

EFFICIENT DATA MINING

R. K. CHAUHAN & ABHISHEK TANEJA

ABSTRACT

To increase the speed of the data mining process is an important goal. In pursuit of that goal, we propose several areas where basic research and development may spur progress to improve the speed of knowledge discovery. We proposed three research areas: automated control of the mining process, facilitated data preparation, and the automatic identification of inadequate models.

Keyword: Data Mining, models, knowledge discovery, process

1. INTRODUCTION

Due to iterative nature of data mining process, performance is the biggest victim. So, the primary objective of this paper is to make the process more straightforward and, to the extent possible, to automate it. This objective is within the view of how data mining is used by the mass of firms and governmental agencies. Data mining is often branded a special project, a unique project to undertaken by specific group. This technology will be mature when it evolves from a one-shot process that is performed by experts to one that is performed without much human intervention and to one that is integrated into business processes of a firm^[1]. Just as data base and OLAP reports are provided daily to apprise executives and others throughout organization of the up-to-the-day state of the business, we anticipate that data mining models will soon be provided on a daily basis to people holding a variety of positions within an organisation. For this to happen, it will be desirable for the knowledge discovery process to be straightened as much as possible, to eliminate unnecessary iterations.

Of course, there are many other considerations in an effort to bring a sophisticated to a large, lay populace, just a for any knowledge management efforts. Our goal here is merely to examine three agenda items where additional research may help to automate and accelerate the mining process from extraction to extraction.

The purpose of this paper is to suggest several topics that would need to be addressed to make data mining better understood and more efficient. This paper is conceived as a vision paper, but for each research topic we offer high-level technical approaches that might be successfully applied to the problem.

2. DATA MINING STEPS

In the commercial and governmental world, there are steps typically taken to go from data to the fielded implementation of a data mining model. Standard steps for a target marketing campaign might be:

1. Extract the data from one or more data sources.
2. Prepare the data: place the data in a format required for input to the data mining software, prepare any meta-data to describe the fields in the data records, impute missing values, remove missing records, etc.
3. Build a predictive propensity model to assign to each record a score that indicates the propensity of establishment to purchase some items of interest.
4. Rank the establishments according to their predictive scores.
5. Determine and extract the top candidates for campaign treatments.
6. Measures the results of the campaign.
7. Measure the results of the campaign.

In an outer, business process loop, measuring results identifies new problems and provides new data to kick start another round of data mining. The inner cycle of returning to a previous step, making changes, and retracing the subsequent steps is the one that reflects the difficulty of data mining. If the analyst realizes that she failed to extract a field from a legacy database, she will have to return to step 1 in the list. If she made an error in assigning classes to the dependent variable, then she will have to return to an early data preparation step, step 2. In general, it is easy to envision a return from any step to any previous one.

The twists of the mining task have been recognized empirically and have been the impetus for work. This system supports a human-analyst-centered view of the mining process, by helping keep track of files, mining results, etc.

Each step in this process certainly affects later steps^[2]. But the degree and the manner to which each step affects later ones are usually unclear. For example, in step1, the effect of including or omitting a data source on the ultimate results of the campaign is unknown an unknowable, in general. So at the same time that data mining analysts go round in circles, retracing previous steps, the effect of those temporary setbacks is really not known. They can only be surmised, based on experience. Our thesis is that if we possessed a model of the entire mining process, we would know whether it was worthwhile to return to earlier stages. This issue is not just an academic nicety. Data mining practitioners are in short supply, their time is pricey, and they are under pressure to produce top results for demanding, paying clients. So it is important

that the steps taken to produce a data mining result are no more time-consuming and expensive than necessary.

The remainder of this paper nominates several research agenda items that may help meet the goals of understanding, accelerating, and automating data mining. In pursuit of our vision of reducing inefficient cycles in the mining process, we suggest three topics: modeling the data mining process, automating data preparation, and mechanically recognizing the inadequacy of models.

3. MODELING OF THE DATA MINING PROCESS

Three possible ways to model the data mining process are: to treat it as (1) a dynamic systems simulation problem; (2) a control problem, suitable for the application of planning methods from Artificial Intelligence(AI), and (3) a control problem, to which learning techniques can be applied.

The first approach views the data mining process as a dynamic process that can be modeled using simulation techniques. Vendors such as Venting Systems provide software for modeling business problems as time-dependent dynamic processes. Simulation has already been applied to the more general problem of modeling project management dynamics including the evaluation of alternative policies to improve performance.

A second way to study the data mining process is to regard it as a planning problem. One standard AI formulation of planning requires (a) an initial state, (b) a goal state and (c) operators. For example, the initial state may correspond to a set of databases that are available for extraction and a formalized description of business problem. The goal state is to achieve some business objectives, such as keeping attrition below 5%. Plan operators correspond to the various data mining actions that can be performed^[3]. The steps may be arranged in a hierarchical plan. The step of applying an algorithm to data, for example, requires several sub-steps, including designation of the target data, the selection of parameters, storing and indexing the results, recording comments about the run, etc.

A third and the related way to represent that data mining process is as a control problem that may be amenable to reinforcement learning and dynamic programming methods. A reinforcement learning problem is characterized by an *agent* that takes *actions*, based on representations of an environment's *state*, and receives a scalar *reward* at each step. Reinforcement learning assumes only that the actions taken by the agent are evaluated, rather than instructing the agent as to the correct response. The possibly noisy and weak supervision provided to a data mining analyst may be more faithfully modeled through the reinforcement learning framework than through classical supervised learning.

There are, nevertheless, hurdles to the applications of reinforcement learning to the control of data mining procedures. For one, reinforcement learning typically has been applied in well-circumscribed domains, e.g., elevator control. Its applications to an area like data mining may be a leap for these techniques, particularly if large amounts of data are to be generated through simulation or otherwise to support the learning of an agent's policy that maps states to the probability of selecting each action.

4. DATA PREPARATION

Anecdotally, it is estimated 70% (plus or minus 20%, say) of the time used in a data mining project is dedicated to data preparation. Data preparation is a lengthy, often tiresome, stage in a mining engagement, and so limiting the iterations of preparation is a particularly attractive goal^[4].

While outlier identification, missing value imputation, discrimination, and other cleansing techniques have received much research attention, other standard preparatory tasks have received less. For example, common is the aggregation or "rolling up" of a customer's transactional data into a single, summary non-transactional record. The transactions of a banking customer may be aggregated so that a single record represents the customer, rather than a series of transactions.

A preliminary question is the level to which the transactional data is aggregated. In a business-to-business setting, a transacting corporation may have one or more physical sites or locations, one or more regional headquarters, one or more affiliated corporations: the appropriate level of representations for a business is not obvious. Situations like this may call for applying a predictive (or other) algorithm at various levels of transactional aggregation and determining the accuracy at each level. In an extreme case, various types of search could be applied to find an appropriate business unit for aggregation^[5]. Exhaustive search, heuristic search, or blind search may be useful to find the level of roll-up that may lead to the greater predictive accuracy, especially if different segments of the training population may be aggregated to different levels.

Another problem that arises in the aggregation of transactional data is the value that is imputed for a variable that summarizes a set of transactional records. For instance, suppose a bank customer makes various automated teller transactions. One may want to use as a derived variable the mean amount of cash withdrawn, a robust mean, or the maximum and the minimum^[6]. Whether to take these (mean, minimum, maximum) or other summaries of the transactional data to the aggregate level is a mining question whose answer may be found through learning or search.

Research into principled ways to choose the right level of aggregation and the most representative aggregated features would help to eliminate inefficient cycles in data preparation.

5. AUTOMATED RECOGNITION OF INEFFICIENT MODELS

One of the drivers of additional knowledge discovery iterations is the recognition that a model is inadequate. To the extent that a model can be automatically recognized as insufficient, the iterations might themselves be automated.

There are several types of models that may be recognized as defective. The first type is trivial model. A classic example is the one-node decision tree. Of course, it is easy computationally to recognize a decision tree with one node. Further, there may be also a set of prototypical mining responses to a trivial model. In the case of a one-node decision tree, for example, one immediate check would confirm that the records are not all assigned to the same class.

At the other end of the spectrum of inadequate models is the model that is trivial because each record is placed in its own partition. For instance, if a unique key is assigned to each record and it is then used as a clustering variable, then one-element clusters may result. Again, this type of over fitting is easy to recognize, and could give rise to automatic responses, such as eliminating a variable that appears with a unique value in too many records.

Other types of model faults may be identified. An empirical way to determine associated responses to inadequacy may be to observe data mining analysts in practice and to discern the heuristics that they use to respond to various model shortcomings.

6. SUMMARY

The objective of this paper has been to suggest several research problems and possible solution paths to make the data mining process less costly by eliminating or automating iterations. We have suggested three areas for additional research: modeling of the mining process, facilitated data preparation, and the mechanical recognition of sub-par models. As the data mining process becomes better understood and more straightforward, it may then be more efficient and more closely integrated into the business processes of organizations.

REFERENCES

- [1] Rakesh Agrawal, Tomasz Imielinski, and Arun Swami, (1993), Mining Association Rules Between Sets of Items in Large Databases, ACM SIGMOD International Conference on Knowledge Discovery and Data Mining, SIGMOD, Proceedings, pp. 207-216.
- [2] Sergey Brin, Rajeev Motwani, Jeffery D Ullman, and Shalom Tsur, (1997), Dynamic Itemset Counting and Implication Rules for Market Basket Data. ACM SIGMOD International Conference on Management of Data, SIGMOD, Proceedings, pp. 255-264.
- [3] Roberto J. Bayardo Jr. and Rakesh Agrawal, (1999), Mining the Most Interesting Rules. ACM SIGMOD International Conference on Knowledge Discovery and Data Mining, SIGMOD, Proceedings, pp. 145-154.

- [4] Tom Brijs, Gilbert Swinnen, Koen Vanhoof, and Geert Wets, (1999), Using Association Rules for Product Assortment Decisions: A Case Study. ACM SIGMOD International Conference on Knowledge Discovery and Data Mining, SIGMOD, Proceedings, pp. 254-260.
- [5] Handley, S., Langley, P., and Rauscher, F. A. (1998). Learning to Predict the Duration of an Automobile Trip. In *Proceedings of 1998 International Conference on KDD and Data Mining (KDD '98)*, pp. 219-223, New York City, USA.
- [6] Agrawal, R. (1999). Data Mining: Crossing the Chasm. Invited Talk at the 5th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '99), San Diego, California.

R. K. Chauhan

Department of Computer Sc. & Applications, Kurukshetra University,
Haryana, INDIA

E-mail: rkckuk@yahoo.com

Abhishek Taneja

Department of Computer Sc. & Applications, DIMT Kurukshetra,
Haryana, INDIA

E-mail: taneja246@yahoo.com