# Initiating and Implementing Data Mining Practices within a Small to Medium-Sized Business Organization

Matthew Baer, Thilini Ariyachandra, and Mark Frolick

*Abstract*—Data analytics, specifically knowledge discovery through data mining has received much attention from industry in recent years. However, much of this attention is mostly focused on implementations in large organizations. This paper discusses the importance of data mining in the context of small to midsized organizations through a brief case study. The paper begins by providing a description of data mining, who or what organizations are involved in data mining, and why data mining is important. Next, this paper describes the successful use of Excel and Visual Basic for Applications to support the data mining of automated test software log files at a local mid-sized manufacturer that is suffering from budgetary constraints and lack of management support. In describing the challenges and the means by which analytics were implemented within the organization, the paper attempts to present a means by which a small to medium business can successfully adopt data mining within the organization.

*Index Terms*—Data analytics, data mining, small to medium business, visual basic, resource constraints.

## I. INTRODUCTION

For the past two decades, organizations have faced an unprecedented influx of data from originating from various sources and stakeholders. Fueled by advances on the Web, the explosion of data volumes within organizations have led to greater sophistication in data analytics solutions available to organizations [1]. Many vendors provide powerful and industrial level analytics solutions to satisfy large corporate data needs. Specifically data mining tools are used to find hidden, valid and actionable information for decision making. Small to medium businesses (SMB) have also amassed data but are unable to use analytics and data mining tools to gain insights they need to effectively run their organizations [2].

SMBs view these solutions as too complex and/or expensive. According to a 2011 survey of small to medium business, more than a quarter of SMBs indicated that getting better insights from data currently within the organization is a top technology challenge in this space [3]. The same survey further revealed that the smaller the company, the less likely they are to use or plan to use analytics/mining solutions. Often times, these organizations find themselves building custom solutions to solve their data mining needs than adopting complex resource intensive vendor tools [4]. The value proposition of implementing data mining solutions in SMBs continues to rise as its benefits to large corporations become established.

Data mining is more frequently used each year in large organizations as more data warehouses are constructed to capture the increasing volumes of data generated by devices and services [5]. These growing volumes of data are employed by a wide variety of industries and organizations for the purpose of Knowledge Discovery in Databases (KDD). Small to medium organizations stand to gain immensely by integrating data mining and KDD into its data analytics arsenal. In order to shed more light on the use of data mining in SMBs, this paper reviews three questions. They are: (1) what is data mining and KDD? (2) Who uses data mining? And finally (3) why should data mining be performed? After discussing the three questions, the paper provide a brief example of an SMBs efforts to begin data mining in its organization's growing volume of test log files using existing commercially available software applications.

## II. DATA MINING

Data mining is succinctly defined as "the science of extracting useful information from large data sets or databases. Its fundamental objective is to provide insight and understanding about the structure of the data and its important features."[6] The two goals of data mining are to describe and summarize information, and to attempt to predict future behavior. Knowledge Discovery in Databases (KDD) "refers to the overall process of discovering useful knowledge from data" [7] using data mining as a step in this process. Data mining describes the general methods used to extract useful and actionable knowledge from databases. Three general steps are typically involved in the data mining process [8]. First, the data must be integrated, cleaned and enhanced before processing. This means that the fields in the data are relevant to each other, and are valuable. Second, algorithms, or data mining methods, must be applied to the data. Some of these data mining methods include classification (mapping data into pre-defined classes), regression (associating data to a variable for the purposes of prediction), clustering (identifying categories of data), dependency modeling, change deviation detection, and summarization. Summarization is a basic method of data mining which describes applying statistical methods to data, such as calculating means, standard deviations, minimums, etc. This data mining technique applied in the case example provided later in this text. Third, the data mining results must be validated and evaluated before knowledge discovery begins so that results can be applied effectively.

M. Baer is with the L-3 Communications Cincinnati Electronics in Mason, Ohio 45040, USA (e-mail: mattbaer@excite.com).

T. Ariyachandra and M. N. Frolick are with the Management Information Systems Department, Xavier University, Cincinnati, OH 45207 USA. (e-mail:ariyachandrat@xavier.edu, frolick@xavier.edu).

Once data mining is defined, who uses data mining and knowledge discovery should be addressed. A wide and growing variety of industries and organizations use data mining and knowledge discovery. The group includes credit card companies, telemarketing and direct marketing firms, airlines, insurance companies, retailers, and financial services and banking companies [9]. Banking and finance use data mining to model and predict fraud, evaluate risk, and analyze profitability. Telephone companies became early adopters of data mining mostly because they have large amounts of quality data, mostly in the form of call records, which enabled them to extract useful information from customers more quickly than other industries [10]. In recent years, scholars, journalists, and intelligence analysts have been employing data mining techniques on text files for knowledge discovery purposes. This includes analyzing large text books, speeches, and phone call transcripts to ascertain patterns and insight [11]. Data mining techniques are also being used by software developers to test for the *correctness* of output in new releases of software [12]. Lastly, various manufacturing companies use data mining to determine appropriate process control parameters and improve the quality of every product, not just a sample [13]. The manufacturing industry subset using data mining includes aluminum processors, semi-conductor manufacturers, electronic assembly and circuit board companies, biotechnology and chemical processing companies, healthcare organizations, and pharmaceutical companies.

Finally, to discuss the value of data mining, the question why should data mining be performed is discussed. In the context of organizations described above, data mining is used for knowledge discovery, and to gain insight to assist in better decision making. The goal of insight and better decision making is to derive economic benefits such as organizational cost savings and increased competitiveness. Others reasons for performing data mining include: adding value to an existing data warehouse, solving a research bottleneck, and increasing operational efficiency [9]. Another reason why data mining should be performed is related to the sheer volume of data being generated. In the past, experts in the industries mentioned above became intimate with the data in their field by constantly and laboriously reviewing and analyzing data to come up with solutions to problems [14]. Since data volumes have increased exponentially, it is now impractical to be an expert analyst without utilizing data mining and knowledge discovery. It is interesting to note that the legal profession struggles with data mining electronically stored information, such as e-mail transactions, during the discovery phase of litigation [15].

Historically, the what's, whose and why's of data mining have pertained to *large* organizations dealing with *large* volumes of data. The impact and implications of data mining on small to mid-size organizations with limited resources is less well known. More research is needed to understand how small to mid-sized organizations effectively and affordably gain insight and knowledge discovery from growing volumes of data.

## III. CASE BACKGROUND

In order to provide a contextual background for use and benefits of data mining for small to midsized organizations, the process flow of a product from a mid-sized manufacturing company in Midwest Ohio in the United States is described. After assembly but prior to shipment, the product unit of the organization is subjected to a series of tests with the purpose of recording various performance parameters, and ensuring proper functionality. In the recent past, product testing within the organization was conducted and recorded manually. Test results were recorded by technicians on sheets of paper and placed in a manila folder, which was subsequently stored in a filing cabinet. As a consequence, product testing was labor intensive and test results were subject to operator recording errors. Also, the performance data collected was not easily accessible and comparing recent performance data to past performance data was, if performed at all, a very laborious task. In an effort to increase testing efficiency, in-house software was developed to automate product testing.

Currently, technicians using the automated test software navigate through a user interface to select appropriate test parameters, and then allow the software to do the rest. It is worth noting that the length of these tests can be anywhere between five and twenty four hours. At the conclusion of testing, technicians are presented with a simple *Pass* or *Fail* test result on the user interface. Results are still recorded manually, but labor hours are significantly reduced, and testing throughput dramatically increased. In addition to the *Pass/Fail* test results, the automated test software also generates result files that contain additional information about product performance. Specifically, four files are generated for every test: a text file that describes the configuration of the test, a text file that records communication between the test software and the product, a comma separated variable file that records one specific product parameter, and a test log file. The test log file is a 40+kilobyte text file containing over 500 lines of text that captures various test parameters and performance results, each line of text with a time/date stamp. Typically, test log files are not reviewed unless the test result is *Fail*. At present, the organization examines the test results and parameters only in reference to failure reports. It is not analyzed or recorded for data comparison purposes.

## IV. AN OPPORTUNITY (WHAT CAN BE DONE BETTER?)

Though automated testing has greatly improved test efficiency and throughput, a great amount of valuable information is not being reviewed, or mined, for the purpose of knowledge discovery. The organization could potentially gain valuable insights by capturing and organizing both the passing and failing units information contained in the test log files for each product tested. Capturing and comparing product parameters and performance results enable members of the manufacturing team to perform, at the very least, the data mining method of summarization [7]. That is, determining arithmetic means and standard deviations to product performance results. Once summarization is established, products can have their test results

cross-referenced. This enables members of the manufacturing team to monitor changes in performance over time and, gain insight into units that are returned for warranty work.

The primary reason the organization has failed to adopt mining methods is time and budgetary constraints. The automated test software developed originally had a time and money limitation attached to it. The first goal of the software developer was to automate the testing. The automated test software was not developed with data collection or analysis in mind. To further develop the software required a significant investment, which was not identified as a management priority at the time. Proving an acceptable return on investment for further software development can be very difficult within a midsized organization where resource constraints can be higher than larger organizations. Gaining management support for this data mining initiative would also be very difficult. As with other IT infrastructure implementations, further development without the support of management leads to a high failure rate [16]. Furthermore, the company experienced a lack of employee buy-in within the production team which further discouraged successful implementation [17]. These potential end users did not see the benefits of data mining, nor demand its implementation. Another barrier to the data mining is sociopolitical in nature. Collecting, summarizing, and analyzing information could shed light on production problems. While from a rational perspective of the organization, one may assume that identifying production problems is a laudable goal. In practice, however, this information tends to put focus on the poor performance of certain manufacturing groups and their managers. Many managers object to initiatives that draw negative attention. This is a sociopolitical reality at many organizations.

## V. AN ALTERNATE SOLUTION

While the value of using a basic data mining technique is evident, data mining and knowledge discovery is currently not encouraged or practiced within the organizations. Resource constraints and management support have prevented the adoption of vendor data mining tools. Instead, the organization can use existing "best-in-breed" applications currently used within the organization such as Microsoft Excel and Visual Basic for Applications (VBA) to accomplish data gathering and analysis [18]. VBA is the structured programming language for Microsoft Office products, including Excel, that uses familiar 'for….Next' loops and 'If….Then….Else' statements [19]. Using VBA with Excel can be a relatively easy and inexpensive method for data mining test log files for products in process, or for data mining existing test log files.

## VI. METHOD

The method for data mining the test log files is presented in block diagram of Fig. 1 in the Excel application environment, a program written in VBA is launched by the user. Consequently, a file dialog box launches that prompts the user to navigate to the folder location of the test log file. Once selected, the VBA program reads each individual line of text

in the test log file, searching for pre-determined parameters and test results. Once discovered, the parameters and test results are written to the Excel spreadsheet. From this point, summarization and analysis can begin regarding the product's performance.
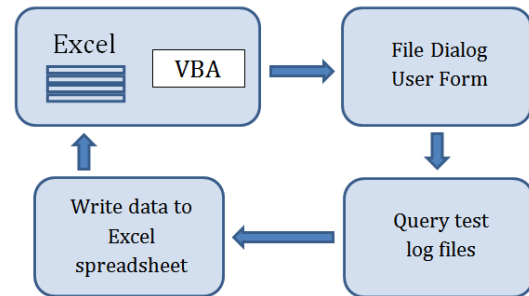


Fig. 1. Using Excel and VBA.

### A. Data Analysis Examples

Some examples of code and results from the organization are described next. The code examples and text have been cleansed and recreated using original code and spreadsheet results. The reason for this manipulation is to prevent the description of proprietary test results and specific test methodologies. The VBA code written for this application was specifically developed to be manipulated by novice programmers/users. Small to midsized manufacturing organizations may not have the necessary knowledge or skills in-house to program and manipulate such applications. Currently, the Web presents various resources which provide general examples of code for novice programmers [20]. It is no longer necessary to re-invent the wheel to develop the VBA code. Fig. 2 displays the sample code for the file dialog box. Fig. 3 displays a sample of the code used to capture data from the test log file. The VBA code uses simplistic 'If….Then' statements to find pre-determined results (labeled generically as R1, R2, R3, etc.) defined in the Excel spreadsheet. Fig. 4 displays the results of the data collected in the Excel spreadsheet. Each row in the spreadsheet represents the test results of one unit. Fig. 4 displays sixteen of the thirty one results that were collected for each unit tested. Note that this collection includes passing and failing examples (see row six, column O).

```
End Sub
Sub Main(strPath)

    Dim fd As FileDialog

    Set fd = Application.FileDialog(msoFileDialogFilePicker)

    Dim vrtSelectedItem As Variant

    With fd
        If .Show = -1 Then

            For Each vrtSelectedItem In .SelectedItems

            Next vrtSelectedItem

        Else
        End If
    End With

    Set fd = Nothing

End Sub
```

Fig. 2 File dialog picker code

### B. An Example of Using Data Collected

Once the data has been collected and organized it becomes much easier to apply statistical methods to the test results, such as correlation, probability distribution, regression, standard deviation, means, etc. This data may also be cross referenced against existing data to solve various

manufacturing issues. For instance, the midsized manufacturing company noticed that various units were failing automated testing due to one particular test result. Intuition led the production team to review a series of suspect components. The initial method of investigation was to manually search each unit's paper work in order to obtain component serial numbers and then correlate those serial numbers to the unit's performance. Once the suspect components were identified and re-worked, the units were sent back to automated testing. Members of the production team were notified that automated test results were being captured and organized using the Excel/VBA application described above. The production team was able to easily and very quickly, verify that all problem units had been identified, re-worked, and successfully tested. The team realized that Information did not have to be "dug up" from the paperwork as stated by a team member. This enabled the production team to wrap up the investigation very quickly. However, the team also realized that difficulty in assessing he return-on-investment to management of the mining technique used. The speed in which the investigation was concluded cannot easily be calculated. "All we could say is that it saved *hours"* stated a team member. Nevertheless, after realizing the value of using the Excel/VBA application to gather information and act on it, four other product lines started using variations of this application to capture product specific test results for both passing and failing units. The midsized manufacturing organization continues to adopt and innovate in their data mining and data analytics efforts within the constraints that they currently have.



Fig. 3. Sample of VBA code.



Fig. 4. Results written to Excel.

## VII. Conclusion

The short case study described above illustrates a very basic use of the data collected from the automated test software test log files that yielded significant results. The organization has yet to experiment with the application of more sophisticated data mining techniques to the information captured. The question remains for the organization as to how many ways can the production process be positively impacted by data mining the test log files of the automated test software. This effort has just begun at this midsized manufacturing organization. They would need to be mindful of the impact of continuously logging test results as time goes on. A strategy for organizing and storing the data must be developed [18]. Once enough data has been gathered and analyzed, a *critical mass* may finally be achieved which helps managers justify the expenditures necessary to develop more sophisticated automated test software. Until that time, the organization will continue to use an Excel/VBA application will continue to serve as the necessary, beginning approach for data mining and knowledge discovery within the organization.

Knowledge discovery maybe in its infancy in small to midsized manufacturing organizations as they slowly begin to understand the value of data mining and analytics to successfully managed business operations. While sophisticated business applications maybe out of reach from small to midsized organizations, it is prudent to be creative in adopting analytics to assess current operations whenever possible. Overtime, it would be possible to see the soft, often intangible, the return on investment from using analytics that would prompt the organization to investment in a more formalized tool for data mining and knowledge discovery in the future.

## References

[1] M. Ohata and A. Kumar, "Big data: A boon to business intelligence," *Financial Executive*, vol. 28, no.7, pp. 63, Sept. 2012.

[2] G. R. Gangadharan, "Business intelligence systems: Design and implementation strategies," in *Proc. 26th Information Technology Interfaces International Conference*, vol. 1, pp. 139–144, 2004.

[3] S. Aggarwal, L. McCabe, and A. Aggarwal. (September, 2011). Small and medium business routes to market study, *SMB Group*, Available: http://www.smb-gr.com/wp-content/uploads/2011/pdfs/RTM_Web_PR_Sept_24_2011.pdf

[4] C. F. Cheung and F. L. Li, "A quantitative correlation coefficient mining method for business intelligence in small and medium enterprises of trading business," *Expert Syst. Appl*. vol. 39 no. 7, pp. 6279-6291, 2012.

[5] Oracle Corporation. (February, 2012). E-Commerce trends for 2012: Mobile and facebook take centerstage as online retailers focus on customers. *Digital Experiences. Oracle Corporation.* [Online] Available: http://www.oracle.com/us/products/applications/ecommerce-trends-2012-1504949.pdf

[6] R. Platon and M. Amazouz, *Application of data mining techniques for industrial process optimization*, Varennes: Natural Resources Canada, July 2007.

[7] U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, "From data mining to knowledge discovery in databases," *AI Magazine*, vol. 17, pp. 37-54, November 1996.

[8] J. Taylor. (April 2012). Information management newsletters article. *Information Management*. [Online] Available: www.information-management.com.

[9] L. Chen, T. Sakaguchi, and M. N. Frolick, "Data mining methods, applications, and tools," *Information Systems Management*, vol. 1, pp. 65-70, 2000.

[10] G. M. Weiss, "Data mining in telecommunications," in *Data Mining and Knowledge Discovery Handbook*, New York: Springer US, 2005, pp. 1189-1201.

[11] A. Don *et al*, "Discovering interereesting usage patterns in text collections: integrating text mining with visualization," in *Proc.CIKM '07*, 2007, pp. 213-222.

[12] M. Last, M. Friedman, and A. Kandel, "The data mining approach to automated software testing," in *Proc.KDD '03*, 2003, pp. 388-396.

[13] A. Kuziak, "Data mining: manufacturing and service applications," *International Journal of Production Research*, vol. 44, pp. 4175-4191, October 2006.

[14] W. Peng, T. Li, and S. Ma, "Mining logs files for data-driven system management," *ACM SIGKDD Explorations Newsletter - Natural language processing and text mining*, vol. 7, no. 1, pp. 44-51, June 2005.

[15] S. Ardisson. (September 2007). Qubit a monthly publication on computer forensic and e-discovery issues. *Bit-X-Bit*. [Online] 7(1). Available: http://www.bit-x-bit.com/sites/default/files/userfiles/files/articles/qubit-sept-07-the-exponential-growth-of-electronic-information-and-its-impact-on-litigation.pdf

[16] K. Lindsey and M. N. Frolick, "Critical factors for data warehouse failure," J*ournal of Data Warehousing*, vol. 1, pp. 48-54, 2003.

[17] R. Hobek, T. R. Ariyachandra, and M. N. Frolick, "The importance of soft skills in business intelligence implementations," *Business Intelligence Journal*, pp. 28-36, January 2009.

[18] M. N. Frolick and M. von Oven, "Taking the repetition out of research and development: The BI collaboration approach," B*usiness Intelligence Journal*, vol. 11, no. 3, pp. 21-26, July 2006.

[19] Microsoft Corporation. (2012). Introducing Visual Basic for Applications. Microsoft Dot Net Web site. [Online]. Available: http://msdn.microsoft.com/en-us/library/office/aa188202%28v=office.10%29.aspx

[20] J. Beaucaire. (January 2011). Forum: Mr. Excel.com. *Mr. Excel.com*. [Online] Available: http://www.mrexcel.com/forum/excel-questions/521441-using-file-name-returned-visual-basic-applications-file-open-dialog.html

**Matthew Baer** earned a Bachelor of Science in Mass Communications Media Management from Miami University, Oxford, Ohio. In 2001 the author earned an Associate of Applied Science in Electro-Mechanical Engineering Technology from Cincinnati State, Cincinnati, Ohio. The author is currently working on his Master's in Business Administration from the Williams College of Business at Xavier University in Cincinnati, Ohio, USA.

He has worked for four years at L-3 Communications Cincinnati Electronics in Mason, Ohio, performing environmental and electro-magnetic interference (EMI) testing in support of manufacturing operations. His current position is Support Engineer for the Environmental Services Group. He also has seven years of manufacturing experience in the semi-conductor equipment industry and two years of experience working for a commercial building service company specializing in energy efficiency. His research interests include business intelligence and data mining.

**Thilini Ariyachandra** is a member of IACSIT and obtained her Ph. D. from the Terry School of Business at the University of Georgia, USA. She is an Associate Professor of Management Information Systems in the Williams College of Business at Xavier University in Cincinnati, Ohio, USA. Her main research area is business intelligence and data warehousing. Specifically she focuses on BI infrastructure architecture and development methodologies, BI agility and success and BI education. Her work has been published in various academic outlets including *Decision Support Systems, Communications of the ACM, and Communications of the AIS*.

**Mark N. Frolick** obtained his Ph.D. from the Terry School of Business at the University of Georgia, USA. He is a Professor of MIS in the Williams College of Business at Xavier University in Cincinnati, Ohio, USA and the holder of the Western & Southern Chair in Management Information Systems. Dr. Frolick was formerly Professor of MIS and Associate Director of the FedEx Center for Cycle Time Research at The University of Memphis. Dr. Frolick has over 20 years' experience in the information systems field. In addition to working for The Southern Company and Georgia Power, he has worked as a consultant for numerous Fortune 500 companies including FedEx, Ford, Hewlett Packard, Medtronic, and Texas Instruments.

He is considered to be a leading authority on business intelligence. His specialties include business performance management, business intelligence, data warehousing, executive information systems, e-business, cycle time reduction, and the diffusion of information technology in organizations.

Dr. Frolick has authored over 130 articles. His research has appeared in such prestigious journals as *MIS Quarterly, Decision Sciences, Journal of Management Information Systems, Decision Support Systems, and Information & Management*. He also worked with Dr. James Wetherbe on the book *Systems Analysis and Design: Best Practices* (West Publishing, 1994). This book was ranked by Computing Newsletter as the top textbook on the topic. Additionally, Dr. Frolick serves as a consulting editor for several publishing companies.