



Data Mining: An Overview

Sabale Swarupa Nivrutti, Research Scholar, Shri JJT University Jhunjhunu, Rajasthan
Dr. Hiren Dand, Research Guide, Department of Computer Science and Engineering,
Shri JJT University Jhunjhunu, Rajasthan

Abstract

One way to think of data mining is as the process of extracting predictive and intelligent models from large amounts of data. It is the skill of taking vast volumes of data and turning it into meaningful information. It processes massive amounts of data by fusing advanced algorithms with conventional data analysis. It is an interdisciplinary field that combines ideas from information theory, computing, statistics, machine learning, database systems, and pattern recognition. It truly has the potential to be taught in electrical engineering programmes. This paper's major goal is to give a concise overview of data mining.

Keywords: Data Mining, Evaluation, Modelling

Introduction

We can now collect and store enormous amounts of data thanks to technology. An estimate of the amount of fresh data generated annually is greater than 15 exabytes. Our ability to create and gather data has greatly expanded due to the widespread usage of the World Wide Web and its related information services, like Google, Yahoo, Excite, InfoSeek, and American Online. We now need methods to help us turn the data into actionable knowledge and information because of its rapid increase.

Facts, language, or numbers that a computer can process are called data. Transactional data, including sales, pricing, payroll, and accounting, may be one way they appear. When you use your credit card to make a purchase or browse the internet, for instance, you create data. Data mining is the process of obtaining valuable information from massive databases, which is a result of the trend of centralising an organization's data.

In recent years, there has been a lot of interest in the concept of data mining. Tom Khabaza launched it in the early 1990s. The process of gleaning knowledge and patterns from massive amounts of data is called data mining, or knowledge discovery databases (KDD). It might also be thought of as the process of combining multiple perspectives to analyse data and synthesise it into actionable insights. It is a branch of computer science that is multidisciplinary. Artificial intelligence, machine learning, database systems, pattern recognition, warehousing, data visualisation, and statistics are all involved in data mining, as Figure 1 illustrates.

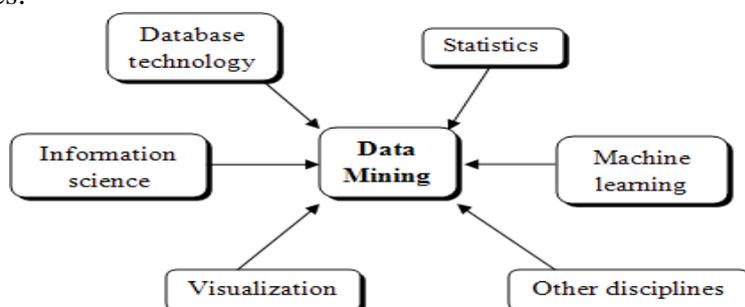


Figure 1 Data mining as an interdisciplinary field (Adapted from Han and Kamber, 2006).

How does data mining work?

Information extraction is comparable to the process of removing metal from ore. It is best to think about data mining as a process. The following steps are involved in the process:

- Effective data processing and storage
- Choosing how many variables to examine
- Visualising and summarising the data

- Using statistics like mean, percentiles, standard deviation, and connection
- Use analytic techniques like k-mean clustering, regression, and closest neighbour algorithms.
- Put the knowledge gleaned from the study into practice.

These are the actions that are typically done when mining data. However, there are certain data mining process models. The Cross- Industry Standard Process for Data Mining (-DM) is the most widely used. This is a data mining open standard. A group of European businesses put forth the proposal in the middle of the 1990s. Figure 2 provides an illustration of it. Every project starts with an understanding of the business and proceeds through the five stages of the procedure.

- Business understanding: this entails creating a project plan and outlining the goals your organisation has for the project.
- Data Understanding: This stage begins with the collection of data and continues with its quality verification to ensure that it meets your objectives.
- Data Preparation: This stage involves choosing any required training and test samples in addition to cleaning the data. The majority of data miners' time is spent in this stage.
- Modelling: During this stage, particular modelling methods are chosen and used with the data. Usually, a given data mining problem can be solved using multiple approaches.
- Evaluation: This entails assessing the models to ascertain how accurately the project's aims and objectives are being met. They return to the modelling stage if they are not satisfied with the model.
- Deployment: This entails putting the acquired knowledge in a format that the client can utilise, like a table or graph. The deployment process will often be completed by the customer rather than the data analyst.

It does not imply that one individual will handle every process. It's a group endeavour.

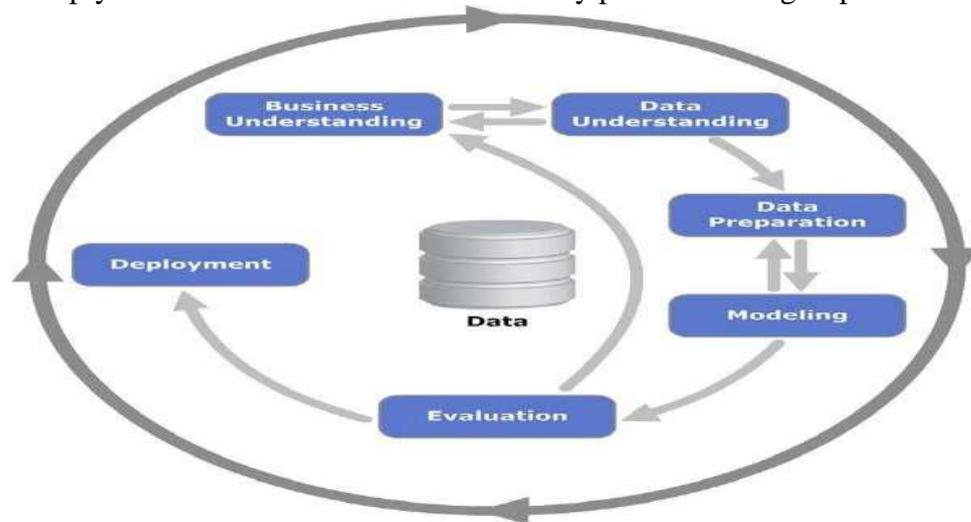


Figure 2 The -DM process.

Applications of Data Mining

The field of data mining has immense promise. Business, the retail sector, telecommunications, intrusion detection, biological data analysis, healthcare, geosciences, and computer security have all seen significant success with the application of data mining. We'll talk about a few of these.

- Intrusion Detection: Any action that jeopardises the availability, integrity, or confidentiality of network resources is considered an intrusion. By creating an intrusion detection data mining algorithm, data mining technology can be used for intrusion detection.
- Telecommunication: This sector offers a wide range of services, including web data transmission, fax, pager, cell phone, Internet messenger, pictures, and email. Within the telecommunications sector, data mining aids in pattern recognition, fraud detection, resource optimisation, and service quality enhancement.
- Business: Finance, particularly in banks and insurance businesses, and e-business are two



of the most significant business domains. Systematic data analysis and data mining are made easier by the generally high quality and dependability of the financial data used in banking and the financial industry.

- Retail business: Due to its ability to gather vast amounts of data on sales, consumer purchase histories, the transportation, consumption, and services of goods, data mining is a fantastic tool in the retail business. Finding client purchasing trends and patterns is made easier with the use of data mining. This results in better client retention and satisfaction as well as higher-quality customer service. Businesses can also benefit from data mining by using it to account for irregular transactions, item flow, and consumption peaks.
- Geosciences: Petrophysical data are used to anticipate reservoirs and identify relationships based on data mining techniques. In complex geological conditions, the logging data are used to identify the effective reservoirs and assess the fuzzy reservoirs.
- Healthcare: Data mining is becoming more and more important in this field. For instance, data mining can assist healthcare organisations in making decisions about customer relationship management, insurers in identifying fraud and abuse, doctors in identifying best practices and treatments, and patients in receiving healthcare services at a lower cost.

Conclusion

Finding patterns in a lot of data that are useful is known as data mining. Although data mining is a relatively new topic of study, it is acknowledged as one that is developing quickly. To be a data miner, one does not need to be an expert in statistics or a computer programmer, however it would not hurt to have some exposure to statistical analysis.

Data has become ubiquitous in nearly every field. It truly has the potential to be taught in electrical engineering programmes. The 2013-founded Society of Data Miners (www.socdm.org) is a useful resource for career and knowledge advancement.

References:

1. D. Braha and A. Shmilovici, "Data mining for improving a cleaning process in the semiconductor industry," *IEEE Transactions on Semiconductor Manufacturing*, vol. 15, no. 1, Feb, 2002, pp. 91-101.
2. I.H. Witten and E. Frank, *Data Mining* (Moran Kaufmann Publishers, Amsterdam, 2005, 2nd ed.).
3. *IEEE Transactions on Knowledge and Data Engineering*, vol. 8, no. 6, Dec. 1996, pp. 866-883.
4. J. Han and M. Kamber, *Data Mining: Concepts and Techniques* (Morgan Kaufmann, San Francisco, CA, 3rd ed., 2011).
5. J. Ledolter, *Data Mining and Business Analytics with R* (John Wiley & Sons, Hoboken, NJ, 2013).
6. M. North, *Data Mining for the Masses* (Global Text, Lexington, KY, 2012).
7. M. S. Brown, *Data Mining for Dummies* (John Wiley & Sons, Hoboken, NJ, 2014).
8. M. S. Chen, J. Han, and P. S. Yu, "Data mining: an overview from a database perspective,"
9. P. N. Tan, M. Steinbach, and V. Kumar, *Introduction to Data Mining*
10. V.K. Deepa and J. R. Geetha, "Rapid development of applications in data mining," *Proceedings of 2013 International Conference on Green High Performance Computing*, Mar. 2013.