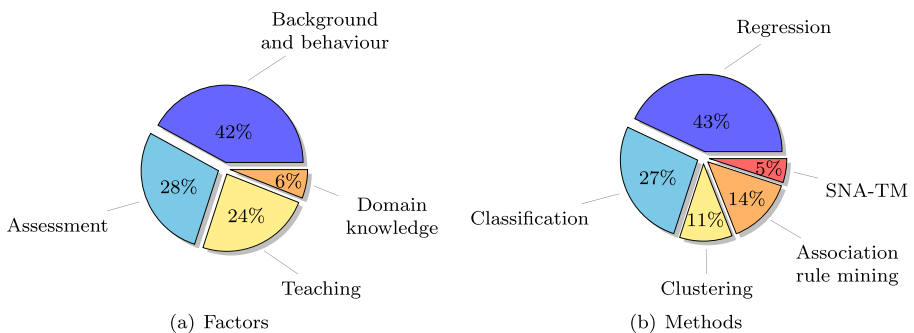
In this paper, we have cited 207 studies out of which 192 are published by EDM researchers either for analysing student performance or they are related to EDM techniques or applications. We have discussed the previous 16 surveys which are related to this topic and reviewed 140 core studies on student performance analysis or prediction in classroom-based education. The primary focus in 79 of them is to find performance influencing factors and 61 studies try to predict the student performance. Figure 3 presents the year wise count of these works. It shows that data mining studies on student performance have gained significant momentum in recent years, especially since 2010. Figure 4 additionally presents the percentage of various factors and methods used by 79 studies on student performance analysis. It seems that student background, behaviour and their internal or external assessments are most popular, whereas the researchers put the least focus on their domain knowledge. Regression and classification are found to be the most popular method used in these student performance related studies. Besides this, Fig. 5 shows the distribution of 61 performance prediction studies based on their objective and time of prediction. We find only 33 % of these studies capable of predicting student performance before course commencement. Furthermore, success or failure prediction turns out to be more popular than the prediction of final grade or score.

As mentioned earlier, predictor identification is one of the most critical tasks in any modelling study, and it is quite prevalent in student performance prediction as well (Helal et al. 2018). A handful of the existing literature explores the student-related factors such as their quality, how they perform in assignment and class test, whether they attend classes, what was their social and academic background etc. We

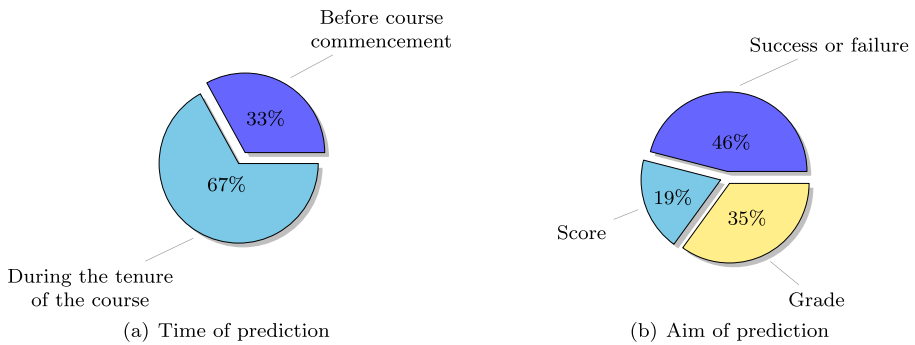


**Fig. 3** Year-wise number of studies

also found some conflicting observation about the influence of gender on student performance. However, we could not find any study which reports a significant influence of gender during prediction. There is no doubt that the quality of a student and their background creates quite a significant impact on student performance. It could be one of the possible reasons that the research community is extremely focused on student-related factors. However, there are many other external factors which could have a significant influence on performance. The teaching quality and previous domain knowledge are two among them. A few existing studies have found a positive impact



**Fig. 4** Percentage of various factors and methods used in student performance analysis



**Fig. 5** Percentage of performance prediction studies as per their aim and time of prediction

of teaching on student motivation. Some of them have even established a strong association between teaching quality and student performance. However, we have found only one prediction study which has used the teaching quality as a predictor. The domain knowledge of a student is also less utilised while predicting student performance. It is essential to mention here that the teaching quality and domain knowledge are difficult to measure. Although student evaluation of teaching excellence (SETE) can be a possible way to measure teaching quality, the research community is sceptical about the validity of the same. The possible bias between student performance and SETE is the primary reason behind this. It is worthy to point out that Khan and Ghosh (2018) have already found a strong association of student performance and SETE where the authors have taken utmost care to reduce the bias. Besides this, researchers have widely used questionnaire, observational techniques etc. for measuring domain knowledge. However, the data collected through these methods have some limitations (Conijn et al. 2017). The readily available course syllabus data can help in this regard.

Importantly, the performance modelling studies have utilised various factors as predictor, and the list of these predictors is not a short one (refer Tables 6 – 9). However, the question is which of them are the candidate of best predictor set. It is still an open research question. In order to estimate an answer, we have compared the effectiveness of various broad predictor types and presented the findings in Table 10. Significantly, the performance prediction studies evaluate their efficiency using various metrics. The prediction accuracy is most popular among them. There are, in fact, 40 out of 61 studies which have used prediction accuracy to measure the model efficiency. We, therefore, consider only these 40 studies for comparison. The findings show that student behaviour measured during the tenure of the course and internal assessment together can predict performance with an average accuracy of 88.02 %. The separate use of student behaviour and internal assessment is even capable of yielding more than 80 % average accuracy. The marks in the assignment, quiz, class test etc., attendance, learning behaviour, activities in various educational support tools are a few examples of such predictors. Although they turn out to be very useful during performance prediction, they can not help in early prediction before course commencement. As per the findings reported in Table 10, student background



**Table 10** Comparison of various predictors used in student performance prediction studies

Predictor type	No of studies *		Avg. Accuracy	
	After	Before	After	Before
Student background	0	1	–	86.00%
Student behaviour	8	1	81.95%	73.30%
Assessment	5	0	84.31%	–
Student background, behaviour	2	1	73.50%	94.42%
Student background, assessment	8	5	80.73%	68.13%
Student behaviour, assessment	5	0	88.02%	–
Student background, behaviour, assessment	1	3	70.00%	62.54%
Overall	29	11	82.07%	71.09%

\*40 out of 61 performance prediction studies measure efficiency in terms of accuracy

and behaviour like emotional skills are most effective while predicting student performance before course commencement. However, only one such study exists which has used these predictors together. The researchers have used external assessments in early prediction studies as well, but they are not as effective as internal assessment marks. The overall average accuracy observed in early prediction is only 71.09 % as compared to 82.07 % for studies predicting performance during the tenure of the course. More comprehensive research is therefore needed to search for more influencing factors which can be estimated early. The semi-structured or unstructured data generated by social media, IoT (internet of things) devices, student movement pattern, and logs of several educational support software may help in this regard. Although existing studies have widely used classification and regression techniques on the structured data, text mining on semi-structured or unstructured data can be the future research direction.

Furthermore, a majority of 46 % modelling studies prefer to classify the performance as success or failure. The efficiency of such classification is usually high due to a fewer number of class labels. A few of them even achieved more than 90 % prediction accuracy. Even so, many of these studies ignore the important concern of class imbalance. If there are only 5 % failure cases in an imbalanced dataset, for example, and none of the predictions is failure, then also the accuracy would be 95 %. Classification in terms of three or four categories is also present in literature. However, the number of such studies are very less. Some modelling studies try to classify the final grade or score as well. Although these studies report significant prediction efficiency during the tenure of the course, the reported efficiency before course commencement is comparatively lower. The final score prediction before course commencement also seems to be a research challenge. Undoubtedly, it is a complex task, as various influencing factors are not yet known. Still, students can get the benefit of such early prediction by taking advantage of the performance-oriented course recommendation (Bodily et al. 2018). It can additionally help the educators to identify at-risk students and take a proactive approach in mitigating the risk of retention (Gašević et al. 2016).

Finally, we also like to point out two crucial aspects observed during this survey. The first one is the widespread use of grade or score for finding performance influencing factors. The examination score is a widely accepted performance indicator which measures the knowledge in general (Pandey and Pal 2011). Many EDM studies, therefore, try to find influencing factors by establishing a correlation with the final grade. However, a superior student is supposed to perform better than a relatively inferior one. The final grade may not be sufficient to analyse the value addition due to external influencing factors. Khan and Ghosh (2018) have proposed a mechanism to measure the value addition, which can help to find the association of several external influencing factors. Secondly, the existing literature primarily focuses on improving the prediction accuracy but ignores the wrongly classified instances. For example, false-negative cases in failure prediction may lead to various problems. In such cases, the teacher may overlook a student who needs extra care but predicted to be successful. The prediction confidence may help in this regard. A statement like ‘I am 90 % confident that the student will fail’ may be more effective than a statement like ‘I think the student will fail’.

## 6 Conclusion

In this paper, we have presented a review of existing EDM literature on student performance analysis and prediction. It primarily focuses on studies related to classroom-based education. As an initial step, we have first identified the current and significant research question in this domain and additionally developed a taxonomy of the research directions as well. This systematic review has utilised the features of Mendeley, a reference management tool provided by Elsevier, for categorising and organising the relevant articles. We have adopted a hybrid approach for searching relevant articles. In addition to the conventional searching mechanism, the alert system of Google Scholar helps in this regard. This review presents a meta-analysis of performance influencing factors identified by the researchers. It also discusses the existing student performance prediction studies in the context of their aim and time of prediction. It seems that the literature on student performance prediction after course commencement is rich. However, the early prediction before course commencement is still an open challenge. Table 11 summarises the critical observations of this survey and also presents the recommendation.

This systematic survey would undoubtedly help the EDM researchers in advancing the field of next-term grade prediction. Such early prediction can facilitate the design of various intelligent computer-based educational systems as well. The performance-aware course recommendation system is one of such example among them. However, our survey includes research papers published in English only. Studies in other languages are not considered here due to our lack of proficiency in those languages. It may be possible that some significant work of student performance prediction before course commencement is available in other languages, but not mentioned in this survey. Moreover, we have used Google Scholar for searching the relevant studies. Any article not listed with Google Scholar is not a part of this survey.

**Table 11** Observations and recommendation for future research directions

---

**Observations and recommendations**

---

What are the factors influencing student performance in classroom learning?

Observation

- Association of final performance with student quality, their behaviour, family background, performance in earlier attended courses, quiz, midterm examination etc. are already established.
- It seems that teaching quality and domain knowledge also influences student performance.

Recommendation

- Future research can focus on measuring other influencing factors before course commencement and finding their impact on student performance.
- A superior student is expected to get better marks. Therefore, proper measure of value addition, in place of grade or marks, can help in establishing the impact of other factors.

Which methods are used for finding these factors?

Observation

- Regression and classification are the two methods widely used for establishing the impact of various factors.
- The usage of semi-structured or unstructured data is comparatively less for this purpose.

Recommendation

- The text mining of unstructured educational data in course syllabus, examination question paper, answer script etc. can further help in understanding various aspects of learning.
- The social network analysis can help in understanding the student behaviour and their impact on performance.

Is it possible to predict before course commencement?

Observation

- Some studies have tried to predict performance before course commencement. However, they are less efficient compared to the prediction studies during the tenure of the course.
- Internal assessment and behaviour seems to be most effective predictor of student performance during the tenure of the course. Unfortunately, these predictors are not available before course commencement.

Recommendation

- Future research may focus on student performance prediction before course commencement.
- Researchers should try to improve early prediction efficiency as well. In addition to student background and past performance, previous teaching quality of the teacher, domain knowledge and recent behaviour of the student can also be explored for this purpose.

Can we predict the actual grade or score?

Observation

- Researchers have primarily focused on classifying student success or grade. Success prediction is more popular and efficient among them.
- Only a few studies are available that tries to predict the final score.

Recommendation

- Future studies can put more focus on final grade or score prediction. Improving the success prediction efficiency before course commencement can be another future research direction.
-

## References

- Abrami, P.C., D'Apollonia, S., Rosenfield, S. (2007). The dimensionality of student ratings of instruction: what we know and what we do not. In *The scholarship of teaching and learning in higher education: an evidence-based perspective* (pp. 385–456): Springer.
- Adjei, S.A., Botelho, A.F., Heffernan, N.T. (2016). Predicting student performance on post-requisite skills using prerequisite skill data: an alternative method for refining prerequisite skill structures. In *Proceedings of the sixth international conference on learning analytics & knowledge* (pp. 469–473): ACM.
- Aghabozorgi, S., Mahrooian, H., Dutt, A., Wah, T.Y., Herawan, T. (2014). An approachable analytical study on big educational data mining. In *International conference on computational science and its applications* (pp. 721–737): Springer.
- Agrawal, R., Imieliński, T., Swami, A. (1993). Mining association rules between sets of items in large databases. *ACM SIGMOD Record*, 22(2), 207–216.
- Ahmed, N.S., & Sadiq, M.H. (2018). Clarify of the random forest algorithm in an educational field. In *2018 international conference on advanced science and engineering (ICOASE)* (pp. 179–184): IEEE.
- Ahmed, S., Paul, R., Hoque, M.L., Sayed, A. (2014). Knowledge discovery from academic data using association rule mining. In *2014 17th international conference on computer and information technology (ICCIT)* (pp. 314–319): IEEE.
- Akçapınar, G. (2015). How automated feedback through text mining changes plagiaristic behavior in online assignments. *Computers & Education*, 87, 123–130.
- Al-Barrak, M.A., & Al-Razgan, M. (2016). Predicting students final GPA using decision trees: a case study. *International Journal of Information and Education Technology*, 6(7), 528–533.
- Al-Obeidat, F., Tubaishat, A., Dillon, A., Shah, B. (2017). Analyzing students' performance using multi-criteria classification. *Cluster Computing*, 21(1), 623–632.
- Angeli, C., Howard, S., Ma, J., Yang, J., Kirschner, P.A. (2017). Data mining in educational technology classroom research: can it make a contribution? *Computers & Education*, 113, 226–242.
- Asif, R., Merceron, A., Ali, S.A., Haider, N.G. (2017). Analyzing undergraduate students' performance using educational data mining. *Computers & Education*, 113, 177–194.
- Backenköhler, M., & Wolf, V. (2017). Student performance prediction and optimal course selection: an MDP approach. In *International conference on software engineering and formal methods* (pp. 40–47): Springer.
- Bahrudinov, B., & Sánchez, E. (2017). Probabilistic classifiers and statistical dependency: the case for grade prediction. In *International work-conference on the interplay between natural and artificial computation* (pp. 394–403): Springer.
- Baker, R.S. (2014). Educational data mining: an advance for intelligent systems in education. *IEEE Intelligent Systems*, 29(3), 78–82.
- Baker, R.S., & Inventado, P.S. (2014). Educational data mining and learning analytics. In *Learning Analytics* (pp. 61–75): Springer.
- Baker, R.S.J.D., & Yacef, K. (2009). The state of educational data mining in 2009 : a review and future visions. *Journal of Educational Data Mining*, 1(1), 3–16.
- Bakhshinategh, B., Zaiane, O.R., ElAtia, S., Ipperciel, D. (2017). Educational data mining applications and tasks: a survey of the last 10 years. *Education and Information Technologies*, 23(1), 537–553.
- Balam, E.M., & Shannon, D.M. (2010). Student ratings of college teaching: a comparison of faculty and their students. *Assessment & Evaluation in Higher Education*, 35(2), 209–221.
- Bayer, J., Bydzovská, H., Géryk, J., Obsivac, T., Popelinsky, L. (2012). Predicting drop-out from social behaviour of students. In *International conference on educational data mining (EDM)*.
- Beemer, J., Spoon, K., He, L., Fan, J., Levine, R.A. (2018). Ensemble learning for estimating individualized treatment effects in student success studies. *International Journal of Artificial Intelligence in Education*, 28(3), 315–335.
- Bendikson, L., Hattie, J., Robinson, V. (2011). Identifying the comparative academic performance of secondary schools. *Journal of Educational Administration*, 49(4), 433–449.
- Berkhin, P. (2006). A survey of clustering data mining techniques. In *Grouping multidimensional data* (pp. 25–71): Springer.
- Bloom, B.S., Englehard, M., Furst, E., Hill, W., Krathwohl, D. (1956). Taxonomy of educational objectives: the classification of educational goals. Handbook I: Cognitive Domain.

- Bodily, R., Ikahihifo, T.K., Mackley, B., Graham, C.R. (2018). The design, development, and implementation of student-facing learning analytics dashboards. *Journal of Computing in Higher Education*, 30(3), 572–598.
- Bogarin, A., Romero, C., Cerezo, R., Sánchez-Santillan, M. (2014). Clustering for improving educational process mining. In *Proceedings of the fourth international conference on learning analytics and knowledge* (pp. 11–15): ACM.
- Bresfelean, V.P., Bresfelean, M., Ghisoiu, N., Comes, C.A. (2008). Determining students' academic failure profile founded on data mining methods. In *ITI 2008 - 30th international conference on information technology interfaces* (pp. 317–322): IEEE.
- Brocato, B.R., Bonanno, A., Ulbig, S. (2015). Student perceptions and instructional evaluations: a multivariate analysis of online and face-to-face classroom settings. *Education and Information Technologies*, 20(1), 37–55.
- Bucos, M., & Drăgulescu, B. (2018). Predicting student success using data generated in traditional educational environments. *TEM Journal*, 7(3), 617–625.
- Buldu, A., & Üçgün, K. (2010). Data mining application on students' data. *Procedia-Social and Behavioral Sciences*, 2(2), 5251–5259.
- Buniyamin, N., Mat, U.B., Arshad, P.M. (2015). Educational data mining for prediction and classification of engineering students achievement. In *IEEE 7th international conference on engineering education ICEED 2015* (pp. 49–53).
- Bydžovská, H. (2016). A comparative analysis of techniques for predicting student performance. In *Proceedings of the 9th international conference on educational data mining*.
- Cakmak, A. (2017). Predicting student success in courses via collaborative filtering. *International Journal of Intelligent Systems and Applications in Engineering*, 5(1), 10–17.
- Campagni, R., Merlini, D., Sprugnoli, R., Verri, M.C. (2015). Data mining models for student careers. *Expert Systems with Applications*, 42(13), 5508–5521.
- Carter, A.S., Hundhausen, C.D., Adesope, O. (2017). Blending measures of programming and social behavior into predictive models of student achievement in early computing courses. *ACM Transactions on Computing Education*, 17(3), 12.
- Cen, L., Ruta, D., Powell, L., Hirsch, B., Ng, J. (2016). Quantitative approach to collaborative learning: Performance prediction, individual assessment, and group composition. *International Journal of Computer-Supported Collaborative Learning*, 11(2), 187–225.
- Centra, J.A. (2003). Will teachers receive higher student evaluations by giving higher grades and less course work? *Research in Higher Education*, 44(5), 495–518.
- Chanlekha, H., & Niramitranon, J. (2018). Student performance prediction model for early-identification of at-risk students in traditional classroom settings. In *Proceedings of the 10th international conference on management of digital ecosystems - MEDES '18* (pp. 239–245): ACM.
- Chatterjee, S., & Hadi, A.S. (2015). *Regression analysis by example*. New York: Wiley.
- Chaturvedi, R., & Ezeife, C. (2013). Mining the impact of course assignments on student performance. In *Educational data mining 2013*.
- Chaturvedi, R., & Ezeife, C.I. (2017). Predicting student performance in an ITS using task-driven features. In *2017 IEEE international conference on computer and information technology (CIT)* (pp. 168–175): IEEE.
- Chen, L., Wang, S., Wang, K., Zhu, J. (2016). Soft subspace clustering of categorical data with probabilistic distance. *Pattern Recognition*, 51, 322–332.
- Chen, W., Brinton, C.G., Cao, D., Mason-singh, A., Lu, C., Chiang, M. (2018). Early detection prediction of learning outcomes in online short-courses via learning behaviors. *IEEE Transactions on Learning Technologies*, 12(1), 44–58.
- Chen, Y., Liu, Q., Huang, Z., Wu, L., Chen, E., Wu, R., et al. (2017). Tracking knowledge proficiency of students with educational priors. In *Conference on information and knowledge management (CIKM)* (pp. 989–998).
- Christian, T.M., & Ayub, M. (2014). Exploration of classification using NBTree for predicting students' performance. In *2014 international conference on data and software engineering (ICODSE)* (pp. 1–6): IEEE.
- Chrysafiadi, K., & Virvou, M. (2013). Student modeling approaches: a literature review for the last decade. *Expert Systems with Applications*, 40(11), 4715–4729.
- Chung, H., & Kim, J. (2016). An ontological approach for semantic modeling of curriculum and syllabus in higher education. *International Journal of Information and Education Technology*, 6(5), 365.

- Conijn, R., Snijders, C., Kleingeld, A., Matzat, U. (2017). Predicting student performance from LMS data: a comparison of 17 blended courses using Moodle LMS. *IEEE Transactions on Learning Technologies*, 10(1), 17–29.
- Cooper, H.M. (1988). Organizing knowledge syntheses: a taxonomy of literature reviews. *Knowledge, Technology & Policy*, 1(1), 104–126.
- Damaševičius, R. (2010). Analysis of academic results for informatics course improvement using association rule mining. In *Information systems development* (pp. 357–363). Berlin: Springer.
- Daud, A., Aljohani, N.R., Abbasi, R.A., Lytras, M.D., Abbas, F., Alowibdi, J.S. (2017). Predicting student performance using advanced learning analytics. In *Proceedings of the 26th international conference on world wide web companion* (pp. 415–421): International World Wide Web Conferences Steering Committee.
- Delen, D., & Crossland, M.D. (2008). Seeding the survey and analysis of research literature with text mining. *Expert Systems with Applications*, 34(3), 1707–1720.
- Dutt, A., Ismail, M.A., Herawan, T. (2017). A systematic review on educational data mining. *IEEE Access*, 5, 15991–16005.
- Dvorak, T., & Jia, M. (2016). Do the timeliness, regularity, and intensity of online work habits predict academic performance? *Journal of Learning Analytics*, 3(3), 318–330.
- Elbadrawy, A., Polyzou, A., Ren, Z., Sweeney, M., Karypis, G., Rangwala, H. (2016). Predicting student performance using personalized analytics. *Computer*, 49(4), 61–69.
- Elouazizi, N., Birol, G., Jandciu, E., Öberg, G., Welsh, A., Han, A., et al. (2017). Automated analysis of aspects of written argumentation. In *Proceedings of the seventh international learning analytics and knowledge conference on - lak '17* (pp. 606–607): ACM.
- Fausett, L.V., & Elwasif, W. (1994). Predicting performance from test scores using backpropagation and counterpropagation. In *IEEE international conference on neural networks*, (Vol. 5 pp. 3398–3402): IEEE.
- Feldman, R., & Sanger, J. (2007). *The text mining handbook: advanced approaches in analyzing unstructured data*. Cambridge: Cambridge University Press.
- Felisoni, D.D., & Godoi, A.S. (2018). Cell phone usage and academic performance: An experiment. *Computers & Education*, 117, 175–187.
- Ferguson, R. (2012). Learning analytics: drivers, developments and challenges. *International Journal of Technology Enhanced Learning*, 4(5/6), 304–317.
- Fernandes, E., Carvalho, R., Holanda, M., Van Erven, G. (2017). Educational data mining: discovery standards of academic performance by students in public high schools in the Federal District of Brazil. In *World conference on information systems and technologies* (pp. 287–296).
- Fernandes, E., Holanda, M., Victorino, M., Borges, V., Carvalho, R., Erven, G.V. (2018). Educational data mining: predictive analysis of academic performance of public school students in the capital of Brazil. *Journal of Business Research*, 94, 335–343.
- Figlio, D.N., & Lucas, M.E. (2004). Do high grading standards affect student performance? *Journal of Public Economics*, 88(9–10), 1815–1834.
- Foley, J., & Allan, J. (2016). Retrieving hierarchical syllabus items for exam question analysis. In *European conference on information retrieval* (pp. 575–586). Cham: Springer.
- Galbraith, C.S., Merrill, G.B., Kline, D.M. (2012). Are student evaluations of teaching effectiveness valid for measuring student learning outcomes in business related classes? A neural network and Bayesian analyses. *Research in Higher Education*, 53(3), 353–374.
- García, E., Romero, C., Ventura, S., Calders, T. (2007). Drawbacks and solutions of applying association rule mining in learning management systems. In *Proceedings of the international workshop on applying data mining in e-learning (ADML 2007), Crete, Greece* (pp. 13–22).
- García, E.P.I., & Mora, P.M. (2011). Model prediction of academic performance for first year students. In *2011 10th Mexican international conference on artificial intelligence* (pp. 169–174): IEEE.
- Gardner, J., & Brooks, C. (2018). Evaluating predictive models of student success: closing the methodological gap. *Journal of Learning Analytics*, 5(2), 105–125.
- Gašević, D., Dawson, S., Rogers, T., Gasevic, D. (2016). Learning analytics should not promote one size fits all: the effects of instructional conditions in predicting academic success. *The Internet and Higher Education*, 28, 68–84.
- Gasevic, D., Jovanovic, J., Pardo, A., Dawson, S. (2017). Detecting learning strategies with analytics: links with self-reported measures and academic performance. *Journal of Learning Analytics*, 4(2), 113–128.

- Gedeon, T.D., & Turner, S. (1993). Explaining student grades predicted by a neural network. In *Proceedings of 1993 international joint conference on neural networks, 1993. IJCNN'93-Nagoya*, (Vol. 1 pp. 609–612): IEEE.
- Golding, P., & Donaldson, O. (2006). Predicting academic performance. In *Frontiers in education conference, 36th Annual* (pp. 21–26): IEEE.
- Goos, M., & Salomons, A. (2016). Measuring teaching quality in higher education: assessing selection bias in course evaluations. *Research in Higher Education*, 58(4), 341–364.
- Gowda, S.M., Baker, R.S., Corbett, A.T., Rossi, L.M. (2013). Towards automatically detecting whether student learning is shallow. *International Journal of Artificial Intelligence in Education*, 23(1–4), 50–70.
- Gray, G., McGuinness, C., Owende, P. (2014). An application of classification models to predict learner progression in tertiary education. In *Advance Computing Conference (IACC), 2014 IEEE International* (pp. 549–554): IEEE.
- Grivokostopoulou, F., Perikos, I., Hatzilygeroudis, I. (2014). Utilizing semantic web technologies and data mining techniques to analyze students learning and predict final performance. In *2014 IEEE international conference on teaching, assessment and learning for engineering (TALE)* (pp. 488–494): IEEE.
- Guarin, C.E.L., Guzman, E.L., Gonzalez, F.A. (2015). A model to predict low academic performance at a specific enrollment using data mining. *Revista Iberoamericana de Tecnologías del Aprendizaje*, 10(3), 119–125.
- Guo, B., Zhang, R., Xu, G., Shi, C., Yang, L. (2015). Predicting students performance in educational data mining. In *International symposium on educational technology, ISET 2015* (pp. 125–128).
- Guruler, H., Istanbulu, A., Karahasan, M. (2010). A new student performance analysing system using knowledge discovery in higher educational databases. *Computers & Education*, 55(1), 247–254.
- Hämäläinen, W., & Vinni, M. (2006). Comparison of machine learning methods for intelligent tutoring systems. In *International conference on intelligent tutoring systems* (pp. 525–534): Springer.
- Hart, S., Daucourt, M., Ganley, C. (2017). Individual differences related to college students' course performance in calculus II. *Journal of Learning Analytics*, 4(2), 129–153.
- Hasan, R., Palaniappan, S., Raziff, A.R.A., Mahmood, S., Sarker, K.U. (2018). Student academic performance prediction by using decision tree algorithm. In *2018 4th international conference on computer and information sciences (ICCOINS)* (pp. 1–5): IEEE.
- Hasheminejad, H., & Sarvmili, M. (2018). S3PSO: students' performance prediction based on particle swarm optimization. *Journal of AI and Data Mining*, 7(1), 77–96.
- Hassan, O.R., & Rasiah, R. (2011). Poverty and student performance in Malaysia. *International Journal of Institutions and Economies*, 3(1), 61–76.
- Hattie, J., & Clinton, J. (2012). Physical activity is not related to performance at school. *Archives of Pediatrics & Adolescent Medicine*, 166(7), 678–679.
- Hattie, J., & Timperley, H. (2007). The power of feedback. *Review of Educational Research*, 77(1), 81–112.
- Helal, S., Li, J., Liu, L., Ebrahimie, E., Dawson, S., Murray, D.J. (2018). Identifying key factors of student academic performance by subgroup discovery. *International Journal of Data Science and Analytics*, 7(3), 227–245.
- Hidayah, I., Permanasari, A.E., Ratwastuti, N. (2013). Student classification for academic performance prediction using neuro fuzzy in a conventional classroom. In *2013 international conference on information technology and electrical engineering (ICITEE)* (pp. 221–225): IEEE.
- Hong, B., Wei, Z., Yang, Y. (2017). Online education performance prediction via time-related features. In *2017 IEEE/ACIS 16th international conference on computer and information science (ICIS)* (pp. 95–100): IEEE.
- Hsiao, I.H., & Lin, Y.L. (2017). Enriching programming content semantics: an evaluation of visual analytics approach. *Computers in Human Behavior*, 72, 771–782.
- Hsiao, I.H., Pandhalkudi Govindarajan, S.K., Lin, Y.L. (2016). Semantic visual analytics for today's programming courses. In *Proceedings of the sixth international conference on learning analytics and knowledge* (pp. 48–53): ACM.
- Hu, X., Cheong, C.W.L., Ding, W., Woo, M. (2017). A systematic review of studies on predicting student learning outcomes using learning analytics. In *Proceedings of the seventh international learning analytics & knowledge conference* (pp. 528–529): ACM.



- Huang, S., & Fang, N. (2013). Predicting student academic performance in an engineering dynamics course: a comparison of four types of predictive mathematical models. *Computers & Education*, 61, 133–145.
- Ibrahim, Z., & Rusli, D. (2007). Predicting students' academic performance: comparing artificial neural network, decision tree and linear regression. In *21st Annual SAS Malaysia Forum* (pp. 1–6).
- Ivančević, V., Čeliković, M., Luković, I. (2010). Analyzing student spatial deployment in a computer laboratory. In *Educational data mining* (p. 2011).
- Jara, M., & Mellar, H. (2010). Quality enhancement for e-learning courses: the role of student feedback. *Computers & Education*, 54(3), 709–714.
- Jishan, S.T., Rashu, R.I., Haque, N., Rahman, R.M. (2015). Improving accuracy of students' final grade prediction model using optimal equal width binning and synthetic minority over-sampling technique. *Decision Analytics*, 2(1), 1.
- Kabakchieva, D. (2012). Student performance prediction by using data mining classification algorithms. *International Journal of Computer Science and Management Research*, 1(4), 686–690.
- Kabra, R.R., & Bichkar, R.S. (2011). Performance prediction of engineering students using decision trees. *International Journal of Computer Applications*, 36(11), 975–8887.
- Kagdi, H., Collard, M.L., Maletic, J.I. (2007). A survey and taxonomy of approaches for mining software repositories in the context of software evolution. *Journal of Software Maintenance and Evolution: Research and Practice*, 19(2), 77–131.
- Kamley, S., Jaloree, S., Thakur, R.S. (2016). A review and performance prediction of students' using association rule mining based approach. *Data Mining and Knowledge Engineering*, 8(8), 252–259.
- Kaufman, L., & Rousseeuw, P.J. (2009). *Finding groups in data: an introduction to cluster analysis* Vol. 344. New York: Wiley.
- Kaviyarasi, R., & Balasubramanian, T. (2018). Exploring the high potential factors that affects students' academic performance. *International Journal of Education and Management Engineering*, 8(6), 15.
- Kesavaraj, G., & Sukumaran, S. (2013). A study on classification techniques in data mining. In *2013 fourth international conference on computing, communications and networking technologies (ICCCNT)* (pp. 1–7): IEEE.
- Khan, A., & Ghosh, S.K. (2016). Analysing the impact of poor teaching on student performance. In *2016 IEEE international conference on teaching, assessment, and learning for engineering (TALE)* (pp. 169–175): IEEE.
- Khan, A., & Ghosh, S.K. (2018). Data mining based analysis to explore the effect of teaching on student performance. *Education and Information Technologies*, 23(4), 1677–1697.
- Khanna, L., Singh, S.N., Alam, M. (2016). Educational data mining and its role in determining factors affecting students academic performance: a systematic review. In *2016 1st India international conference on information processing (IICIP)* (pp. 1–7): IEEE.
- Kim, B.H., Vizitei, E., Ganapathi, V. (2018). GritNet: student performance prediction with deep learning. arXiv:1804.07405.
- Koedinger, K.R., D'Mello, S., McLaughlin, E.A., Pardos, Z.A., Rosé, C.P. (2015). Data mining and education. *Wiley Interdisciplinary Reviews: Cognitive Science*, 6(4), 333–353.
- Kotsiantis, S.B., & Pintelas, P.E. (2005). Predicting students marks in hellenic open university. In *Fifth IEEE international conference on advanced learning technologies, 2005. ICAALT 2005* (pp. 664–668): IEEE.
- Koutina, M., & Kermanidis, K.L. (2011). Predicting postgraduate students' performance using machine learning techniques. In *IFIP international conference on artificial intelligence applications and innovations* (pp. 159–168): Springer.
- Kumar, D.A., Selvam, R.P., Kumar, K.S. (2018). Review on prediction algorithms in educational data mining. *International Journal of Pure and Applied Mathematics*, 118(8), 531–537.
- Kumar, M., Singh, A.J., Handa, D. (2017). Literature survey on student's performance prediction in education using data mining techniques. *International Journal of Education and Management Engineering*, 6, 40–49.
- Li, K.F., Rusk, D., Song, F. (2013). Predicting student academic performance. In *2013 seventh international conference on complex, intelligent, and software intensive systems (cisis)* (pp. 27–33): IEEE.
- Lin, C.H., Kwon, J.B., Zhang, Y. (2018). Online self-paced high-school class size and student achievement. *Educational Technology Research and Development*, pp 1–20.



- Lipovetsky, S., & Conklin, W.M. (2015). Predictor relative importance and matching regression parameters. *Journal of Applied Statistics*, 42(5), 1017–1031.
- Liu, Q., Huang, Z., Huang, Z., Liu, C., Chen, E., Su, Y., et al. (2018a). Finding similar exercises in online education systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1821–1830): ACM.
- Liu, Q., Wu, R., Chen, E., Xu, G., Su, Y., Chen, Z., et al. (2018b). Fuzzy cognitive diagnosis for modelling examinee performance. *ACM Transactions on Intelligent Systems and Technology*, 9(4), 48.
- Livieris, I.E., Drakopoulou, K., Mikropoulos, T.A., Tampakas, V., Pintelas, P. (2018). An ensemble-based semi-supervised approach for predicting students' performance. In *Research on e-Learning and ICT in Education* (pp. 25–42): Springer.
- Loh, C.S., & Sheng, Y. (2015). Measuring the (dis-)similarity between expert and novice behaviors as serious games analytics. *Education and Information Technologies*, 20(1), 5–19.
- Lorenzetti, C., Maguitman, A., Leake, D., Menczer, F., Reichherzer, T. (2016). Mining for topics to suggest knowledge model extensions. *ACM Transactions on Knowledge Discovery from Data*, 11(2), 23.
- Lu, O.H.T., Huang, A.Y.Q., Huang, J.C., Lin, A.J.Q., Ogata, H., Yang, S.J.H. (2018). Applying learning analytics for the early prediction of students' academic performance in blended learning. *Journal of Educational Technology and Society*, 21(2), 220–232.
- Macfadyen, L.P., Dawson, S., Prest, S., Gašević, D. (2015). Whose feedback? A multilevel analysis of student completion of end-of-term teaching evaluations. *Assessment & Evaluation in Higher Education*, 41(6), 821–839.
- Márquez-Vera, C., Cano, A., Romero, C., Ventura, S. (2013). Predicting student failure at school using genetic programming and different data mining approaches with high dimensional and imbalanced data. *Applied Intelligence*, 38(3), 315–330.
- Marsh, H.W. (1984). Students' evaluations of university teaching: dimensionality, reliability, validity, potential biases, and utility. *Journal of Educational Psychology*, 76(5), 707.
- Marsh, H.W. (2007). Students' evaluations of university teaching: Dimensionality, reliability, validity, potential biases and usefulness. In *The scholarship of teaching and learning in higher education: An evidence-based perspective* (pp. 319–383): Springer.
- Martinez, D. (2001). Predicting Student Outcomes Using Discriminant Function Analysis.
- Mat, U.B., Buniyamin, N., Arsad, P.M., Kassim, R. (2013). An overview of using academic analytics to predict and improve students' achievement: a proposed proactive intelligent intervention. In *2013 IEEE 5th conference on engineering education (ICEED)* (pp. 126–130): IEEE.
- McCarthy, K.S., & Goldman, S.R. (2017). Constructing interpretive inferences about literary text: the role of domain-specific knowledge. *Learning and Instruction*, 60, 245–251.
- Meier, Y., Xu, J., Atan, O., van der Schaar, M. (2016). Predicting grades. *IEEE Transactions on Signal Processing*, 64(4), 959–972.
- Miguéis, V., Freitas, A., Garcia, P.J., Silva, A. (2018). Early segmentation of students according to their academic performance: a predictive modelling approach. *Decision Support Systems*, 115, 36–51.
- Mimis, M., El Hajji, M., Es-saady, Y., Ouelid Guejdi, A., Douzi, H., Mammass, D. (2018). A framework for smart academic guidance using educational data mining. *Education and Information Technologies*, 24(2), 1379–1393.
- Mishra, T., Kumar, D., Gupta, S. (2014). Mining students' data for prediction performance. In *International conference on advanced computing and communication technologies, ACCT* (pp. 255–262).
- Moher, D., Liberati, A., Tetzlaff, J., Altman, D.G. (2009). Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Annals of Internal Medicine*, 151(4), 264–269.
- Montuschi, P., Lamberti, F., Gatteschi, V., Demartini, C. (2015). A semantic recommender system for adaptive learning. *IT Professional*, 17(5), 50–58.
- Moore, S., & Kuol, N. (2005). Students evaluating teachers: exploring the importance of faculty reaction to feedback on teaching. *Teaching in Higher Education*, 10(1), 57–73.
- Moos, D.C., & Azevedo, R. (2008). Self-regulated learning with hypermedia: the role of prior domain knowledge. *Contemporary Educational Psychology*, 33(2), 270–298.
- Mueen, A., Zafar, B., Manzoor, U. (2016). Modeling and predicting students' academic performance using data mining techniques. *International Journal of Modern Education and Computer Science*, 8(11), 36.
- Nakayama, M. (2016). Lexical analysis of syllabi in the area of technology enhanced learning. In *2016 15th international conference on information technology based higher education and training (ITHET)* (pp. 1–5): IEEE.

- Natek, S., & Zwilling, M. (2014). Student data mining solution-knowledge management system related to higher education institutions. *Expert Systems with Applications*, 41(14), 6400–6407.
- Nikolic, S., Ritz, C., Vial, P.J., Ros, M., Stirling, D. (2015). Decoding student satisfaction: How to manage and improve the laboratory experience. *IEEE Transactions on Education*, 58(3), 151–158.
- O’Connell, K.A., Wostl, E., Crosslin, M., Berry, T.L., Grover, J.P. (2018). Student ability best predicts final grade in a college algebra course. *Journal of Learning Analytics*, 5(3), 167–181.
- Ogor, E.N. (2007). Student academic performance monitoring and evaluation using data mining techniques. In *Electronics, robotics and automotive mechanics conference* (pp. 354–359): IEEE.
- Ornelas, F., & Ordonez, C. (2017). Predicting student success: a naïve Bayesian application to community college data. *Technology, Knowledge and Learning*, 22(3), 299–315.
- Osmanbegović, E., & Suljić, M. (2012). Data mining approach for predicting student performance. *Economic Review*, 10(1), 3–12.
- Ostrow, K., Donnelly, C., Heffernan, N. (2015). Optimizing partial credit algorithms to predict student performance. In *International conference on educational data mining (EDM)*.
- Pal, S., & Chaurasia, V. (2017). Is alcohol affect higher education students performance: searching and predicting pattern using data mining algorithms. *International Journal of Innovations & Advancement in Computer Science IJIACS ISSN*, 6(4), 2347–8616.
- Pandey, M., & Taruna, S. (2016). Towards the integration of multiple classifier pertaining to the Student’s performance prediction. *Perspectives in Science*, 8, 364–366.
- Pandey, U.K., & Pal, S. (2011). A data mining view on class room teaching language. arXiv:1104.4164.
- Papamitsiou, Z.K., Terzis, V., Economides, A.A. (2014). Temporal learning analytics for computer based testing. In *Proceedings of the fourth international conference on learning analytics and knowledge* (pp. 31–35): ACM.
- Parack, S., Zahid, Z., Merchant, F. (2012). Application of data mining in educational databases for predicting academic trends and patterns. In *2012 IEEE international conference on technology enhanced education (ICTEE)* (pp. 1–4): IEEE.
- Pardos, Z.A., Heffernan, N.T., Anderson, B., Heffernan, C.L., Schools, W.P. (2010). Using fine-grained skill models to fit student performance with Bayesian networks. *Handbook of Educational Data Mining*, 417.
- Peña-Ayala, A. (2014). Educational data mining: a survey and a data mining-based analysis of recent works. *Expert Systems with Applications*, 41(4), 1432–1462.
- Polyzou, A., & Karypis, G. (2016). Grade prediction with course and student specific models. In *Pacific-asia conference on knowledge discovery and data mining* (pp. 89–101). Cham: Springer.
- Polyzou, A., & Karypis, G. (2019). Feature extraction for next-term prediction of poor student performance. *IEEE Transactions on Learning Technologies*, 12(2), 237–248.
- Pong-Inwong, C., & Rungworawut, W. (2012). Teaching evaluation using data mining on moodle LMS forum. In *2012 6th international conference on new trends in information science, service science and data mining (ISSDM2012)* (pp. 550–555): IEEE.
- Price, L., Svensson, I., Borell, J., Richardson, J.T.E. (2017). The role of gender in students’ ratings of teaching quality in computer science and environmental engineering. *IEEE Transactions on Education*, 60(4), 281–287.
- Quadri, M.M.N., & Kalyankar, N.V. (2010). Drop out feature of student data for academic performance using decision tree techniques. *Global Journal of Computer Science and Technology*, 10(2).
- Quille, K., & Bergin, S. (2018). Programming: Predicting student success early in CS1. A re-validation and replication study. In *Proceedings of the 23rd annual ACM conference on innovation and technology in computer science education* (pp. 15–20): ACM.
- Quinlan, J.R. (1990). Decision trees and decision-making. *IEEE Transactions on Systems, Man, and Cybernetics*, 20(2), 339–346.
- Ramesh, V., Parkavi, P., Ramar, K. (2013). Predicting student performance: a statistical and data mining approach. *International Journal of Computer Applications*, 63(8), 35–39.
- Rani, S., & Kumar, P. (2017). A sentiment analysis system to improve teaching and learning. *Computer*, 50(5), 36–43.
- Rekha, R., Angadi, A., Pathak, A., Kapur, A., Gosar, H., Ramanathan, M., et al. (2012). Ontology driven framework for assessing the syllabus fairness of a question paper. In *2012 IEEE international conference on technology enhanced education (ICTEE)* (pp. 1–5): IEEE.
- Romero, C., López, M.I., Luna, J.M., Ventura, S. (2013). Predicting students’ final performance from participation in on-line discussion forums. *Computers and Education*, 68, 458–472.

- Romero, C., & Ventura, S. (2007). Educational data mining: a survey from 1995 to 2005. *Expert Systems with Applications*, 33(1), 135–146.
- Romero, C., & Ventura, S. (2010). Educational data mining: a review of the state of the art. *IEEE Transactions on Systems, Man, and Cybernetics, Part C: Applications and Reviews*, 40(6), 601–618.
- Romero, C., Ventura, S., Espejo, P.G., Hervás, C. (2008). Data mining algorithms to classify students. In *Educational data mining 2008*.
- Saarela, M., & Kärkkäinen, T. (2015). Analysing student performance using sparse data of core bachelor courses. *Journal of Educational Data Mining*, 7(1), 3–32.
- Sandoval, A., Gonzalez, C., Alarcon, R., Pichara, K., Montenegro, M. (2018). Centralized student performance prediction in large courses based on low-cost variables in an institutional context. *The Internet and Higher Education*, 37, 76–89.
- Santana, M.A., Costa, E.B., Fonseca, B., Rego, J., de Araújo, F.F. (2017). Evaluating the effectiveness of educational data mining techniques for early prediction of students' academic failure in introductory programming courses. *Computers in Human Behavior*, 73, 247–256.
- Saxena, P.S., & Govil, M.C. (2009). Prediction of student's academic performance using clustering. In *National conference on cloud computing & big data* (pp. 1–6).
- Shahiri, A.M., Husain, W., Rashid, A.N. (2015). A review on predicting student's performance using data mining techniques. *Procedia Computer Science*, 72, 414–422.
- She, H.C., Cheng, M.T., Li, T.W., Wang, C.Y., Chiu, H.T., Lee, P.Z., et al. (2012). Web-based undergraduate chemistry problem-solving: the interplay of task performance, domain knowledge and web-searching strategies. *Computers & Education*, 59(2), 750–761.
- Shingari, I., & Kumar, D. (2018). A survey on various aspects of education data mining in predicting student performance. *Journal of Applied Science and Computations*, 5(6), 38–42.
- Siemens, G., & Baker, R.S.J.D. (2012). Learning analytics and educational data mining: towards communication and collaboration. In *Proceedings of the 2nd international conference on learning analytics and knowledge* (pp. 252–254): ACM.
- Sivakumar, S., & Selvaraj, R. (2018). Predictive modeling of students performance through the enhanced decision tree. In *Advances in electronics, communication and computing* (pp. 21–36). Singapore: Springer.
- Stronge, J.H., Ward, T.J., Tucker, P.D., Hindman, J.L. (2007). What is the relationship between teacher quality and student achievement? An exploratory study. *Journal of Personnel Evaluation in Education*, 20(3–4), 165–184.
- Sullivan, W., Marr, J., Hu, G. (2017). A predictive model for standardized test performance in michigan schools. In *Applied computing and information technology* (pp. 31–46): Springer.
- Superby, J.F., Vandamme, J.P., Meskens, N. (2006). Determination of factors influencing the achievement of the first-year university students using data mining methods. In *Workshop on educational data mining*, (Vol. 32 p. 234).
- Sweeney, M., Rangwala, H., Lester, J., Johri, A. (2016). Next-term student performance prediction: a recommender systems approach. *Journal of Educational Data Mining*, 8(1), 22–51.
- Tair, M.M.A., & El-Halees, A.M. (2012). Mining educational data to improve students' performance: a case study. *International Journal of Information*, 2(2), 140–146.
- Thai-Nghe, N., Busche, A., Schmidt-Thieme, L. (2009). Improving academic performance prediction by dealing with class imbalance. In *Ninth international conference on intelligent systems design and applications* (pp. 878–883): IEEE.
- Uddin, M.F., & Lee, J. (2017). Proposing stochastic probability-based math model and algorithms utilizing social networking and academic data for good fit students prediction. *Social Network Analysis and Mining*, 7(1), 29.
- Üstünlüoğlu, E. (2016). Teaching quality matters in higher education: a case study from Turkey and Slovakia. *Teachers and Teaching*, 23(3), 367–382.
- Uttl, B., White, C.A., Gonzalez, D.W. (2017). Meta-analysis of faculty's teaching effectiveness: Student evaluation of teaching ratings and student learning are not related. *Studies in Educational Evaluation*, 54, 22–42.
- Van Inwegen, E., Adjei, S., Wang, Y., Heffernan, N. (2015). An analysis of the impact of action order on future performance: the fine-grain action model. In *Proceedings of the fifth international conference on learning analytics and knowledge* (pp. 320–324): ACM.

- VeeraManickam, M.R.M., Mohanapriya, M., Pandey, B.K., Akhade, S., Kale, S.A., Patil, R., et al. (2018). Map-reduce framework based cluster architecture for academic student's performance prediction using cumulative dragonfly based neural network. *Cluster Computing*, 22(1), 1259–1275.
- Walters, W.H. (2007). Google scholar coverage of a multidisciplinary field. *Information Processing and Management*, 43(4), 1121–1132.
- Wang, Y., Ostrow, K., Adjei, S., Heffernan, N. (2016). The opportunity count model: a flexible approach to modeling student performance. In *Proceedings of the Third (2016) ACM Conference on Learning@Scale* (pp. 113–116): ACM.
- Widyahastuti, F., & Tjhin, V.U. (2017). Predicting students performance in final examination using linear regression and multilayer perceptron. In *10th international conference on human system interactions (HSI)* (pp. 188–192): IEEE.
- Willoughby, T., Anderson, S.A., Wood, E., Mueller, J., Ross, C. (2009). Fast searching for information on the Internet to use in a learning context: the impact of domain knowledge. *Computers & Education*, 52(3), 640–648.
- Wook, M., Yusof, Z.M., Nazri, M.Z.A. (2016). Educational data mining acceptance among undergraduate students. *Education and Information Technologies*, 22(3), 1195–1216.
- Xing, W., Guo, R., Petakovic, E., Goggins, S. (2015). Participation-based student final performance prediction model through interpretable Genetic Programming: Integrating learning analytics, educational data mining and theory. *Computers in Human Behavior*, 47, 168–181.
- Xiong, X., Adjei, S., Heffernan, N. (2014). Improving retention performance prediction with prerequisite skill features. In *Educational data mining 2014*.
- Xu, M., Liang, Y., Wu, W. (2017). Predicting honors student performance using RBFNN and PCA method. In *International Conference on database systems for advanced applications* (pp. 364–375): Springer.
- Yang, S.J.H., Lu, O.H.T., Huang, A.Y.Q., Huang, J.C.H., Ogata, H., Lin, A.J.Q. (2018). Predicting students' academic performance using multiple linear regression and principal component analysis. *Journal of Information Processing*, 26, 170–176.
- Yim, S., & Warschauer, M. (2017). Web-based collaborative writing in L2 contexts: methodological insights from text mining. *Language Learning & Technology*, 21(1), 146–165.
- Yin, H., Wang, W., Han, J. (2016). Chinese undergraduates' perceptions of teaching quality and the effects on approaches to studying and course satisfaction. *Higher Education*, 71(1), 39–57.
- Yoo, J., & Kim, J. (2014). Can online discussion participation predict group project performance? Investigating the roles of linguistic features and participation patterns. *International Journal of Artificial Intelligence in Education*, 24(1), 8–32.
- Yu, L., Lee, C., Pan, H., Chou, C., Chao, P., Chen, Z., et al. (2018). Improving early prediction of academic failure using sentiment analysis on self evaluated comments. *Journal of Computer Assisted Learning*, 34(4), 358–365.
- Zabaleta, F. (2007). The use and misuse of student evaluations of teaching. *Teaching in Higher Education*, 12(1), 55–76.
- Zacharis, N.Z. (2015). A multivariate approach to predicting student outcomes in web-enabled blended learning courses. *The Internet and Higher Education*, 27, 44–53.
- Zaugg, H., West, R.E., Tateishi, I., Randall, D.L. (2011). Mendeley: Creating communities of scholarly inquiry through research collaboration. *TechTrends*, 55(1), 32–36.
- Zimmermann, J., Brodersen, K.H., Heinimann, H.R., Buhmann, J.M. (2015). A model-based approach to predicting graduate-level performance using indicators of undergraduate-level performance. *Journal of Educational Data Mining*, 7(3), 151–176.
- Zhang, X., Sun, G., Pan, Y., Sun, H., He, Y., Tan, J. (2018). Students performance modeling based on behavior pattern. *Journal of Ambient Intelligence and Humanized Computing*, 9(5), 1659–1670.
- Zollanvari, A., Kizilirmak, R.C., Kho, Y.H., Hernández-Torrano, D. (2017). Predicting students' GPA and developing intervention strategies based on self-regulatory learning behaviors. *IEEE Access*, 5, 23792–23802.
- Zorilla, M.E., García-Saiz, D., Balcázar, J.L. (2010). Towards parameter-free data mining: mining educational data with yacaree. In *Educational data mining 2011*.
- Zuber, M. (2014). A survey of data mining techniques for social network analysis. *International Journal of Research in Computer Engineering & Electronics*, 3(6), 1–8.