# Incorporating domain knowledge into data mining classifiers: An application in indirect lending

Atish P. Sinha, Huimin Zhao *

*Sheldon B. Lubar School of Business, University of Wisconsin-Milwaukee, P. O. Box 742, Milwaukee, WI 53201-0742, United States*

### A R T I C L E   I N F O

### A B S T R A C T

Data mining techniques have been applied to solve classification problems for a variety of applications such as credit scoring, bankruptcy prediction, insurance underwriting, and management fraud detection. In many of those application domains, there exist human experts whose knowledge could have a bearing on the effectiveness of the classification decision. The lack of research in combining data mining techniques with domain knowledge has prompted researchers to identify the fusion of data mining and knowledge-based expert systems as an important future direction. In this paper, we compare the performance of seven data mining classification methods—naive Bayes, logistic regression, decision tree, decision table, neural network, $k$-nearest neighbor, and support vector machine—with and without incorporating domain knowledge. The application we focus on is in the domain of indirect bank lending. An expert system capturing a lending expert's knowledge of rating a borrower's credit is used in combination with data mining to study if the incorporation of domain knowledge improves classification performance. We use two performance measures: misclassification cost and AUC (area under the curve). A 2×7 factorial, repeated-measures ANOVA, with the two factors being domain knowledge (present or absent) and data mining method (seven methods), as well as a special statistical test for comparing AUCs, is used for analyzing the results. Analysis of the results reveals that incorporation of domain knowledge significantly improves classification performance with respect to both misclassification cost and AUC. There is interaction between classification method and domain knowledge. Incorporation of domain knowledge has a higher influence on performance for some methods than for others. Both measures—misclassification cost and AUC—yield similar results, indicating that the findings of the study are robust.

© 2008 Elsevier B.V. All rights reserved.

## 1. Introduction

Data mining techniques have been applied to solve classification problems for a variety of applications, including credit scoring, bankruptcy prediction, insurance underwriting, and management fraud detection. These techniques automatically induce prediction models, called *classifiers*, based on historical data about previously solved problem cases. The classifiers can then be applied to recommend solutions to new problem cases.

In many of the application domains that have been studied by data mining researchers, there exist human experts who have developed their expertise through years of experience in solving problems in those domains. An expert's knowledge tends to be heuristic in nature. Because experts often find it difficult to articulate the heuristics or rules of thumb that they use to efficiently solve a problem, acquiring their expertise is usually a difficult and challenging task. This phenomenon is commonly referred to as the *knowledge acquisition bottleneck* [20].

A major benefit of using a data mining technique is that it bypasses the knowledge acquisition bottleneck. By unearthing the patterns or knowledge from the data itself, data mining methods obviate the need for eliciting knowledge from a

---

* Corresponding author. Tel.: +1 414 229 6524; fax: +1 414 229 5999.
*E-mail addresses:* sinha@uwm.edu (A.P. Sinha), hzhao@uwm.edu (H. Zhao).

human expert. Clearly, data mining lends itself naturally to domains in which there is a dearth of human expertise or in which domain knowledge cannot be easily formalized. However, there are domains that have large bodies of *domain knowledge* encapsulated in the form of human expertise. Also, in some of those domains, there exist large volumes of data. But very little research has been conducted to examine if domain knowledge can be incorporated into data mining for better performance. As Dybowski et al. [10] (p. 293) stated:

> "At present, knowledge engineering and machine learning remain largely separate disciplines, yet, in many fields of endeavor, substantial human expertise exists alongside data archives. When both data and domain knowledge are available, how can these two resources effectively be combined to construct decision support systems?"

In this paper, we address the question by examining if such a fusion of domain knowledge and data could improve classifier performance in the domain of indirect bank lending. We compare the performance of seven data mining classification methods—naive Bayes, logistic regression, decision tree, decision table, neural network, *k*-nearest neighbor, and support vector machine—with and without incorporating domain knowledge. An expert system capturing a lending expert's knowledge of rating a borrower's credit is used in combination with data mining to study whether the incorporation of domain knowledge improves classification performance. We use two measures—misclassification cost and AUC (area under the curve)—for evaluating classifier performance.

Our study makes an important contribution to existing research in data mining by empirically investigating whether domain expertise improves performance of classifiers built using different methods. Also, instead of using classification accuracy or error rate as the sole performance measure—as has been the norm in fusion research—we evaluate the classifiers with respect to misclassification cost under a range of cost ratios, and with respect to AUC, an aggregate measure.

The paper is organized as follows. Section 2 reviews the related work in the area. Section 3 describes the domain knowledge and Section 4 describes the seven classification methods used in this study. Section 5 presents the theoretical framework and research questions. Section 6 describes the research design and methodology. Section 7 presents the results and Section 8 provides a discussion of the results and their implications. Section 9 summarizes the contributions of this study and outlines directions for future research.

## 2. Background

The process of building an expert system [5] is known as *knowledge engineering*. It involves a knowledge engineer eliciting procedures, strategies, and rules of thumb from a domain expert and transferring this heuristic knowledge into a computer program [48]. The resulting expert system solves problems in much the same way as the expert does, by using shortcuts and tricks and ignoring irrelevant information. It uses the stored knowledge to achieve high performance [18]. Expert systems have been applied to solve different types of problems, including prediction, diagnosis, design, planning, debugging, and control.

Human experts develop their knowledge through years of experience in solving problems in a narrow area. Because experts usually find it difficult to articulate the heuristics or rules of thumb that they apply, knowledge engineers need to overcome what is known as the *knowledge acquisition bottleneck*. According to Johnson [25], "the more competent domain experts become, the less able they are to describe the knowledge they use to solve problems!" Several knowledge acquisition techniques—such as interviews, protocol analysis, observation, and focus groups—are available for facilitating knowledge transfer from experts.

The field of data mining has its origins in statistics and machine learning. Several data mining methods are available for classification problems, including statistical techniques such as naive Bayes, discriminant analysis, and logistic regression, and machine learning techniques such as decision tree/rule induction and neural network. Data mining classifiers have been developed in several application domains, such as bankruptcy prediction [27,44,47], financial performance prediction [31], bond rating analysis [22], credit evaluation [46], credit risk assessment [9], and network intrusion detection [56].

In knowledge engineering, the focus is on the *knowledge* of a human expert in a specific problem area. On the other hand, the focus of data mining is on the *data* available in an organization. As noted before, these two fields have largely remained independent of one another, despite the fact that expert systems and data mining methods could play complementary roles in situations where both knowledge and data are available. Fayyad et al. [14] contended that the use of domain knowledge is important in all stages of the knowledge discovery process.

In one of the earliest studies on the subject, Pazzani and Kibler [34] developed a general purpose relational learning algorithm called FOCL, which combines explanation-based and inductive learning. In a later study, Pazzani et al. [35] conducted an experiment comparing FOCL with a domain theory to FOCL without a domain theory. A partial knowledge base of an expert system was used as the domain theory. They found that incorporating domain theory significantly reduced misclassification costs when larger training sets were used.

Hirsh and Noordewier [19] used background knowledge of molecular biology to re-express data in terms of higher-level features. Using C4.5 decision trees and backprop neural networks on DNA sequence learning tasks, they conducted experiments with and without the higher-level features. For both learning methods, the use of higher-level features resulted in significantly lower error rates.

In another study, Ambrosino and Buchanan [2] examined whether the addition of domain knowledge improved the learning of a rule induction program for predicting the risk of mortality in patients with community-acquired pneumonia. The augmented models performed significantly better (lower percent mean error) than the models without domain knowledge. The method preferred by the subjects for incorporating domain knowledge was addition of new attributes, which were derived from existing attributes.

Data mining techniques are good at generating useful statistics and finding patterns in large volumes of data, but "they are not very smart in *interpreting* these results, which is crucial for turning them into *interesting, understandable and actionable knowledge*" [37]. Pohle [37] viewed the lack of