



Auto insurance fraud detection using unsupervised spectral ranking for anomaly

Ke Nian^{a,1}, Haofan Zhang^{a,1}, Aditya Tayal^{a,1}, Thomas Coleman^{b,2}, Yuying Li^{a,*}

^a Cheriton School of Computer Science, University of Waterloo, Waterloo, ON, N2L 3G1, Canada

^b Combinatorics and Optimization, University of Waterloo, Waterloo, ON, N2L 3G1, Canada

Received 29 February 2016; accepted 2 March 2016

Available online 9 March 2016

Abstract

For many data mining problems, obtaining labels is costly and time consuming, if not practically infeasible. In addition, unlabeled data often includes categorical or ordinal features which, compared with numerical features, can present additional challenges. We propose a new unsupervised spectral ranking method for anomaly (SRA). We illustrate that the spectral optimization in SRA can be viewed as a relaxation of an unsupervised SVM problem. We demonstrate that the first non-principal eigenvector of a Laplacian matrix is linked to a bi-class classification strength measure which can be used to rank anomalies. Using the first non-principal eigenvector of the Laplacian matrix directly, the proposed SRA generates an anomaly ranking either with respect to the majority class or with respect to two main patterns. The choice of the ranking reference can be made based on whether the cardinality of the smaller class (positive or negative) is sufficiently large. Using an auto insurance claim data set but ignoring labels when generating ranking, we show that our proposed SRA significantly surpasses existing outlier-based fraud detection methods. Finally we demonstrate that, while proposed SRA yields good performance for a few similarity measures for the auto insurance claim data, notably ones based on the Hamming distance, choosing appropriate similarity measures for a fraud detection problem remains crucial.

© 2016, China Science Publishing & Media Ltd. Production and hosting by Elsevier on behalf of KeAi Communications Co. Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Keywords: Unsupervised learning; Fraud detection; Rare class ranking; Similarity measure; Kernels; Spectral clustering; One-class svm

1. Introduction

The main objective of the paper is to propose a method for fraud detection by detecting anomaly of interdependence relation, captured by a kernel similarity, among the feature variables. The method includes a new ranking

* Corresponding author.

E-mail addresses: knian@uwaterloo.ca (K. Nian), haofan.zhang@uwaterloo.ca (H. Zhang), amtayal@uwaterloo.ca (A. Tayal), tfc Coleman@uwaterloo.ca (T. Coleman), yuying@uwaterloo.ca (Y. Li).

Peer review under responsibility of China Science Publishing & Media Ltd.

¹ All authors acknowledge funding from the National Sciences and Engineering Research Council of Canada.

² This author acknowledges funding from the Ophelia Lazaridis University Research Chair. The views expressed herein are solely from the authors.

<http://dx.doi.org/10.1016/j.jfds.2016.03.001>

2405-9188/© 2016, China Science Publishing & Media Ltd. Production and hosting by Elsevier on behalf of KeAi Communications Co. Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

scheme, with the top of the ranked list indicating the most suspicious case, based on spectral analysis of the Laplacian of the similarity kernel. To illustrate, the proposed method is applied to both synthetic data sets and an auto insurance application.

Fighting against insurance fraud is a challenging problem both technically and operationally. It is reported that approximately 21%–36% auto-insurance claims contain elements of suspected fraud but only less than 3% of the suspected fraud is prosecuted.^{1,2} Traditionally, insurance fraud detection relies heavily on auditing and expert inspection. Since manually detecting fraud cases is costly and inefficient and fraud need to be detected prior to the claim payment, data mining analytics is increasingly recognized as a key in fighting against fraud. This is due to the fact that data mining and machine learning techniques have the potential to detect suspicious cases in a timely manner, and therefore potentially significantly reduce economic losses, both to the insurers and policy holders. Indeed there is great demand for effective predictive methods which maximize the true positive detection rate, minimize the false positive rate, and are able to quickly identify new and emerging fraud schemes.

Fraud detection can be approached as an *anomaly ranking* problem. Anomaly detection encompasses a large collection of data mining problems such as disease detection, credit card fraud detection, and detection of any new pattern amongst the existing patterns. In addition, comparing to simply providing binary classifications, providing a ranking, which represents the degree of relative abnormality, is advantageous in cost and benefit evaluation analysis, as well as in turning analytic analysis into action.

If anomaly can be treated as a rare class, many methods for supervised rare class ranking exist in the literature. RankSVM³ can be applied to a bi-class rare class prediction problem. However, solving a nonlinear kernel RankSVM problem is computationally prohibitive for large data mining problems. Using SVM ranking loss function, a rare class based nonlinear kernel classification method, RankRC,^{4,5} is proposed.

Unfortunately it may not be feasible or desirable to use supervised anomaly ranking for fraud detection, since obtaining clearly fraudulent (and non-fraudulent) labels is very costly, if not impossible. Even if one ignores human investigative costs, it is quite common to find fraud investigators to differ in their claim assessments. This raises additional reliability issues in data (specifically in labels). In contrast, unsupervised learning has the advantages of being more economical and efficient application of knowledge discovery. Moreover, the need to detect fraudulent claims before payments are made and to quickly identify new fraud schemes essentially rule out supervised learning as a candidate solution to effective fraud detection in practice. Therefore, an unsupervised anomaly ranking is more appropriate and beneficial here.

Standard unsupervised anomaly detection methods include clustering analysis and outlier detection. Many outlier detection methods have been proposed in the literature. Examples include *k*-Nearest Neighbor (*k*-NN) outlier detection, one-class Support Vector Machine (OC-SVM) (including kernel-based), and density-based methods, e.g., Local Outlier Factor (LOF). The effectiveness of these methods has been investigated in numerous application domains, including network intrusion detection, credit card fraud detection, and abnormal activity detection in electronic commerce. However, many standard outlier detection methods, e.g., one-class SVM, are only suitable for detecting outliers with respect to a single global cluster, which we refer to as global outliers in this paper. An implicit assumption in this case is that normal cases are generated from one mechanism and abnormal cases are generated from other mechanisms. Density-based methods can be effective in detecting both global outliers and local outliers; but assumption here is that data density is the only discriminant for abnormality. Density-based methods fail when small dense clusters also constitute abnormality. In addition, density based methods, e.g., LOF, often require users to define a parameter which specifies a neighborhood to compare the density. Tuning these parameters can often be challenging.

Understandably, unsupervised learning is much more difficult than supervised learning, since learning targets are not available to guide the learning process. In practice, difficulty in unsupervised learning is further exacerbated by the additional challenge in identifying relevant features for unsupervised learning methods. Due to these challenges, the existing literature on auto insurance fraud detection typically formulates the problem as a supervised learning problem.^{6,7}

The literature on *unsupervised* auto insurance fraud detection is extremely sparse. To the best of our knowledge, it includes the self-organizing feature map method,⁸ and PRIDIT analysis.^{9–11} PRIDIT is based on RIDIT scores and Principal Component Analysis. We note these studies^{8–11} have been conducted using the single Personal Injury Protection (PIP) data set,⁹ which is provided by Automobile Insurance Bureau (AIB) from Massachusetts. This data set has been preprocessed by auditors and fraud detection inspectors. Specifically, data features are the red flags specified by domain experts, with attribute values for fraudulent instances in the data set typically smaller than that of

Download English Version:

<https://daneshyari.com/en/article/1002109>

Download Persian Version:

<https://daneshyari.com/article/1002109>

[Daneshyari.com](https://daneshyari.com)