

Data Mining With Application of Big Data



Aqeel Abbood Rahmah

Master of Science (Information System),
Nizam College (Autonomous), O.U.,
Basheer Bagh, Hyderabad.

Abstract:

Big Data concern large-volume, complex, growing data sets with multiple, autonomous sources. With the fast development of networking, data storage, and the data collection capacity, Big Data are now rapidly expanding in all science and engineering domains, including physical, biological and biomedical sciences. This paper presents a HACE theorem that characterizes the features of the Big Data revolution, and proposes a Big Data processing model, from the data mining perspective. This data-driven model involves demand-driven aggregation of information sources, mining and analysis, user interest modelling, and security and privacy considerations. We analyze the challenging issues in the data-driven model and also in the Big Data revolution. The most fundamental challenge for Big Data applications is to explore the large volumes of data and extract useful information or knowledge for future actions. To support Big Data mining, high-performance computing platforms are required, which impose systematic designs to unleash the full power of the Big Data.

INTRODUCTION:

Data mining (the analysis step of the “Knowledge Discovery in Databases” process, or KDD), [1] an interdisciplinary subfield of computer science, [2][3][4] is the computational process of discovering patterns in large data sets involving methods at the intersection of artificial intelligence, machine learning, statistics, and database systems. [2] The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use.

[2] Aside from the raw analysis step, it involves database and data management aspects, data pre-processing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of discovered structures, visualization, and online updating. [2] The term is a misnomer, because the goal is the extraction of patterns and knowledge from large amount of data, not the extraction of data itself. [5] It also is a buzzword [6] and is frequently applied to any form of large-scale data or information processing (collection, extraction, warehousing, analysis, and statistics) as well as any application of computer decision support system, including artificial intelligence, machine learning, and business intelligence. The popular book “Data mining: Practical machine learning tools and techniques with Java” [7] (which covers mostly machine learning material) was originally to be named just “Practical machine learning”, and the term “data mining” was only added for marketing reasons. [8] Often the more general terms “(large scale) data analysis”, or “analytics” – or when referring to actual methods, artificial intelligence and machine learning – are more appropriate.

The actual data mining task is the automatic or semi-automatic analysis of large quantities of data to extract previously unknown interesting patterns such as groups of data records (cluster analysis), unusual records (anomaly detection) and dependencies (association rule mining). This usually involves using database techniques such as spatial indices. These patterns can then be seen as a kind of summary of the input data, and may be used in further analysis or, for example, in machine learning and predictive analytics. For example, the data mining step might identify multiple groups in the data, which can then be used to obtain more accurate prediction results by a decision support system.

Neither the data collection, data preparation, nor result interpretation and reporting are part of the data mining step, but do belong to the overall KDD process as additional steps. The related terms data dredging, data fishing, and data snooping refer to the use of data mining methods to sample parts of a larger population data set that are (or may be) too small for reliable statistical inferences to be made about the validity of any patterns discovered. These methods can, however, be used in creating new hypotheses to test against the larger data populations.

DATA MINING:

Data Mining is defined as the procedure of extracting information from huge sets of data. In other words, we can say that data mining is mining knowledge from data. The tutorial starts off with a basic overview and the terminologies involved in data mining and then gradually moves on to cover topics such as knowledge discovery, query language, classification and prediction, decision tree induction, cluster analysis, and how to mine the Web. There is a huge amount of data available in the Information Industry.

This data is of no use until it is converted into useful information. It is necessary to analyze this huge amount of data and extract useful information from it. Extraction of information is not the only process we need to perform; data mining also involves other processes such as Data Cleaning, Data Integration, Data Transformation, Data Mining, Pattern Evaluation and Data Presentation. Once all these processes are over, we would be able to use this information in many applications such as Fraud Detection, Market Analysis, Production Control, Science Exploration, etc.

What is Data Mining?

Data Mining is defined as extracting information from huge sets of data. In other words, we can say that data mining is the procedure of mining knowledge from data. The information or knowledge extracted so can be used for any of the following applications –

- Market Analysis
- Fraud Detection
- Customer Retention
- Production Control
- Science Exploration

Data Mining Applications:

Data mining is highly useful in the following domains –

- Market Analysis and Management
- Corporate Analysis & Risk Management
- Fraud Detection

Apart from these, data mining can also be used in the areas of production control, customer retention, science exploration, sports, astrology, and Internet Web Surf-Aid.

Market Analysis and Management:

Listed below are the various fields of market where data mining is used –

- Customer Profiling – Data mining helps determine what kind of people buy what kind of products.
- Identifying Customer Requirements – Data mining helps in identifying the best products for different customers. It uses prediction to find the factors that may attract new customers.
- Cross Market Analysis – Data mining performs association/correlations between product sales.
- Target Marketing – Data mining helps to find clusters of model customers who share the same characteristics such as interests, spending habits, income, etc.
- Determining Customer purchasing pattern – Data mining helps in determining customer purchasing pattern.
- Providing Summary Information – Data mining provides us various multidimensional summary reports.

Corporate Analysis and Risk Management:

Data mining is used in the following fields of the Corporate Sector –

- Finance Planning and Asset Evaluation – It involves cash flow analysis and prediction, contingent claim analysis to evaluate assets.

- Resource Planning – It involves summarizing and comparing the resources and spending.

- Competition – It involves monitoring competitors and market directions.

Fraud Detection:

Data mining is also used in the fields of credit card services and telecommunication to detect frauds. In fraud telephone calls, it helps to find the destination of the call, duration of the call, time of the day or week, etc. It also analyzes the patterns that deviate from expected norms.

LITERATURE SURVEY:

“Algorithms for Mining the Evolution of Conserved Relational States in Dynamic Networks, AUTHORS: R. Ahmed and G. Karypis”

Dynamic networks have recently being recognized as a powerful abstraction to model and represent the temporal changes and dynamic aspects of the data underlying many complex systems. Significant insights regarding the stable relational patterns among the entities can be gained by analyzing temporal evolution of the complex entity relations. This can help identify the transitions from one conserved state to the next and may provide evidence to the existence of external factors that are responsible for changing the stable relational patterns in these networks. This paper presents a new data mining method that analyzes the time-persistent relations or states between the entities of the dynamic networks and captures all maximal non-redundant evolution paths of the stable relational states. Experimental results based on multiple datasets from real-world applications show that the method is efficient and scalable.

“Novel Approaches to Crawling Important Pages Early AUTHORS: M.H. Alam, J.W. Ha, and S.K. Lee”

Web crawlers are essential to many Web applications, such as Web search engines, Web archives, and Web directories, which maintain Web pages in their local repositories. In this paper, we study the problem of crawl scheduling that biases crawl ordering toward important pages.

We propose a set of crawling algorithms for effective and efficient crawl ordering by prioritizing important pages with the well-known Page Rank as the importance metric. In order to score URLs, the proposed algorithms utilize various features, including partial link structure, inter-host links, page titles, and topic relevance. We conduct a large-scale experiment using publicly available data sets to examine the effect of each feature on crawl ordering and evaluate the performance of many algorithms. The experimental results verify the efficacy of our schemes. In particular, compared with the representative Rank Mass crawler, the FPR-title-host algorithm reduces computational overhead by a factor as great as three in running time while improving effectiveness by 5 % in cumulative Page Rank

“Identifying Influential and Susceptible Members of Social Networks AUTHORS: S. Aral and D. Walker”

Identifying social influence in networks is critical to understanding how behaviours spread. We present a method that uses in vivo randomized experimentation to identify influence and susceptibility in networks while avoiding the biases inherent in traditional estimates of social contagion. Estimation in a representative sample of 1.3 million Face book users showed that younger users are more susceptible to influence than older users, men are more influential than women, women influence men more than they influence other women, and married individuals are the least susceptible to influence in the decision to adopt the product offered. Analysis of influence and susceptibility together with network structure revealed that influential individuals are less susceptible to influence than noninfluential individuals and that they cluster in the network while susceptible individuals do not, which suggests that influential people with influential friends may be instrumental in the spread of this product in the network.

“Big Privacy: Protecting Confidentiality in Big Data AUTHORS: A. Machanavajhala and J.P. Reiter “

A tremendous amount of data about individuals – e.g., demographic information, internet activity, energy usage, communication patterns and social interactions –

are being collected and analyzed by many national statistical agencies, survey organizations, medical centers, and Web and social networking companies. Wide dissemination of microdata (data at the granularity of individuals) facilitates advances in science and public policy, helps citizens to learn about their societies, and enables students to develop skills at data analysis. Often, however, data producers cannot release microdata as collected, because doing so could reveal data subjects' identities or values of sensitive attributes. Failing to protect confidentiality (when promised) is unethical and can cause harm to data subjects and the data provider. It even may be illegal, especially in government and research settings. For example, if one reveals confidential data covered by the U. S. Confidential Information Protection and Statistical Efficiency Act, one is subject to a maximum of \$250,000 in fines and a five year prison term.

“Analyzing Collective Behavior from Blogs Using Swarm Intelligence AUTHORS: S. Banerjee and N. Agarwal”

With the rapid growth of the availability and popularity of interpersonal and behavior-rich resources such as blogs and other social media avenues, emerging opportunities and challenges arise as people now can, and do, actively use computational intelligence to seek out and understand the opinions of others. The study of collective behavior of individuals has implications to business intelligence, predictive analytics, customer relationship management, and examining online collective action as manifested by various flash mobs, the ArabSpring (2011) and other such events.

In this article, we introduce a natureinspired theory to model collective behavior from the observed data on blogs using swarm intelligence, where the goal is to accurately model and predict the future behavior of a large population after observing their interactions during a training phase. Specifically, an ant colony optimization model is trained with behavioral trend from the blog data and is tested over real-world blogs. Promising results were obtained in trend prediction using ant colony based pheromone classifier and CHI statistical measure. We provide empirical guidelines for selecting suitable parameters for the model, conclude with interesting observations, and envision future research directions.

PROBLEM STATEMENT:

- » The rise of Big Data applications where data collection has grown tremendously and is beyond the ability of commonly used software tools to capture, manage, and process within a “tolerable elapsed time.” The most fundamental challenge for Big Data applications is to explore the large volumes of data and extract useful information or knowledge for future actions. In many situations, the knowledge extraction process has to be very efficient and close to real time because storing all observed data is nearly infeasible.
- » The unprecedented data volumes require an effective data analysis and prediction platform to achieve fast response and real-time classification for such Big Data.

Drawbacks:

- The challenges at Tier I focus on data accessing and arithmetic computing procedures. Because Big Data are often stored at different locations and data volumes may continuously grow, an effective computing platform will have to take distributed large-scale data storage into consideration for computing.
- The challenges at Tier II center around semantics and domain knowledge for different Big Data applications. Such information can provide additional benefits to the mining process, as well as add technical barriers to the Big Data access (Tier I) and mining algorithms (Tier III).
- At Tier III, the data mining challenges concentrate on algorithm designs in tackling the difficulties raised by the Big Data volumes, distributed data distributions, and by complex and dynamic data characteristics.

PROBLEM DEFINITION:

We propose a HACE theorem to model Big Data characteristics. The characteristics of HACE make it an extreme challenge for discovering useful knowledge from the Big Data.

The HACE theorem suggests that the key characteristics of the Big Data are

- 1) huge with heterogeneous and diverse data sources,

2) autonomous with distributed and decentralized control, and
 3) complex and evolving in data and knowledge associations.

To support Big Data mining, high-performance computing platforms are required, which impose systematic designs to unleash the full power of the Big Data.

ADVANTAGES:

Provide most relevant and most accurate social sensing feedback to better understand our society at real-time.

IMPLEMENTATION:

Implementation is the stage of the project when the theoretical design is thrtturned out into a working system. Thus it can be considered to be the most critical stage in achieving a successful new system and in giving the user, confidence that the new system will work and be effective. The implementation stage involves careful planning, investigation of the existing system and it's constraints on implementation, designing of methods to achieve changeover and evaluation of changeover methods.

Integrating and mining bio data:

We have integrated and mined bio data from multiple sources to decipher and utilize the structure of biological networks to shed new insights on the functions of biological systems. We address the theoretical underpinnings and current and future enabling technologies for integrating and mining biological networks. We have expanded and integrated the techniques and methods in information acquisition, transmission, and processing for information networks. We have developed methods for semantic-based data integration, automated hypothesis generation from mined data, and automated scalable analytical tools to evaluate simulation results and refine models.

Big Data Fast Response:

We propose to build a stream-based Big Data analytic framework for fast response and real-time decision making. Designing Big Data sampling mechanisms to reduce Big Data volumes to a manageable size for processing

- Building prediction models from Big Data streams. Such models can adaptively adjust to the dynamic changing of the data, as well as accurately predict the trend of the data in the future; and
- A knowledge indexing framework to ensure real-time data monitoring and classification for Big Data applications.

Pattern matching and mining:

We perform a systematic investigation on pattern matching, pattern mining with wildcards, and application problems as follows:

- Exploration of the NP-hard complexity of the matching and mining problems,
- Multiple patterns matching with wildcards,
- Approximate pattern matching and mining, and
- Application of our research onto ubiquitous personalized information processing and bioinformatics.

Key technologies for integration and mining:

We have performed an investigation on the availability and statistical regularities of multisource, massive and dynamic information, including cross-media search based on information extraction, sampling, uncertain information querying, and cross-domain and cross-platform information polymerization. To break through the limitations of traditional data mining methods, we have studied heterogeneous information discovery and mining in complex inline data, mining in data streams, multigranularity knowledge discovery from massive multisource data, distribution regularities of massive knowledge, quality fusion of massive knowledge.

Group influence and interactions:

- Employing group influence and information diffusion models, and deliberating group interaction rules in social networks using dynamic game theory.
- Studying interactive individual selection and effect evaluations under social networks affected by group emotion, and analyzing emotional interactions and influence among individuals and groups, and
- Establishing an interactive influence model and its computing methods for social network groups, to reveal the interactive influence effects and evolution of social networks.

CONCLUSION:

Driven by real-world applications and key industrial stakeholders and initialized by national funding agencies, managing and mining Big Data have shown to be a challenging yet very compelling task. While the term Big Data literally concerns about data volumes, our HACE theorem suggests that the key characteristics of the Big Data are 1) huge with heterogeneous and diverse data sources, 2) autonomous with distributed and decentralized control, and 3) complex and evolving in data and knowledge associations. Such combined characteristics suggest that Big Data require a “big mind” to consolidate data for maximum values. To explore Big Data, we have analyzed several challenges at the data, model, and system levels. To support Big Data mining, high-performance computing platforms are required, which impose systematic designs to unleash the full power of the Big Data. At the data level, the autonomous information sources and the variety of the data collection environments, often result in data with complicated conditions, such as missing/uncertain values. In other situations, privacy concerns, noise, and errors can be introduced into the data, to produce altered data copies. Developing a safe and sound information sharing protocol is a major challenge.

At the model level, the key challenge is to generate global models by combining locally discovered patterns to form a unifying view. This requires carefully designed algorithms to analyze model correlations between distributed sites, and fuse decisions from multiple sources to gain a best model out of the Big Data. At the system level, the essential challenge is that a Big Data mining framework needs to consider complex relationships between samples, models, and data sources, along with their evolving changes with time and other possible factors. A system needs to be carefully designed so that unstructured data can be linked through their complex relationships to form useful patterns, and the growth of data volumes and item relationships should help form legitimate patterns to predict the trend and future. We regard Big Data as an emerging trend and the need for Big Data mining is arising in all science and engineering domains. With Big Data technologies, we will hopefully be able to provide most relevant and most accurate social sensing feedback to better understand our society at real time. We can further stimulate the participation of the public audiences in the data production circle for societal and economical events. The era of Big Data has arrived.

FUTURE ENHANCEMENT:

Open nature of Big Data systems exposes them to malicious activity. Building trust relationships among peers can mitigate attacks of malicious peers. In future we can enhance the system with distributed algorithms that enable a peer to reason about trustworthiness of other peers based on past interactions and recommendations. Peers create their own trust network in their proximity by using local information available and do not try to learn global trust information. Two contexts of trust, service, and recommendation contexts are defined to measure trustworthiness in providing services and giving recommendations. Interactions and recommendations are evaluated based on importance, recency, and peer satisfaction parameters. Additionally, recommender’s trustworthiness and confidence about a recommendation are considered while evaluating recommendations. Simulation experiments on a file sharing application show that the proposed model can mitigate attacks on 16 different malicious behaviour models. In the experiments, good peers were able to form trust relationships in their proximity and isolate malicious peers.

REFERENCES:

- [1] R. Ahmed and G. Karypis, “Algorithms for Mining the Evolution of Conserved Relational States in Dynamic Networks,” *Knowledge and Information Systems*, vol. 33, no. 3, pp. 603-630, Dec. 2012.
- [2] M.H. Alam, J.W. Ha, and S.K. Lee, “Novel Approaches to Crawling Important Pages Early,” *Knowledge and Information Systems*, vol. 33, no. 3, pp 707-734, Dec. 2012.
- [3] S. Aral and D. Walker, “Identifying Influential and Susceptible Members of Social Networks,” *Science*, vol. 337, pp. 337-341, 2012.
- [4] A. Machanavajjhala and J.P. Reiter, “Big Privacy: Protecting Confidentiality in Big Data,” *ACM Crossroads*, vol. 19, no. 1, pp. 20-23, 2012.
- [5] S. Banerjee and N. Agarwal, “Analyzing Collective Behavior from Blogs Using Swarm Intelligence,” *Knowledge and Information Systems*, vol. 33, no. 3, pp. 523-547, Dec. 2012.

- [6] E. Birney, "The Making of ENCODE: Lessons for Big-Data Projects," *Nature*, vol. 489, pp. 49-51, 2012.
- [7] J. Bollen, H. Mao, and X. Zeng, "Twitter Mood Predicts the Stock Market," *J. Computational Science*, vol. 2, no. 1, pp. 1-8, 2011.
- [8] S. Borgatti, A. Mehra, D. Brass, and G. Labianca, "Network Analysis in the Social Sciences," *Science*, vol. 323, pp. 892-895, 2009.
- [9] J. Bughin, M. Chui, and J. Manyika, *Clouds, Big Data, and Smart Assets: Ten Tech-Enabled Business Trends to Watch*. McKinsey Quarterly, 2010.
- [10] D. Centola, "The Spread of Behavior in an Online Social Network Experiment," *Science*, vol. 329, pp. 1194-1197, 2010.
- [11] E.Y. Chang, H. Bai, and K. Zhu, "Parallel Algorithms for Mining Large-Scale Rich-Media Data," *Proc. 17th ACM Int'l Conf. Multimedia, (MM '09)*, pp. 917-918, 2009.
- [12] R. Chen, K. Sivakumar, and H. Kargupta, "Collective Mining of Bayesian Networks from Distributed Heterogeneous Data," *Knowledge and Information Systems*, vol. 6, no. 2, pp. 164-187, 2004.
- [13] Y.-C. Chen, W.-C. Peng, and S.-Y. Lee, "Efficient Algorithms for Influence Maximization in Social Networks," *Knowledge and Information Systems*, vol. 33, no. 3, pp. 577-601, Dec. 2012.
- [14] C.T. Chu, S.K. Kim, Y.A. Lin, Y. Yu, G.R. Bradski, A.Y. Ng, and K. Olukotun, "Map-Reduce for Machine Learning on Multicore," *Proc. 20th Ann. Conf. Neural Information Processing Systems (NIPS '06)*, pp. 281-288, 2006.
- [15] G. Cormode and D. Srivastava, "Anonymized Data: Generation, Models, Usage," *Proc. ACM SIGMOD Int'l Conf. Management Data*, pp. 1015-1018, 2009.
- [16] S. Das, Y. Sismanis, K.S. Beyer, R. Gemulla, P.J. Haas, and J. McPherson, "Ricardo: Integrating R and Hadoop," *Proc. ACM SIGMOD Int'l Conf. Management Data (SIGMOD '10)*, pp. 987-998. 2010.
- [17] P. Dewdney, P. Hall, R. Schilizzi, and J. Lazio, "The Square Kilometre Array," *Proc. IEEE*, vol. 97, no. 8, pp. 1482-1496, Aug. 2009.
- [18] P. Domingos and G. Hulten, "Mining High-Speed Data Streams," *Proc. Sixth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (KDD '00)*, pp. 71-80, 2000.
- [19] G. Duncan, "Privacy by Design," *Science*, vol. 317, pp. 1178-1179, 2007.
- [20] B. Efron, "Missing Data, Imputation, and the Bootstrap," *J. Am. Statistical Assoc.*, vol. 89, no. 426, pp. 463-475, 1994.

AUTHORS BIOGRAPHY:

Aqeel Abbood Rahmah, pursuing his Master of Science in Information System, from Nizam College (Autonomous), O.U, Basheer Bagh, Hyderabad, India.