



An application of supervised and unsupervised learning approaches to telecommunications fraud detection [☆]

Constantinos S. Hilas*, Paris As. Mastorocostas

Department of Informatics and Communications, Technological Educational Institute of Serres, Terma Magnisias, GR-62124 Serres, Greece

ARTICLE INFO

Article history:

Received 30 July 2007

Accepted 24 March 2008

Available online 31 March 2008

Keywords:

Fraud detection

Telecommunications

User profiling

Supervised learning

Unsupervised learning

ABSTRACT

This paper investigates the usefulness of applying different learning approaches to a problem of telecommunications fraud detection. Five different user models are compared by means of both supervised and unsupervised learning techniques, namely the multilayer perceptron and the hierarchical agglomerative clustering. One aim of the study is to identify the user model that best identifies fraud cases. The second task is to explore different views of the same problem and see what can be learned from the application of each different technique. All data come from real defrauded user accounts in a telecommunications network. The models are compared in terms of their performances. Each technique's outcome is evaluated with appropriate measures.

© 2008 Elsevier B.V. All rights reserved.

1. Introduction

Telecommunications fraud can be simply described as any activity by which telecommunications service is obtained without intention of paying [10]. Telecommunications fraud has certain characteristics that make it particularly attractive to fraudsters. The main one is that the danger of localization is small. This is because all actions are performed from a distance which in conjunction with the mesh topology and the size of networks makes the process of localization time-consuming and expensive. Additionally, no particularly sophisticated equipment is needed, if one is needed at all. The simple knowledge of an access code, which can be acquired even with methods of social engineering, makes the implementation of fraud feasible. Finally, the product of telecommunications fraud, a phone call, is directly convertible to money [16].

Several categories of telecommunications fraud have been reported. The main are the technical fraud, the contractual fraud, the hacking fraud, and the procedural fraud [10]. In [1] 12 distinct fraud types are identified while combinations of them have also been reported [13]. The most common fraud scenario in private networks is the superimposed fraud. This is the case of an employ-

ee, the fraudster, who uses another employee's authorization code to access outgoing trunks and costly services. Thus, the fraudster's activity is superimposed over the legitimate user's one.

Telecommunications fraud has drawn the attention of many researchers in recent years not only due to the huge economic burden on companies' accountings but also due to the interesting aspect of user behavior characterization. Fraud detection techniques involve the monitoring of users' behavior in order to identify deviations from some expected or normal norm. Research in telecommunications fraud detection is mainly motivated by fraudulent activities in mobile technologies [1,4,10,20,24,31]. The techniques used come from the area of statistical modeling like rule discovery [2,7,24,30], clustering [3,27], Bayesian rules [4], visualization methods [5], Markov models [31] or neural network classification [14,21,24,32]. Combinations of more than one method have also been proposed [17,28,31]. In [8] one can find a bibliography on the use of data mining and machine learning methods for automatic fraud detection. The site is updated up to November 2004. Most of the aforementioned approaches use a combination of legitimate user behavior examples and some fraud examples. The aim is to detect any usage changes in the legitimate user's history.

The industry's interest in fraud detection problems is also stressed by the high number of relevant patents. A quick search with the keywords "fraud detection" in an online search-engine, [9], on July 2007, revealed 76 patents, 22 of them being relevant to telecommunications.

In general, all fraud cases can actually be viewed as fraud scenarios which are related to the way the access to the network was acquired. Detection techniques tailored to one case may fail to detect other types of fraud. For example, velocity traps which can identify

[☆] The authors thank the staff of the Telecommunications Center of the Aristotle University of Thessaloniki for their contribution of data. This work was supported in part by the Research Committee of the Technological Educational Institute of Serres, Greece.

* Corresponding author. Tel.: +30 23210 49185; fax: +30 23210 49128.

E-mail addresses: chilas@teiser.gr (C.S. Hilas), mast@teiser.gr (P.As. Mastorocostas).

the use of a cloned cell phone will fail to detect a case of contractual fraud. So, fraud detection focuses on the analysis of users' activity. The related approaches are divided into two main subcategories. The absolute analysis that searches for thresholds between legal and fraudulent behavior, and the differential approach that tries to detect extreme changes in a user's behavior. In both cases, analysis is achieved by means of statistical and probabilistic methods, neural networks and rule based systems. However, the use of indicators of excessive usage has been criticized as they may not only imply fraud but may also point to the best customers [25].

In the present paper, we are interested in the different lessons than can be learned from the application of different learning algorithms on different user behavior representations (profiles). Both supervised and unsupervised learning methods are applied. One would expect the findings of one method to be used as inputs to the other one, e.g. first use the unsupervised method and then apply the supervised one in order to boost the learning process. However, this is not the case in the present work. Each method is applied independently from the other and is expected to reveal different aspects of the modeling approach. The main task is to cross-check the effectiveness of different user profiles to discriminate between legitimate and fraudulent activity and additionally identify the elements that are important in the learning process and compare the conclusions from the application of the two methods.

The paper proceeds as follows. In the next section, the data that were used are described along with the user modeling approach. In Section 3 the experimental procedure, i.e., the learning methods, is presented. The experimental results are given in Section 4. In the last section conclusions are drawn.

2. Models of user behavior

Traditionally, in computer security, user modeling is achieved by means of appropriate user profiles. The main idea behind a user's profile is that the user's past behavior can be accumulated. Profiles are constructed based on any basic usage characteristic such as resources consumed, login location, typing rate and counts of particular commands. In telecommunications, user profiles can be constructed from appropriate usage characteristics. The aim is to distinguish a normal user from a fraudster. The latter is, in most cases, a user of the system who knows and mimics normal user behavior. The data that can be used to monitor the usage of a telecommunications network are contained in the call detail record (CDR) of any private branch exchange (PBX). The CDR contains data such as: the caller ID, the chargeable duration of the call, the called

party ID, the date and the time of the call, etc [15]. In mobile telephone systems, such as GSM, the data records that contain details of every mobile phone attempt are the Toll Tickets.

Our experiments are based on real data extracted from a database that holds the CDRs from an organization's PBX for a period of eight years. According to the organization's charging policy, only calls to national, international and mobile destinations are charged. Calls to local destinations are not charged so they are not included in the examples. In order to properly charge users, for the calls they place, a system of Authorization Codes is used. Each user owns a unique authorization code which enables telephone sets to "unlock" and access outgoing trunks. If anyone (e.g. a fraudster) finds a code he can use it to place his own calls from any telephone set within the organization.

Several user accounts, which have been defrauded, have been identified. All contain both examples of legitimate and fraudulent activity. The detailed daily accounts were examined by a field expert and each phone call was marked as either normal or defrauded.

Each profile from each user was labeled according to two different ways. The first one was to identify the fraudster's first day of activity. Then each user's account was split into two sets, one pre-fraud (example of legitimate use) and one post-fraud (example of fraudulent use). Pre- and post- are relevant to the first day that the fraudulent activity appeared. The other labeling approach was more detailed. If no fraudulent activity was present during a day, then the whole day was marked as normal. If at least one call from the fraudster was present then the whole day was marked as fraud.

For each user, three different profile types are constructed. The first one (Profile1) is build up from the accumulated weekly behavior of the user. The profile consists of seven fields which are the mean and the standard deviation of the number of calls per week (calls), the mean and the standard deviation of the duration (dur) of calls per week, the maximum number of calls, the maximum duration of one call and the maximum cost of one call (Fig. 1). All maxima are computed within a week's period.

The second profile (Profile2) is a detailed daily behavior of a user which is constructed by separating the number of calls per day and their corresponding duration per day according to the called destination, i.e., national (nat), international (int), and mobile (mob) calls, and the time of the day, i.e., working hours (w), afternoon hours (a), and night (n) (Fig. 2).

Last, the third profile (Profile3) is an accumulated per day behavior (Fig. 3). It consists of the number of calls and their corresponding duration separated only according to the called destination, that is, national, international and mobile calls.

mean(calls)	std(calls)	mean(dur)	std(dur)	max(calls)	max(dur)	max(cost)
-------------	------------	-----------	----------	------------	----------	-----------

Fig. 1. Profile1 of telephone calls.

nat_calls_w	nat_dur_w	nat_calls_a	nat_dur_a	nat_calls_n	nat_dur_n
mob_calls_w	mob_dur_w	mob_calls_a	mob_dur_a	mob_calls_n	mob_dur_n
int_calls_w	int_dur_w	int_calls_a	int_dur_a	int_calls_n	int_dur_n

Fig. 2. Profile2 of telephone calls.

nat_calls	nat_dur	mob_calls	mob_dur	int_calls	int_dur
-----------	---------	-----------	---------	-----------	---------

Fig. 3. Profile3 of telephone calls.

Download English Version:

<https://daneshyari.com/en/article/405391>

Download Persian Version:

<https://daneshyari.com/article/405391>

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