

OVERVIEW

Big data analytics—A review of data-mining models for small and medium enterprises in the transportation sector

Siti Aishah Mohd Selamat¹  | Simant Prakoonwit¹ | Reza Sahandi² | Wajid Khan¹ |
Manoharan Ramachandran²

¹Department of Creative Technology, Faculty of Science and Technology, Bournemouth University, Poole, UK

²Department of Computing, Faculty of Science and Technology, Bournemouth University, Poole, UK

Correspondence

Siti Aishah Mohd Selamat, Department of Creative Technology, Faculty of Science and Technology, Bournemouth University, Poole, England.
Email: aishah@bournemouth.ac.uk

The need for small and medium enterprises (SMEs) to adopt data analytics has reached a critical point, given the surge of data implied by the advancement of technology. Despite data mining (DM) being widely used in the transportation sector, it is staggering to note that there are minimal research case studies being done on the application of DM by SMEs, specifically in the transportation sector. From the extensive review conducted, the three most common DM models used by large enterprises in the transportation sector are identified, namely “Knowledge Discovery in Database,” “Sample, Explore, Modify, Model and Assess” (SEMMA), and “CRoss Industry Standard Process for Data Mining” (CRISP-DM). The same finding was revealed in the SMEs’ context across the various industries. It was also uncovered that among the three models, CRISP-DM had been widely applied commercially. However, despite CRISP-DM being the de facto DM model in practice, a study carried out to assess the strengths and weakness of the models reveals that they have several limitations with respect to SMEs. This paper concludes that there is a critical need for a novel model to be developed in order to cater to the SMEs’ prerequisite, especially so in the transportation sector context.

This article is categorized under:

Application Areas > Business and Industry

Application Areas > Industry Specific Applications

KEYWORDS

data mining, knowledge discovery database, CRISP-DM, SEMMA, SMEs, transportation, big data

1 | INTRODUCTION

Intelligent transportation systems (ITS) utilize advanced technologies and systems to provide efficient and safe transportation services, while minimizing the operational cost and environmental impacts (An, Lee, & Shin, 2011). The ITS evolution has seen a dramatic development in the last two decades—whereby, from the 1970s to 1980s the primary area of development was concentrated in curbing traffic congestion (Taylor, Stott, Parker, & Bradley, 2015). From the 1980s to 1990s, the building of Intelligent Infrastructure and Vehicles was the core focus of development (Taylor et al., 2015). With the advancement of technology in the 21st century, data are increasingly collected every hour, every minute, and every second causing a data explosion era. The International Data Corporation (IDC) forecast that the volume of data is expected to grow up to 50 zettabytes globally (equivalent to 50 billion terabytes) by the year 2020 (Zeng & Lusch, 2013). These can revolutionize the development of ITS, by shaping a traditional technology-driven system (Zhang et al., 2011) into a more robust ITS ecosystem (Zeng & Lusch, 2013). The influx of data can only become an asset to the organization if they are implicitly intelligible to

translate useful knowledge for small and medium enterprise (SME) organizations (Derrick, 2012; Iansiti & Lakhani, 2014). As shared by Lyons in a thought-provoking editorial piece in *Transport Reviews*—in a progressive environment load with robustness, interconnectedness, it is yet “uncertain.” Therefore, there is a crucial need to evaluate the relevancy of the transport analysis purposes (Lyons, 2016).

In a conclusive report by European Commission (2017), the key economic drivers of growth in the European continent is the SMEs—contributed 3.9 trillion euros to the economy in 2015. This is twice as much in comparison to the large enterprises (Eurostat, 2009). The transportation and the storage enterprise made up of 5% of the 22.3 million of the nonfinancial business economy in 2012 (Eurostat, 2017). SMEs can further reap two to three times growth rate through the exploitation of advanced technologies (such as social media, big data, cloud computing, and mobile). Eurostat identified that less than 7% of the European SMEs have employed data analytics in their business, making a need for digital transformation as a high priority for the European Union (EU, 2013a) in this data explosion era. In an in-depth report by IDC on IDC European Vertical Market Survey 2012, it was ascertained that only 3% (an estimation of 1,500 out of 50,000 organizations) of the SMEs in the transportation and storage sector have deployed data analytics in their business (European Commission, 2013b). The transport and storage SMEs have been relentlessly labeled as laggards in the adoption of big data technologies. SMEs can become 5–6% more productive through the utilization of data analytics in the business—as evident in the larger transport companies (European Commission, 2013a). Despite the momentous potential benefits of utilization of big data analytics, the transport and storage SMEs are still dawdling in their adoption efforts. In 2012, the adoption rate of big data analytics of SMEs in the United Kingdom stood at 0.2% compared to the large enterprise, with an uptake of 25% (European Commission, 2017). This is indeed an alarming figure, as the fast adoption rate by the large enterprise may eventually implicate SMEs to become irrelevant and absolute. Therefore, there is an urgent need for SMEs to begin exploring the implementation of big data analytic and data mining (DM).

To ensure the relevance of this study, only articles published in the last 10 years were included. The selections of the literature were divided into three categories. The first category includes papers related to big data analytics for SMEs. The second and third category encompass papers related to DM models in the transportation sector from the SMEs and the large enterprises subsequently. The full text of each article was screened in order to validate its relevancy and applicability. Upon screening, only suitable articles were included in this study. This paper presents the outcome of a critical study of the big data analytics literature, in respect of DM models for SMEs in the transportation sector in particular. This information was extracted from online databases, such as ACM Digital Library, Science Direct, Springer, EBSCOhost EJS, Semantic Scholar, Google Scholar (search engine), and IGI Global. The research aims to provide researchers, transportation business leaders, and policy makers’ eminent findings of the big data analytics research studies. It is anticipated that this paper will magnify the emergence of big data technologies aiding SME’s understanding and capitalization to facilitate and spur business growth.

In section 2, the paper provides background information on the big data analytics challenges and problems faced by SMEs. Sections 3 and 4 cover the DM application case studies in both the transportation sector and SMEs’ context. Section 5 presents a comparative study completed for the most commonly used DM models. This is followed by a discussion on the model’s strengths and deficiencies. Finally, the paper concludes with several key points and a discourse of the future research work to be undertaken.

2 | BACKGROUND: SMEs’ ADOPTION BARRIERS IN BIG DATA ANALYTICS

Aside from the transport and storage sector, the SME groups are straggling with the implementation of big data analytics in their businesses. This raises an alarm as to the hindering factor(s) that is curbing the SMEs from advancing with the evolution of big data technology. In a recent in-depth study by Coleman et al. (2016), they uncover several core factors contributing to the slow acceptance of big data analytics by SMEs on the European continent. The factors are listed below:

1. Minimal cognizance in the big data analytics domain.
2. Little or no interest in new management trends.
3. Insufficient in-house data analytics experts.
4. Increasing shortage of competent data analyst in the labor market.
5. Lack of exemplary successful case studies for SMEs to refer to.
6. Lack of effective analytics consulting services.
7. Highly complex analytics solutions in the software market.
8. Data security concerns.
9. Data protection and privacy concerns.
10. Lack of financial access to invest in new technologies.

TABLE 1 Classification of SME big data analytic barriers

Areas of concerns	List of challenges and barriers
Resources	<ul style="list-style-type: none"> • Insufficient in-house data analytics experts • Increasing shortage of competent data analyst in the labor market • Lack of effective and value for money analytics consulting services available • Highly complex analytics solution in the software market • Lack of financial access to invest in new technologies
Knowledge	<ul style="list-style-type: none"> • Low understanding of big data analytics domain • Little or no interest in new management trends • Lack of exemplary successful case studies for the SMEs to refer to
Data management	<ul style="list-style-type: none"> • Data security concerns • Data protection and privacy concerns

SME = small and medium enterprise.

2.1 | Classifications of area of concerns

To further comprehend the nature of the identified barriers, the SMEs' areas of concerns are classified into three groups—resources, knowledge, and data management (as outlined in Table 1). Building on the research findings of Coleman et al. (2016), the three classified groups will constitute several examples from other research studies to support and validate these key barriers as identified.

2.1.1 | Resources

In respect of grouping classification, the SMEs' nucleus area of concern is mostly in relation to the subject of resources. This is followed by knowledge and data management concerns. It can be deduced that the lack of an in-house data analytics specialist is an implication caused by the insufficient number of available qualified data analytics talent. For instance, in the United States, it is envisaged that by 2018, there will be a shortage of close to 190,000 skilled analytical talents and also a shortfall of 1.5 million analysts and managers with the relevant competency to derive strategic decision(s) from big data analysis (Manyika et al., 2011). A survey carried out among recruiters in the United Kingdom revealed that up to 57% of the recruiters are facing obscurity in filling up the big data analytics gaps—this is inclusive of the large companies (UK e-skills, 2013). The scarcity of qualified data scientists would deter the analytics development scene in the European market (Probst et al., 2014). Given the shortfall of talent supply, it is expected that the existing analytics software, readily available in the market would aid in curbing the impending expertise gap. There are plenty of analytics solutions available in the market. Nonetheless, to find a solution that is both user friendly and embedded with robust analytical capabilities is scarce. The need for an instinctive user interface is critical in shortening the user learning curve (Probst et al., 2014)—allowing faster implementation for the SMEs. As evaluation platforms may tend to be vendor biased; an end user with minimal or zero proficiency in analytics may find it difficult to select a solution with a decent price–performance ratio. With an innumerable number of studies highlighting that financial limitations are the SMEs' major hindrance block (Bartlett & Bukvič, 2001; Fuller-Love, 2006), a lower price–performance ratio solution would be more desirable by the SMEs.

2.1.2 | Knowledge

In the area of knowledge concerns, a survey conducted in the United Kingdom reveals that SMEs' personnel has an exceptionally low comprehension of the big data analytics domain (UK e-skills, 2013). A similar survey conducted in Germany shared an identical result (Coleman et al., 2016). To a great extent, the SMEs are uncertain of their datasets potentiality, in turn drawing hesitation on the need to invest in data science capabilities to reap the intended benefits as affirmed by the various analytics *connoisseur*. In spite of the fact that there are guidelines available for the SMEs to make reference to, there is still a lack of exemplary research case studies that propagate the successful implementation of analytics in the SME sphere (Coleman et al., 2016). The existing big data use cases generated in the EU often do not correlate with the SMEs' points of interests (Wadhwa, 2014). More case studies are needed to possibly fill the knowledge gap and jumpstart the SMEs' enthusiasm to take more interest in the big data analytics domain.

2.1.3 | Data management

Finally, yet importantly, data security, protection, and privacy are the SMEs' key concerns in the area of data management. Close to 50% of SMEs identified data security and protection as the key barrier to big data analytics—in a worldwide survey of more than 82 SME companies (Coleman et al., 2016). In comparison to larger companies, SMEs lack the competency to scale up their IT security level (Lacey & James, 2010). The use of obsolete and nonupdated database management system raised a critical IT security gap for SMEs, consequently making SMEs less resilient against cyber-attacks and intrusions. The processing and analysis of customer's data by SMEs, on the other hand, has to abide by EU legality on data protection and

privacy. The lengthy EU data protection law (Rights EUAff, 2017), mooted in 2012, creates an added constraint for SMEs. Predominantly, in view of the lack of financial resources, SMEs could not meet the expenses to engage a legal expert in order to fully grasp the EU's data protection legislation requirements.

2.2 | Supplemental intangible barriers

In addition to the discussions in respect of adoption barriers outlined in sections 2 and 2.1, it is worth considering the supplemental intangible barriers that may also hinder the adoption of analytics for SMEs. These intangible barriers relate to SME's organizational culture, organization structure, and decision making. First and foremost, in terms of organizational culture—given that SMEs are highly domain specialized, they have little or no interest in new management trends that might be beneficial for the organization (Goebel, Norman, & Karanasios, 2015). This culture of intrinsic conservatism is leading the SME's attitude in taking big data analytics as a management hype instead of an opportunistic viewpoint. The second aspect of organizational structure implies the need to have a fitting management concept within an organization, in order to create an economic success on the adoption of analytics (McAfee & Brynjolfsson, 2012). Unlike the large enterprise, the organizational structures of most SMEs are flat with few or no levels of middle management between the executives and staffs (Capgemini, 2012). The flat organizational structure of SMEs would in turn impact on the way the organization makes its decisions. The decision makers in SMEs are often the business owners, which are in a way usually tied up with the owner's identity and life (Goebel et al., 2015). The decision making within the large enterprise tends to be more rational because of the complexity of the organization structure and decision-making units (Culkin & Smith, 2000). In view of the scarcity of resources and expertise, SMEs would be in a limiting position to make a complex decision, as it is often reliant on the business owner's intuition (Culkin & Smith, 2000). In other words, if the business owner is not personally attuned with the latest business trends of analytic adoption, it will be an intricate barrier to overcome.

2.3 | What DM can mean for SMEs?

The explosion of data is deemed critical for SMEs because, during the DM process, organizations can radically learn more about their business and translate the new knowledge into better decision making and performance (McAfee & Brynjolfsson, 2012). In other words, DM has the potential to transform traditional SMEs to organizations with a competitive advantage. For instance, suppose an SME aims to mine its customer data, the potential benefit would entail creating cross-selling avenues at a higher margin, improving its customer retention and satisfaction rates, identifying the most profitable customer group and last but not least, enhancing the SME's marketing and sales strategy (Tan, Yeoh, Boo, & Liew, 2013). In the mining of inventory data, on the other hand, SMEs can gain an advantage by forecasting the inventory that will help reduce the total value of stock held. This would create a positive implication in allowing timely inventory purchase from the supplier, creating a supplier lock-in, leading the SME to a better trading agreement (Waller & Fawcett, 2013). It is therefore evident that DM has the capability to create business opportunities to enable SMEs to stay ahead of their competition and leverage on the possibilities.

3 | DM IN TRANSPORTATION SECTOR

In the area of transportation, there have been several pieces of research developing novel approaches to traffic management, motorist and commuter safety, transport mobility, road accident management, and much more—with the application of DM. Table 2 is a compilation of applications of DM in the transportation sector.

This table illustrates that DM is widely applied in all three transportation modes—land, air, and sea. An example of a DM application in land transportation is research carried out by Giovanni et al., in which DM is used to predict the railroad demands to facilitate operational and manpower planning for Malha Regional Sudeste (MRS) Logistica (Viglioni, Cury, & Silva, 2007). Cristobal et al.'s research denotes a similar area of interest in the effective management of resource planning in predicting passenger demands for Gran Canaria Island Public Transport (Moreno-Díaz, Pichler, & Quesada-Arencibia, 2015). An example of sea transportation DM application is Greis et al.'s research, which the study involves in applying DM to identify high-risk shipments reaching the U.S. ports (Greis & Nogueira, 2011). Finally, Lukáčová, Babič, and Paralič's (2014) research adopt DM to assist the Federal Aviation Administration (Abdul Rahman, Shamsuddin, Hassan, & Abu Haris, 2016) to predict potential incidents and implications.

From the compilation of DM application in the transportation sector, Table 2 reflects a distinctive commonality, whereby DM was applied to extract new information/knowledge for prediction capabilities. Second, the three recurring DM models adopted by the enterprises were Knowledge Discovery in Database (KDD), Sample, Explore, Modify,

TABLE 2 Case studies on DM applications in the transportation sector

Studies	Sector type	Company	Transportation mode/type	Data type	Model adopted	Application
Viglioni and Cury (2007)	Private	MRS Logistica	Land/train	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	Prediction of railroad demands to facilitate operation and manpower planning
Wong and Chung (2007)	Private	Taiwanese Domestic Airline	Air/airflight	<ul style="list-style-type: none"> Historical data Structured data 	KDD	Mining passengers' demographic, travel behavior and core service quality information for customer retention initiatives
Haluzová (2008)	Public	Prague Public Transit Company	Land/bus	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	Identification of the accident influences between car and tram on the electric tramway net
Shin, Park, Saha, and Kim (2009)	Private	Jeju Taxi Service	Land/taxi	<ul style="list-style-type: none"> Historical data Unstructured data 	SEMMA	Analyzing passenger pick-up location patterns to proposed potential pick-up locations for empty taxis
Mirabadi and Sharifian (2010)	Public	Iranian Railways	Land/train	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	Analyzing historical accident data to discover the unsafe condition contributing factors
Zhang, Huang, and Zong (2010)	Public	China railways	Land/train	<ul style="list-style-type: none"> Historical data Structured data 	KDD	Deriving intelligent-based decision making in accident treatments
Greis and Nogueira (2011)	Public	US Seaport (Department of Homeland Security)	Sea/shipping cargo	<ul style="list-style-type: none"> Real-time data Structured data Unstructured data 	CRISP-DM	Identification of high-risk shipments reaching the US ports
de Almeida and Ferreira (2013)	Public	BUS public transport	Land/bus	<ul style="list-style-type: none"> Historical data Structured data 	SEMMA	Identification of the most fuel-efficient resources in route operation and areas of resources for improvements
Lukáčová et al. (2014)	Public	Federal aviation administration	Air/airflight	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	Analyzing the aviation historical incident data to predict potential incidents and implications
Moreno-Díaz et al. (2015)	Public	Gran Canaria Island Public Transport	Land/bus	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	Predicting passenger demand for efficient resource planning and deployment

CRISP-DM = CRoss Industry Standard Process for Data Mining; DM = Data mining; KDD = Knowledge Discovery in Database; SEMMA = Sample, Explore, Modify, Model, and Assess.

Model, and Assess (SEMMA), and CRoss Industry Standard Process for Data Mining (CRISP-DM). Out of the 10 industrial examples quoted, 6 enterprises had adopted the CRISP-DM model and the remaining 4 enterprises had used the SEMMA and KDD model equally. Of the table, it is evident that CRISP-DM marks as the most commonly used model. And these findings correlate with the industrial polls conducted by KdNuggets.Com (Kdnuggets, 2002, 2004, 2007). On the types of data used, the majority of the enterprises are leveraging on their historical data for data processing and analyzing. The data are in the form of structured data—referring to data that are organized in a relational database that is structured in columns (fields) and rows (record) (Duan & Xiong, 2015). On a different note, one key prominent finding derived from Table 2 indicates that little is known of research studies on SMEs in the transportation sector.

4 | DM IN THE SMEs' CONTEXT

Having understood that there is a minimal study on SMEs in the transportation sector, this section aims to encapsulate the use of DM in the SMEs' context across various industries. In all, 10 recent case studies have been tabulated, as seen in Table 3. The case studies reflect the application of DM in the food and beverage, tourism, information technology (IT), financial, aviation, trading, and manufacturing industries. From the 10 case studies, DM is primarily used for prediction and improvement of the decision-making process. For instance, under the finance industry, Mandala, Nawangpalupi, and Praktikto's (2012) research is involved in assessing the credit risk of loan lenders. On the other hand, Koyuncugil and Ozgulbas (2012), in their research, employ DM to predict to detect financial risks for SMEs. In respect of IT, SMEs are also utilizing DM for various usages. In a research conducted by Bozdogan and Zincir-Heywood (2012), DM is used to facilitate the SMEs in collecting its public resources automatically to create a knowledge base to support its IT management. In another study conducted by Ibukun et al., the objective for applying DM is to identify customer segmentation in order to carry out target marketing (Afolabi et al., 2016). A recent study by Packianather et al. (2017) for SMEs in the manufacturing industry, they applied DM to generate unique and new knowledge for forecasting and strategic decision making. The case studies for application of DM by SMEs as depicted in Table 3 insinuate that DM can be applied across many industries. The CRISP-DM is a widely used model. Our in-

TABLE 3 Case studies on DM applications in the SME context

Studies	Research context	Industry	Data type	Model adopted	Application
Raju et al. (2016)	UK-based SME wholesaler	Food and beverage	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	Forecasting freshly produced (short shelf life) product demands
Rebón, Castander, Argandoña, Gerrikagoitia, and Alzua-Sorzabal (2015)	Tourism SMEs	Tourism	<ul style="list-style-type: none"> Real-time data Structured data 	KDD	Enhancing the analysis technique to improve decision-making process of credit fraud transaction detection
Pytel et al. (2013)	Project planning for SME	Information technology	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	To predict the cost and effort estimation for small-sized software projects
Mandala et al. (2012)	Rural bank	Financial	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	To develop credit assessment to in order to classify lenders as performing or nonperforming loan risk
Koyuncugil and Ozgulbas (2012)	SMEs	Financial	<ul style="list-style-type: none"> Historical data Structured data Unstructured data 	KDD	To predict financial risk detection for SMEs
Bozdogan and Zincir-Heywood (2012)	IT management for SMEs	Information technology	<ul style="list-style-type: none"> Historical data Structured data 	SEMMA	To automatically generate an IT management support knowledge base from public resources
Cheung and Li (2012)	SMEs	Trading	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	To uncover hidden patterns in the sales and market domain
Packianather et al. (2017)	SMEs	Manufacturing	<ul style="list-style-type: none"> Historical data Structured data 	KDD	To generate unique and new knowledge for forecasting and strategic decision making
Young et al. (2010)	SMEs	Aviation	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	To decide how and where aircraft maintenance process can be enhanced or amended
Afolabi, Worlu, and Uwadia (2016)	SMEs	Information technology	<ul style="list-style-type: none"> Historical data Structured data 	CRISP-DM	To identify customer segmentation in order to carry out target marketing

CRISP-DM = CRoss Industry Standard Process for Data Mining; DM = data mining; KDD = Knowledge Discovery in Database; SEMMA = Sample, Explore, Modify, Model, and Assess; SME = small and medium enterprise.

depth review (see Tables 2 and 3) exhibits that CRISP-DM is the most frequently used DM model. Out of the 10 case studies considered, 6 studies utilized the CRISP-DM model, 3 used the KDD model, and 1 adopted the SEMMA model. Another synonymous result reflects that the SMEs are in most cases using their historical structured data for DM applications—as reflected in Table 2. Details of the elements and functions of the most frequently used DM models will be discussed in the following section.

5 | THE EVOLUTION OF DM

DM refers to the science of identifying valuable and unique information from a substantial size of datasets or databases. The mining procedures involve intensive data analysis (Xindong, Xingquan, Gong-Qing, & Wei, 2014). The core goal of DM is processing a large amount of data to generate new knowledge (Fan & Bifet, 2013). DM serves two primary goals—(a) uncover new insights and (b) generate predictions (Fayyad, Piatetsky-Shapiro, & Smyth, 1996). DM is a process for KDD (Fayyad et al., 1996). The term KDD refers to a set of broad processes used to discover valuable knowledge from a set of data collection (Piatetsky-Shapiro, 1991). The emergence of KDD was sparked by the rising establishment of big databases in the varying number of organizations in the early 1990s (Mariscal, Marbán, & Fernández, 2010). This created a paradigm shift for the need to develop DM algorithms with the capabilities to unearth gainful insights from the big volume of data that are residing in companies' databases. Figure 1 depicts the overall evolution of the DM process models with KDD as its foundation and CRISP-DM as the core focal point of the evolution. This section only discusses the three most applied models, KDD, SEMMA, and CRISP-DM, in depth. A critical comparison of these models has been made which is discussed in the following section.

5.1 | KDD

The KDD process can be defined as an un-superficial way of distinguishing potentially useful, valid and conclusively understandable patterns from the data (Marbán, Mariscal, & Segovia, 2009). The term process refers to the many stages that are involved in the KDD process. KDD can also be described as the overall approach of uncovering valuable knowledge from data (Fayyad et al., 1996). It also entails the evaluation and (perhaps) the interpretation of the new insights and knowledge

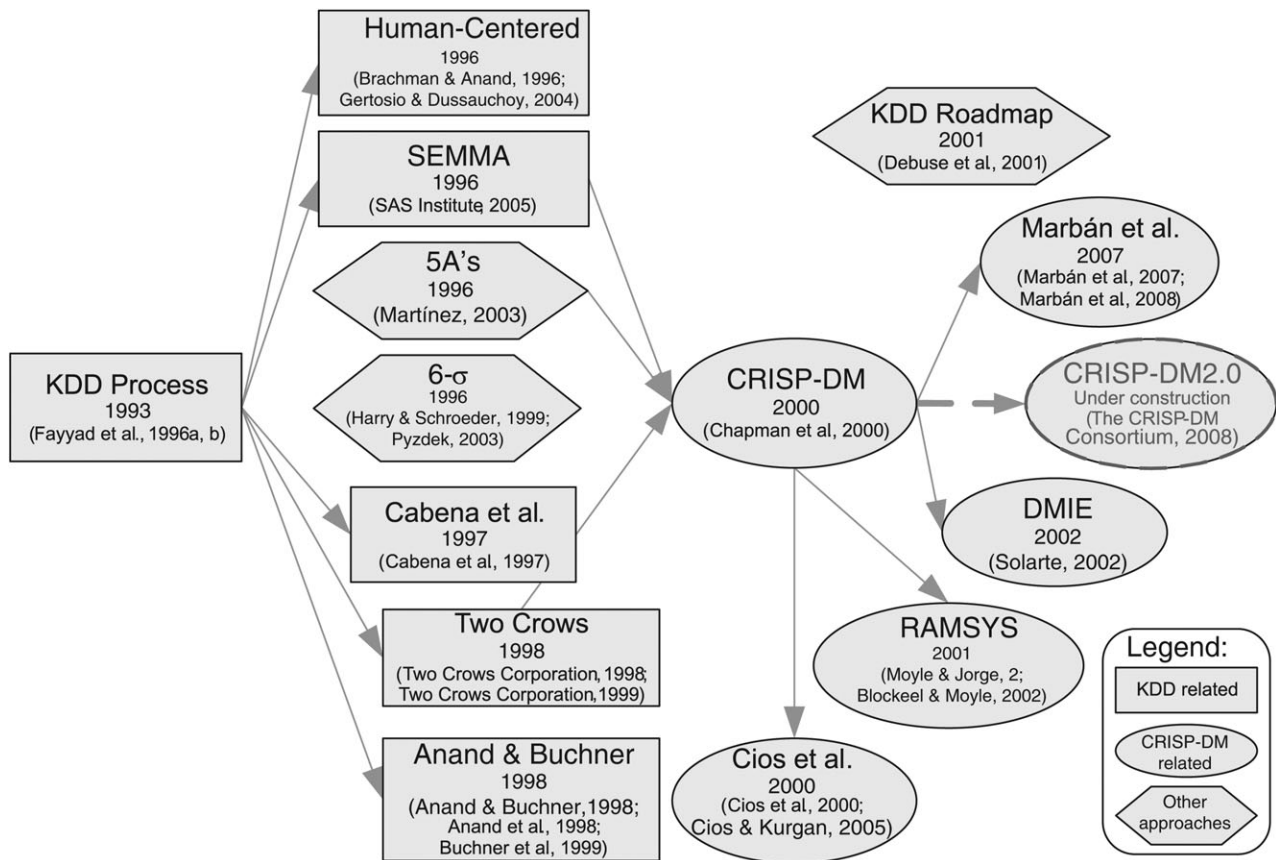


FIGURE 1 Evolution of data-mining methodologies. The figure represents the evolution of data-mining methodologies from the 1990s to 2000s. Source: Mariscal et al. (2010)

for decision making. Outlined in Figure 2 is the overall overview of the KDD process from the data viewpoint—interactive, iterative, and with many feedback loop points. In all, the KDD process encompasses nine steps (Fayyad et al., 1996).

The outline of the KDD steps is as follows:

1. Understanding the application domain.
2. Constructing a target dataset.
3. Data cleaning and pre-processing.
4. Data transformation.
5. DM function selection.
6. DM algorithm selection.
7. DM.

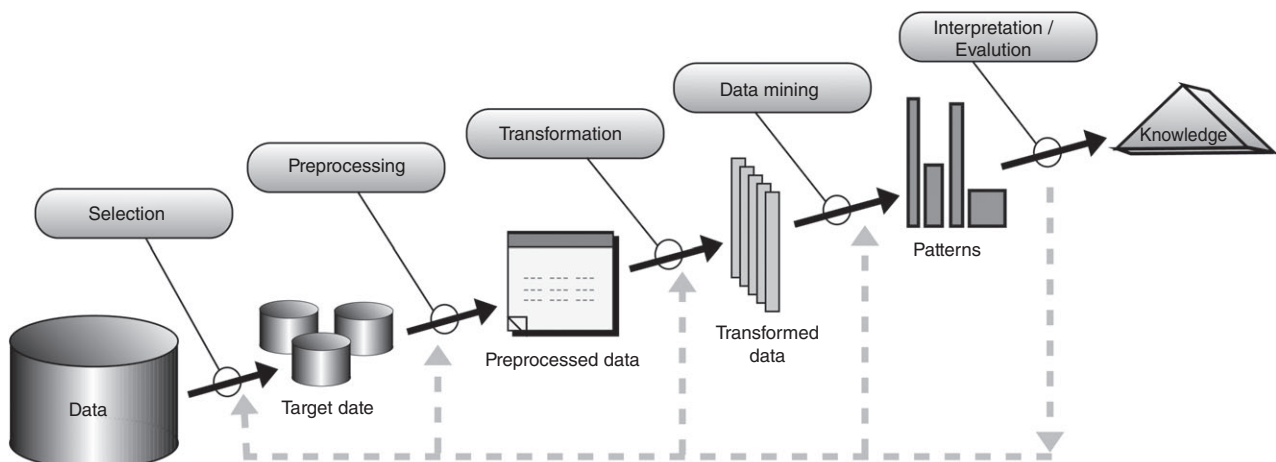


FIGURE 2 Knowledge Discovery in Database. The figure represents the overall KDD methodology process. Source: Fayyad et al. (1996)

8. Examination and evaluation of mined data.
9. Employing newly discovered knowledge.

5.2 | SEMMA

Developed by the Suite of Analytics (SAS) Institute, SEMMA refers to the systematic tool set of SAS enterprise miner for delivering the DM core tasks (Inc SI, 2005). The SEMMA model can only function with the enterprise miner tool, which had been developed by the SAS Institute. The KDD process, on the other hand, is an open-source process that can be administered in various environments. The SEMMA model's principle focus is on its model development point of DM (Shmueli, Patel, & Bruce, 2011). Figure 3 illustrates the five SEMMA steps. The steps consist of Sample, Explore, Modify, Model, and Assess.

5.3 | CRISP-DM

The CRISP-DM model was mooted together by highly acclaimed organizations such as Teradata, Daimler-Chrysler, SPSS, and OHRA in the mid-1990s (Chapman et al., 2000). It is considered as a de facto standard for establishing DM projects. The popularity of CRISP-DM was contributed to the fact that the model is applicable across all industries (Mariscal et al., 2010). Unlike the KDD and SEMMA models, the CRISP-DM process model renders a continuous lifecycle *modus operandi*. In addition, in each phase of the project, it corresponds to the designated tasks and interrelation between each task. As depicted in Figure 4, the overall cycle of the CRISP-DM DM project comprises of six stages. The chain of cycle in CRISP-



FIGURE 3 Sample, Explore, Modify, Model, Assess methodology. The figure represents the overall SEMMA methodology process. Source: Mariscal et al. (2010)

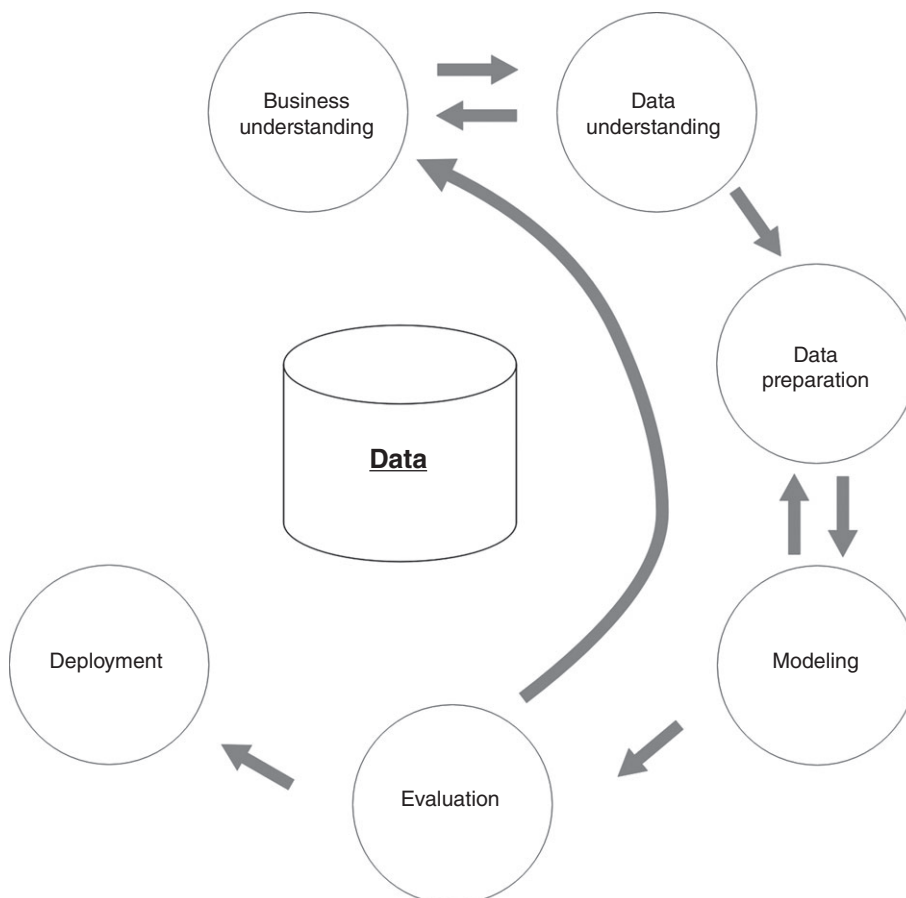


FIGURE 4 CRISP-DM methodology. The figure represents the overall CRISP-DM methodology process. Adapted from Chapman et al. (2000)

DM is flexible, allowing the end user to move back and forth freely. The chain of sequence is really dependent on the result of the specific task of the concerning phase.

The six phases of CRISP-DM are

1. Business understanding
2. Data understanding
3. Data preparation
4. Modeling
5. Evaluation
6. Deployment

6 | DETAILED DISCUSSIONS OF DM MODELS

In this section, the DM models will be discussed in two parts. The first part consists of a critical comparison of the KDD, CRISP-DM, and SEMMA models from various aspects. The strengths and limitations of the DM model in the SMEs' context are discussed. The objective of this discussion is to synthesize the findings, which would be key in the proposal for a new DM model, especially for SMEs in the transportation sector, in which further elaboration are provided.

6.1 | Comparison of models

Quantitative and qualitative comparisons of the three models are shown in Table 4. For illustration purposes, the first and second rows outline the key facts of the models, namely the creator(s) and the year when each model was first introduced. The third to fifth rows indicate the models' functions, the total number of implementation steps, and a brief description of each phase. Subsequently, the sixth to the tenth rows specify the industry involvement of each model, the requirement for background knowledge in DM, the status of the software tool supported, availability of model documentation for users' reference, and finally, the status whether the model can support an open-source tool. The 11th to the final rows consist of industry-related components like total case studies in the transportation sector, total case studies in the context of SMEs, overall case studies across all sectors, application areas, and finally, the KDnuggets poll results for 2007 and 2014. The core motivation for carrying out the model comparisons is shown in Table 4. It is worth pointing out that despite the variation of phases entailed by each DM model, the three models are entrusted with the same core functions. The eventual intended outcome of the three DM models is to uncover new insights and to generate predictions. One prominent difference among the three models is the need for an initial understanding (in the domain or business) of the first phase of DM projects. Unlike SEMMA, both KDD and CRISP-DM require this particular phase. For SEMMA, the DM methodology focal point is in its technical characteristic that is involved during the development process, starting with data sampling. Another key difference that differentiates SEMMA from the other two models is its shortfall in applying and deploying the new knowledge, which is uncovered. In terms of the development of the models, only KDD had not had any involvement from industry and has no supporting documentation for the users to make reference to. A common component that all the three models share is the need to have prior knowledge in DM in order to be able to apply the model in practice. With reference to software support, unlike KDD and CRISP-DM, SEMMA does not support open-source tools. As indicated in Table 3, CRISP-DM is the most applied DM model by SMEs in the transport sector. This finding is also affirmed by a poll conducted by KDnuggets—a leading online resource on DM. The first poll conducted on 200 respondents in 2007 indicates that 42% of the respondents used the CRISP-DM model, 19% created their own model, 13% used SEMMA, 4% used other nondomain-specific models, 7.3% used the KDD process, and the remaining 5.3% used other models. A second poll conducted in 2014 on the same number of respondents shows that 43% of the respondents used the CRISP-DM model, 27.5% created their own model, 8.5% used SEMMA, 8% used other nondomain-specific models, 7.5% used the KDD process, and the remaining 5.5% used other models. The two polls observed an increase in the usage of the CRISP-DM model. It is worth noting that almost one-quarter of the respondents were using their own model in their own individual domain. This pre-eminent finding may suggest that CRISP-DM is not ultimately the de facto DM model for all domains. Nonetheless, as shown in Table 4, CRISP-DM outweighs the KDD and SEMMA models in view of it being industrially attuned. This is credited to its development on real-world knowledge discovery experience (Chapman et al., 2000).

6.2 | Strengths and limitations of models in the SME context

Following on from the model's comparison, the strengths and limitations of each model are discussed in this section. This is done in the context of SMEs in the various sectors. Table 5 shows strengths and limitations of each model according to the

TABLE 4 Model comparison—KDD, CRISP-DM, and SEMMA

MODEL	KDD	CRISP-DM	SEMMA
Developed by	Fayyad et al.	CRISP-DM Consortium	SAS Institute
Year of model introduced	1996	1996 (officially released in 2000)	1997
Functions	1. Uncover new and unique insights 2. Generate predictions		
Total steps	9	6	5
Phase	1. Application domain understanding	1. Business understanding	-
	2. Creating a target data set	2. Data understanding	1. Sample
	3. Data cleaning and		2. Explore
	4. Data transformation	3. Data preparation	3. Modify
	5. Data-mining method selection	4. Modeling	4. Model
	6. Data-mining algorithm selection		
	7. Data-mining application		
	8. Discovered patterns interpretation	5. Evaluation	5. Assessment
	9. Using discovered knowledge	6. Deployment	-
Industry involvement	No	Yes	Yes
		Consortium of companies involving Teradata, Daimler-Chrysler, SPSS and OHRA	Individually by SAS Institute
Requires background knowledge in DM	Yes	Yes	Yes
Software tool support	Yes	Yes	Yes
	Mineset	SPSS Clementine	SAS
Documentation	No	Yes	Yes
Open-source tool support	Yes	Yes	No
Total case studies in the transportation sector count	2	6	2
Total case studies in the SME context count	3	6	1
Overall case studies count	5	12	3
Application areas	Aviation, rail, tourism, financial, manufacturing	Logistic, cargo, aviation, rail, public transport, software, financial, marketing and sales, trading	Software, public transport, street taxis
Kdnuggets poll results for 2007 (200 votes total) Kdnuggets (2014)	7.3%	42%	13%
Kdnuggets poll results for 2014 (200 votes total) Kdnuggets (2014)	7.5%	43%	8.5%

case studies listed. Table 5 shows, the consistent list of strengths of CRISP-DM illustrates that the model is applicable to the industry, addressing business objectives and issues, providing structured approaches and processes, as well as having the flexibility to use any DM tool. This holistically indicates that CRISP-DM is practically applicable to the business environment. The compiled list of limitations, on the other hand, is as follows:

1. Long and arduous process with detailed steps to be undertaken in each process.
2. Requires DM knowledge.
3. Explicit need for detailed DM requirements.
4. Challenge faced in deriving how and when the selection of data is necessary or irrelevant.
5. The late selection of DM technique affects the data format compatibility causing to return to the data analysis stage.
6. Inadequate knowledge on domain expert's terminology.

Based on the CRISP-DM limitations, perhaps, the outcome of the case studies is suggesting that the user is experiencing an exhaustive process when applying the model in view of the detail steps that each phase contains. Further, with insufficient knowledge on DM, the user may face challenges in grasping the CRISP-DM concept and mode of application; for instance, in the study by Cheung and Li (2012), the limitation encountered during the identification of the necessary or irrelevant data for analysis. In the study by Ibukun et al., the issue faced was the inadequacy knowledge of DM's terminologies (Afolabi et al., 2016). The technical aspect as encountered by Young, Fehskens, Pujara, Burger, and Edwards (2010) was the delay in the selection of the DM technique. When the selected DM techniques do not correspond with the selected data format, incompatibility will occur, causing the user to return again to the data analysis stage. The strengths of the second most

TABLE 5 DM model strengths and limitations—SME context

Model	Case studies	SME industry	Strengths	Limitations
CRISP-DM	Raju et al. (2016)	Food and beverage	<ul style="list-style-type: none"> • Applicable to industry context • Addresses business objective and issues • Structured approach and process 	<ul style="list-style-type: none"> • Long and arduous process with detail steps to be undertaken in each process
	Pytel et al. (2013)	Information technology	<ul style="list-style-type: none"> • Applicable to industry context • Having the flexibility of using any DM tool 	<ul style="list-style-type: none"> • Requires DM knowledge • Explicit need for detailed DM requirements
	Mandala et al. (2012)	Financial	<ul style="list-style-type: none"> • Applicable to industry context • Addresses business objective and issues 	<ul style="list-style-type: none"> • Requires DM knowledge
	Cheung and Li (2012)	Trading	<ul style="list-style-type: none"> • Applicable to industry context • Structured approach and process • Having the flexibility of using any DM tool 	<ul style="list-style-type: none"> • Challenge faced in deriving how and when the selection of data is necessary or irrelevant • Long and arduous process with detail steps to be undertaken in each process
	Young et al. (2010)	Aviation	<ul style="list-style-type: none"> • Applicable to industry context • Having the flexibility of using any DM tool • Structured approach and process 	<ul style="list-style-type: none"> • Requires DM knowledge • The late selection of DM technique affects the data format compatibility causing to return to the data analysis stage
	Afolabi et al. (2016)	Information technology	<ul style="list-style-type: none"> • Applicable to industry context • Having the flexibility of using any DM tool 	<ul style="list-style-type: none"> • Inadequate knowledge on domain expert's terminology
KDD	Rebón et al. (2015)	Tourism	<ul style="list-style-type: none"> • Interactive and iterative process • Can be applied to industry context 	<ul style="list-style-type: none"> • Requires background knowledge in DM
	Koyuncugil and Ozgulbas (2012)	Financial	<ul style="list-style-type: none"> • User-friendly process • Contains feedback loops in each process 	<ul style="list-style-type: none"> • Requires background knowledge in DM • Insuitability of tool with the prepared data cause to make unnecessary loop back to the earlier process
	Packianather et al. (2017)	Manufacturing	<ul style="list-style-type: none"> • Interactive and iterative process • Can be applied to industry context 	<ul style="list-style-type: none"> • Wrong selection of data resulted to the wrong results
SEMMA	Bozdogan and Zincir-Heywood (2012)	Information technology	<ul style="list-style-type: none"> • Applicable to industry context • Robust support by DM provider 	<ul style="list-style-type: none"> • Highly technical centric process • No clear indication on how to apply the new knowledge

applied model, KDD, were accredited to its highly interactive process, containing feedback loops in each process. Like CRISP-DM, the KDD model is endorsed for being industrially applicable in the business environment. The limitations of KDD on the contrary share similar aspects to that of CRISP-DM whereby, the incompatibility of tool and datasets would require the KDD users to return to the KDD process again. Finally, the strengths of the least applied model, SEMMA, are attributed to the availability of its DM robust user support. Its drawback is in the case of the model's highly technical centric process. This may pose a great challenge if an organization has assigned a user who is not technically equipped in the SEMMA domain. In addition, in comparison to CRISP-DM and KDD, the SEMMA model falls short on the knowledge application phase.

7 | CONCLUSIONS AND FUTURE WORK

The need for SMEs to deploy data analytics has reached a point of criticality; with the immense surge of data collected via the advancement of technologies. In accordance to Eurostat, SMEs will reap a higher productivity level of up to 6% through the utilization of data analytics in the business. Ignoring the call for technology advancement can risk SMEs falling behind large enterprises that are more forthcoming toward the adoption of the new technology. As identified in section 2, resource, knowledge, and data management are the key areas of concerns that are hindering SMEs from adopting analytics—all of which needs to be addressed. In this paper, the applications of DM models are examined in the transportation sector in the context of both SMEs and large organizations. This paper reveals a compelling finding (see Tables 2 and 3) that CRISP-DM is the most prominent DM model that is widely used by SMEs in the transportation sector. The model is mainly used for prediction and facilitation of decision-making processes. Another commonality in the findings is the types of data being used to run the CRISP-DM model. The majority of the businesses are leveraging on their historical structured data. Finally, this paper highlights that there is limited case study research on DM application by SMEs in the transportation sector in particular. Three common DM processes (KDD, SEMMA, and CRISP-DM) were critically compared. The comparison was made against the overall implementation processes, the model's strengths and limitations. In all, the findings show the reason why CRISP-DM has been commercially adopted. In addition, our research shows that the model is flexible to suit any business using any DM tool.

Despite CRISP-DM being the de facto DM model for businesses to adopt—as examined in the study of the model's strengths and limitations in the SME context, there are several shortfalls that require addressing. The core limitation is the principal expectation of the need to have background knowledge on DM in order to fully grasp the terms, concept, and application of DM for the organization. The second limitation relates to the intense and exhaustive process that the CRISP-DM entails for applying the model in practice. Finally, the delay due to the selection of a DM technique may implicate on the data format compatibility affecting the DM overall process. Following up from this paper, the future research work aims to develop a novel DM model to suit SMEs in the transportation sector. Taking CRISP-DM as the foundation model, an Intelligent Transportation Analytical Model (ITAM) is to be developed. The ITAM aims to conduct an intelligent analysis with the objectives of churning out new insights, showing hidden patterns and relationship within the existing datasets to aid business decision making. This would undertake the impending limitations of the CRISP-DM model and at the same time taking, into consideration the impending SMEs' constraints as learned—primarily in terms of time and human-capacity constraints. Transportation SMEs sector will be identified. The companies' datasets will be collected for evaluation that is to understand the dataset characteristics. Following that, the ITAM will be proposed, tested in a real-life application and undergoes evaluations.

CONFLICT OF INTEREST

The authors have declared no conflicts of interest for this article.

ORCID

Siti Aishah Mohd Selamat  <http://orcid.org/0000-0003-2844-9806>

REFERENCES

- Abdul Rahman, F., Shamsuddin, S. M., Hassan, S., & Abu Haris, N. (2016). A Review of KDD-data mining framework and its application in logistics and transportation. *International Journal of Supply Chain Management*, 5, 77–84.
- Afolabi, I. T., Worlu, R. E., & Uwadia, O. C. (2016). Data mining approach for target marketing SMEs in Nigeria. *Covenant Journal of Informatics and Communication Technology*, 4, 1–5.
- de Almeida, J., & Ferreira, J. C. (2013). BUS public transportation system fuel efficiency patterns. In: *International Conference on Machine Learning and Computer Science*, Kuala Lumpur, Malaysia.
- An, S.-H., Lee, B.-H., & Shin, D.-R. (2011). A survey of intelligent transportation systems. In: *Computational Intelligence, Communication Systems and Networks (CICSyN), 2011 Third International Conference on Computational Intelligence, Communication Systems and Networks*, Bali, Indonesia. 332–337. <https://doi.org/10.1109/CICSyN.2011.76>
- Bartlett, W., & Bukvič, V. (2001). Barriers to SME growth in Slovenia. *MOST: Economic Policy in Transitional Economies*, 11, 177–195.
- Bozdogan, C., & Zincir-Heywood, N. (2012). Data mining for supporting it management. In: *Network Operations and Management Symposium (NOMS), 2012 IEEE: IEEE*.
- Capgemini. (2012). Measuring organizational maturity in predictive analytics: The first step to enabling the vision resource. https://www.capgemini.com/resource-fileaccess/resource/pdf/Measuring_Organizational_Maturity_in_Predictive_Analytics_the_First_Step_to_Enabling_the_Vision.pdf [accessed 1 October 2016].
- Chapman, P., Clinton, J., Kerber, R., Khabaza, T., Reinartz, T., Shearer, C., & Wirth, R. (2000). CRISP-DM 1.0 step-by-step data mining guide. Technical report, CRISP-DM.
- Cheung, C., & Li, F. (2012). A quantitative correlation coefficient mining method for business intelligence in small and medium enterprises of trading business. *Expert Systems with Applications*, 39, 6279–6291.
- Coleman, S., Göb, R., Manco, G., Pievatolo, A., Tort-Martorell, X., & Reis, M. S. (2016). How can SMEs benefit from big data? Challenges and a path forward. *Quality and Reliability Engineering International*, 32, 2151–2164. <https://doi.org/10.1002/qre.2008>
- European Commission. (2013a). Strategic Policy Forum on Digital Entrepreneurship Fuelling Digital Entrepreneurship in Europe Background paper. <https://ec.europa.eu/docsroom/documents/5313/attachments/1/translations/en/renditions/native> [accessed 22 October 2016].
- European Commission. (2013b). Business opportunities: Big data. https://ec.europa.eu/growth/tools-databases/dem/sites/default/files/pagefiles/big_data_v1.1.pdf [accessed 22 October 2016].
- European Commission. (2017). Annual report on European SMEs 2015/2016—SME recovery continues. https://ec.europa.eu/jrc/sites/jrcsh/files/annual_report_-_eu_smes_2015-16.pdf [accessed 22 October 2016].
- Culkin, N., & Smith, D. (2000). An emotional business: A guide to understanding the motivations of small business decision takers. *Qualitative Market Research: An International Journal*, 3, 145–157.
- Derrick, H. (2012). The C-level is coming around on big data [infographic]. Retrieved from <http://search.ebscohost.com/login.aspx?direct=true&db=edsnbk&AN=146FD1536214CEC8&site=eds-live&scope=site>
- Duan, L., & Xiong, Y. (2015). Big data analytics and business analytics. *Journal of Management Analytics*, 2, 1–21. <https://doi.org/10.1080/23270012.2015.1020891>
- E-skills UK. (2013). *Big data analytics: adoption and employment trends, 2012–2017*. Adelaide, Australia: VOCEDplus.
- Eurostat. (2009). SMEs were the main drivers of economic growth between 2004 and 2006. Retrieved from: <http://ec.europa.eu/eurostat/en/web/products-statistics-in-focus/-/KS-SF-09-071>
- Eurostat. (2017). Structural business statistics overview—Statistics explained. Retrieved from http://ec.europa.eu/eurostat/statistics-explained/index.php/Structural_business_statistics_overview-Main_statistical_findings
- Fan, W., & Bifet, A. (2013). Mining big data: Current status, and forecast to the future. *SIGKDD Explorations*, 14, 1.
- Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). From data mining to knowledge discovery in databases. *AI Magazine*, 17, 37.
- Fuller-Love, N. (2006). Management development in small firms. *International Journal of Management Reviews*, 8, 175–190.

- Goebel, R., Norman, A., & Karanasios, S. (2015). Exploring the value of business analytics solutions for SMEs. Association for Information Systems AIS Electronic Library (AISeL). *UK Academy for Information Systems Conference Proceedings 2015*. <http://aisel.aisnet.org/ukais2015/22>
- Greis, N. P., & Nogueira, M. L. (2011). Use of data mining for validation and verification of maritime cargo movement. *Institute for Homeland Security Solutions (Research Brief)*. North Carolina: University of North Carolina.
- Haluzová, P. (2008). Effective data mining for a transportation information system. *Acta Polytechnica*, 48, 24–26.
- Iansiti, M., & Lakhani, K. R. (2014). Digital ubiquity: How connections, sensors, and data are revolutionizing business. *Harvard Business Review*, 92, 90–99.
- Inc SI. (2005). SAS version 9.1. SAS. North Carolina, USA.
- Kdnuggets. (2002). What main methodology are you using for data mining? Retrieved from <http://www.kdnuggets.com/polls/2002/methodology.htm> [accessed 25 September 2016].
- Kdnuggets. (2004). Data mining methodology. Retrieved from: http://www.kdnuggets.com/polls/2004/data_mining_methodology.htm [accessed 25 September 2016].
- Kdnuggets. (2007). Data mining methodology. Retrieved from http://www.kdnuggets.com/polls/2007/data_mining_methodology.htm [accessed 25 September 2016].
- Kdnuggets. (2014). What main methodology are you using for your analytics, data mining, or data science projects? Poll. Retrieved from <http://www.kdnuggets.com/polls/2014/analytics-data-mining-data-science-methodology.html> [accessed 25 September 2016].
- Koyuncugil, A. S., & Ozgulbas, N. (2012). Financial early warning system model and data mining application for risk detection. *Expert Systems with Applications*, 39, 6238–6253. <https://doi.org/10.1016/j.eswa.2011.12.021>
- Lacey, D., & James, B. E. (2010). *Review of availability of advice on security for small/medium sized organisations*, Cheshire, UK: Information Commissioner's Office (ICO). 2:2013.
- Lukáčová, A., Babič, F., & Paralič, J. (2014). Building the prediction model from the aviation incident data. In *Applied Machine Intelligence and Informatics (SAMi)*, 2014 I.E. Herľany, Slovakia: 12th International Symposium on IEEE.
- Lyons, G. (2016). Transport analysis in an uncertain world. *Transport Reviews*, 36, 553–557.
- Mandala, I. G. N. N., Nawangpalupi, C. B., & Praktikto, F. R. (2012). Assessing credit risk: An application of data mining in a rural bank. *Procedia Economics and Finance*, 4, 406–412. [https://doi.org/10.1016/S2212-5671\(12\)00355-3](https://doi.org/10.1016/S2212-5671(12)00355-3)
- Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., Byers, A. H. (2011). *Big data: The next frontier for innovation, competition, and productivity*. New York: McKinsey Global institute.
- Marbán, Ó., Mariscal, G., & Segovia, J. (2009). A data mining & knowledge discovery process model. *Data Mining and Knowledge Discovery in Real Life Applications*, 2009, 8.
- Mariscal, G., Marbán, Ó., & Fernández, C. (2010). A survey of data mining and knowledge discovery process models and methodologies. *The Knowledge Engineering Review*, 25, 137–166. <https://doi.org/10.1017/s0269888910000032>
- McAfee, A., & Brynjolfsson, E. (2012). Big data: The management revolution. *Harvard Business Review*, 90, 60–68.
- Mirabadi, A., & Sharifian, S. (2010). Application of association rules in Iranian railways (RAI) accident data analysis. *Safety Science*, 48, 1427–1435. <https://doi.org/10.1016/j.ssci.2010.06.006>
- Moreno-Díaz, R., Pichler, F., & Quesada-Arencibia, A. (2015). Using data mining to improve the public transport in gran Canaria Island (Vol. 9520). Switzerland: Springer International Publishing. <https://doi.org/10.1007/978-3-319-27340-2>
- Taylor, N., Stott, I., Parker, J., & Bradley, J. (2015). Integrated Transport Planning Ltd., Consulting AGWW, Sustainability CTA, Institute JMHDER. The Transport Data Revolution. *Catapult*.
- Packianather, M. S., Davies, A., Harraden, S., Soman, S., & White, J. (2017). Data mining techniques applied to a manufacturing SME. *Procedia CIRP*, 62, 123–128.
- Piatetsky-Shapiro, G. (1991). Discovery, analysis, and presentation of strong rules. *Knowledge Discovery in Databases*, 248, 255–264.
- Probst, L., Frideres, L., Demetri, D., Vomhof, B., Lonkeu, O.-K., & Luxembourg, P. (2014). *Business innovation observatory—Customer experience*. European Union. Brussels, Belgium: European Commission.
- Pytel, P., Britos, P., & García-Martínez, R. (2013). A proposal of effort estimation method for information mining projects oriented to SMEs. In G. Poels (Ed.), *Enterprise Information Systems of the Future: 6th IFIP WG 8.9 Working Conference, CONFENIS 2012, Ghent, Belgium, September 19–21, 2012, Revised Selected Papers* (pp. 58–74). Berlin, Heidelberg: Springer.
- Raju, Y., Kang, P. S., Moroz, A., Clement, R., Hopwell, A., & Duffy, A. (2016). *Investigating the demand for short-shelf life food products for SME wholesalers*. Leicester, UK: De Montfort University.
- Rebón, F., Castander, I., Argandoña, J., Gerrikagoitia, J. K., & Alzua-Sorzabal, A. (2015). An antifraud system for tourism SMEs in the context of electronic operations with credit cards. *American Journal of Intelligent Systems*, 5, 27–33. <https://doi.org/10.5923/j.ajis.20150501.03>
- Rights EUAF. (2017). *Handbook on European data protection law*. Vienna, Austria: FRA (European Union Agency for Fundamental Rights). Retrieved from <http://fra.europa.eu/en/publication/2014/handbook-european-data-protection-law>
- Shin, I.-H., Park, G.-L., Saha, A., Kwak, H.-Y., & Kim, H. Analysis of moving patterns of moving objects with the proposed framework. *Computational Science and Its Applications—ICCSA 2009*, 2009, 5593:443–452. doi:https://doi.org/10.1007/978-3-642-02457-3_38
- Shmueli, G., Patel, N. R., & Bruce, P. C. (2011). *Data mining for business intelligence: Concepts, techniques, and applications in Microsoft Office Excel with XLMiner*. Hoboken, NJ: John Wiley and Sons).
- Tan, D.-W., Yeoh, W., Boo, Y. L., & Liew, S.-Y. (2013). The impact of feature selection: A data-mining application in direct marketing. *Intelligent Systems in Accounting, Finance and Management*, 20, 23–38. <https://doi.org/10.1002/isaf.1335>
- Viglioni, G. M. C., Cury, M. V. Q., & Silva, P. A. L. D. (2007). Methodology for railway demand forecasting using data mining. In *Proceedings of the SAS Global Forum, Brazil*.
- Wadhwa, K. (2014). *BYTE: Big data roadmap and cross-disciplinary community for addressing societal externalities*, European Data Forum, Athens, Greece.
- Waller, M. A., & Fawcett, S. E. (2013). Data science, predictive analytics, and big data: A revolution that will transform supply chain design and management. *Journal of Business Logistics*, 34, 77–84.
- Wong, J.-Y., & Chung, P.-H. (2007). Managing valuable Taiwanese airline passengers using knowledge discovery in database techniques. *Journal of Air Transport Management*, 13, 362–370. <https://doi.org/10.1016/j.jairtraman.2007.07.001>
- Xindong, W., Xingquan, Z., Gong-Qing, W., & Wei, D. (2014). Data mining with big data. *IEEE Transactions on Knowledge and Data Engineering*, 26, 97–107. <https://doi.org/10.1109/tkde.2013.109>
- Young, T., Fehskens, M., Pujara, P., Burger, M., & Edwards, G. (2010). Utilizing data mining to influence maintenance actions. In *AUTOTESTCON, 2010 IEEE*, Orlando, FL, USA.
- Zeng, D., & Lusch, R. (2013). Big data analytics perspective shifting from transactions to ecosystems. *IEEE Access*, 28(2), 2–5.
- Zhang, C., Huang, Y., & Zong, G. (2010). Study on the application of knowledge discovery in data bases to the decision making of railway traffic safety in China. In: 2010 International Conference on Management and Service Science (MASS). <https://doi.org/10.1109/ICMSS.2010.5577012>
- Zhang, J., Wang, F.-Y., Wang, K., Lin, W.-H., Xu, X., & Chen, C. (2011). Data-driven intelligent transportation systems: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 12, 1624–1639. <https://doi.org/10.1109/tits.2011.2158001>

How to cite this article: Mohd Selamat SA, Prakoonwit S, Sahandi R, Khan W, Ramachandran M. Big data analytics—A review of data-mining models for small and medium enterprises in the transportation sector. *WIREs Data Mining Knowl Discov*. 2018;8:e1238. <https://doi.org/10.1002/widm.1238>