

Mining Fraudsters and Fraudulent Strategies in Large-Scale Mobile Social Networks

S.Deepak¹, D.Bharath², R. Ram Reddy³, R. Bharath kumar⁴

^{1,2,3}Scholar, Department of CSE in Jayamukhi institute of technological sciences, Narsampet, Warangal, Telangana, India

⁴Assistant professor, Department of CSE in Jayamukhi institute of technological sciences, Narsampet, Warangal, Telangana, India

ABSTRACT

The fast development of contemporary communications technologies, in particular communications through mobile phones, has significantly eased human social connections and the flow of information. However, the advent of telemarketing scams may greatly waste individual fortune and society wealth, which can result in a possible slowdown or harm to economic conditions. In this study, we offer a method for detecting fraudulent telemarketing practises, with a primary focus on illuminating the phenomena of "precise fraud" and the tactics that are used by con artists to accurately pick targets for their scams. We utilise a one-month comprehensive dataset of telecommunication information in Shanghai with 54 million users and 698 million call records in order to conduct this research. The dataset covers a period of one month. As a result of our research, we have discovered that the information of users may well have been badly leaked, but also that fraudster have a preference over the age of the target user as well as their activity on mobile networks. In addition, we provide an innovative semi-supervised learning approach for the purpose of differentiating fraudsters from other types of users. Our method exceeds numerous other state-of-the-art algorithms in terms of accuracy when it comes to identifying fraudsters, as shown by experimental findings on data taken from the real world (for example, +0.278 in terms of F1 on average). We think that the findings of our research have the potential to impact policy decisions made by both the government and companies that supply mobile services.

1 INTRODUCTION

The advent of technology that enables worldwide connection in recent years has led to a significant increase in the frequency of fraudulent actions. The terrible effects of scams are felt by millions of individuals.

For example, in China, the issue of phone fraud has been recognised as being of great importance. According to projections made by both Qihoo1 and Tencent2, there were more over 500 million instances of phone fraud in 2016, which resulted in losses of around 16.4 billion USD. In the meanwhile, less than three percent of these cases are closed successfully.

On August 29th, 2016, it was revealed that a college professor in Beijing had lost \$2.67 million USD to a phone fraudster who represented himself to be a court official. The fraudster made the allegation over the phone.

The effects of phone scams have been far more devastating, and in some cases even life-threatening, than the financial impact they have had on individual victims.

The prevention of fraud is receiving a significant amount of attention. However, due to the limited availability of data and the high level of sensitivity, this field of study has remained largely unexplored by academics. The vast majority of the currently available research on fraud detection [1, 2, 3, 4] builds experiments using either simulated data or real information with a restricted scope. In this article, we focus on a large-scale mobile social network that exists in the real world. Our research encompasses a full set of call records in Shanghai and spans a period of 30 days, beginning September 1 and ending September 30, 2016. The beginning and ending times of each discussion are noted, in addition to the

anonymous phone numbers that are logged for each call log entry. Annotations of scammers that were produced by the audience are also obtained by us. However, there are still a lot of additional obstacles to overcome. The first difficulty arises from the sensitive nature of the data, which prevents us from accessing the information included in each call log's content. Monitoring the content of conversations for certain subjects, such as money transfers, would make it much simpler to identify those who commit fraud. We are compelled to rely on meta information in order to draw conclusions about the inference since we do not have access to the content information.

How is it that intelligent individual, such as the college professor in the example given above, may be taken advantage of? Our research demonstrates that there is a significant possibility that users' information has been compromised, and that fraudsters do not choose their victims at random but rather follow a predetermined plan (See details in Section 3). The second difficulty is to figure out how to expose fraudulent approach in order to have a better understanding of fraud.

The third obstacle is the asymmetry in the labelling system.

According to the findings of our research, more than 95.2 percent of users do not engage in fraudulent activity. Although the imbalance issue has been addressed by systems that identify fraudulent activity on credit cards [5] and insurance policies [6], to the best of the researcher, it has not been thoroughly researched in the context of telecommunications.

In order to solve the first and second challenges, we have designed and constructed a number of exploratory analyses that we have run on our actual mobile network in order to examine the behaviour patterns of fraudsters. Following the results of our trials, we discuss a number of fraudulent tactics. For instance, we see that con artists have a penchant for young individuals, particularly those who are engaged in phone interactions on a regular basis. We have also found that it is best for us

to hang up the bogus phone call quickly rather than spending more time slugging off from the fraudster in order to prevent getting additional fraudulent phone calls. This is because we are trying to avoid receiving more fraudulent phone calls.

In light of these findings, we devise an innovative factorgraph-based model that we call FFD in order to identify dishonest individuals. To be more explicit, our approach takes into account the structural knowledge and inclination of fraudsters when it comes to selecting targets. We go one step further and provide a framework for semi-supervised learning that makes use of both known and unknown labels in order to tackle the problem of label sparsity. According to the results of our tests, we can observe that our model achieves an improvement on F1 of 0.278 compared to a few other approaches that are considered to be state-of-the-art.

It is important to bring attention to our efforts in the following ways:

- Using data from actual phone conversations, we reveal how those who commit fraud as opposed to those who don't commit fraud behave differently on mobile networks.
- We are doing research into the "exact fraudulent technique" and are making an appeal to all parties involved to ensure that the security of personal information has been given top priority.
- We provide an innovative framework that may differentiate mobile network fraudsters from other users of a specific mobile network.
- We demonstrate the viability of our approach by testing it on a large-scale phone service inside the physical world.

II DATA AND PROBLEM

The numbers and statistics. We make use of a mobile dataset that is comprised of comprehensive records of users' telecommunications with one another in Shanghai over the course of one month, beginning September 1, 2016, and ending September 30, 2016. (1 month). The

information was supplied by China Telecom, one of the most prominent mobile service providers in China, and includes 698,811,827 call records that were exchanged between 54,169,476 individuals.

Each entry in the call log includes the time the call began, the time it ended, the number of the caller, and the number of the callee. We are also able to acquire various personal features of each user, including age, sex, and birthplace, since personal identification is necessary to obtain a mobile number in China. In addition, we are able to retrieve this information. China Telecom took steps to protect users' privacy by concealing their identities in our dataset.

We do not divulge any personally identifying information about any of the people in this dataset and instead solely report on the general statistics from this dataset throughout this whole study. It is important to be aware that it is unusual for a single person in China to have more than one phone number. This is due to the fact that it is not easy to get a phone number in China. Because of this, we consider each individual phone number to correspond to a separate user.

Data labelling. After that, we explain how we get our hands on the ground truth data. In general, we get the label data from Baidu3 and Qihoo 3604, both of which have set up report telephone to collect fraudulent callers' phone numbers. To be more explicit, when we have a user and her telephone number, we utilise the application programming interfaces (APIs) of both Baidu and Qihoo 360 to determine whether or not the person has ever been reported as a fraudster. If the person has been reported as a fraudster to Baidu or Qihoo 360 by another user, we will classify the user as a fraudster. These facts on the ground truth originate from a significant number of complaints and, as a result, a very high degree of trust may be assigned to them. By doing this, we are able to acquire annotations of a total of 340,550 people, of whom there are 15,660 fraudulent individuals (around 4.6 percent). The next step is to outline the issue of fraudster mining in order to uncover further

fraudsters by automated means, expand the artefact annotations, and so on.

The framing of the issue. We build a mobile communication network by using the call records that are included inside our dataset. Formally, we construct a directed graph denoted by the equation $G = (V, E)$, where V represents the collection of users, and each directed edge $e_{ij} \in E$ denotes that the user v_i calls v_j ($v_i, v_j \in V$). Every user in V has been assigned a label in the range $y_i \in Y$, which indicates whether or not she is a fraudster ($y_i = 1$), a regular user ($y_i = 0$), or if we do not yet know what her identity is ($y_i = ?$).

During the process of building the network, we came across an unusual phenomenon: scammers phone the number "200" far more often than regular people do (Figure 1). More than seventieth of all fraudulent phone calls are routed via the "200" area code. According to the findings of a research, the number "200" is a transit number that con artists use to hide their real telephone numbers and save money on telephone fees.

According to this finding, for two calls made at the same time, one is from user A to "200," and another is from "200" to another user B, we integrate them as a unified call log from A to B. The first call is made from user A to "200," and the second call is made from "200" to another user B. The following is a formulation of our problem:

Definition 1. Fraudster mining. Given a mobile communication network $G = (V, E)$, and an identity vector $Y = \{Y_L, Y_U\}$ with missing values, where Y_L denotes labelled identity information of users in G and Y_U stands for unknown identities, our objective is to infer the missing values in Y , i.e., to identify fraudsters that may be hiding among other users in the network.

III PROPOSED SYSTEM

Methods that are based on classifiers. The detection of fraud is presented as a problem of binary classification in a significant amount of published material. In other words, given a list of phone numbers, determine if each number is

active or not is natural or bogus. Take, for instance, Weatherford et al. [18] in this regard. Make advantage of user accounts that may retain information for an extended period of time and train neural networks to distinguish between fraudulent conduct and normal behaviour .typical one. Yus off [4] proposes an alternative that uses neural networks. a model that uses the Gaussian mixed model (GMM) as the underlying statistical classifier. Dominik makes use of a classification algorithm of the threshold variety [19].

The most significant disadvantage of classifier-based approaches is the fact that annotations have a significant impact on its performance, and will suffer as a result of the lack of information on the label. Within the scope of this study, we present a framework for semi-supervised learning to advance Improve the performance by using labels that are unknown to you. Methods that are based on graphs. The primary focus of these types of techniques is on discover scammers by recognising surprisingly dense zone related to a network. For instance, Hooi and colleagues [1] concentrate on locating.con artists who operate in the presence of disguises and who put forth the an algorithm named FRAUDAR. Tseng et al. [2] construct a network by using weighted edges to reflect the length of each call between each pair of user and the frequency of calls. After that, they carry out a weighted HITS. algorithm 11 on the network to figure out the trustworthiness of a certain node. phone number and identify potentially fraudulent phone calls based on the to the reliability of the value. Identical concepts are mandated in a number of the other works already in existence [20], [21], [22], [23], [24]. Other pieces of work use outlier detection methods to discover anomalous user profiles [25].The probabilistic model is the foundation for the approach that has been suggested. graphical model, which has also been used in fraud prevention and detection. The identification of reviews (26), the extraction of social events (27), and the signal processing [28], among other things A great number of graph-based approaches only take into account a select few categories. pertaining to

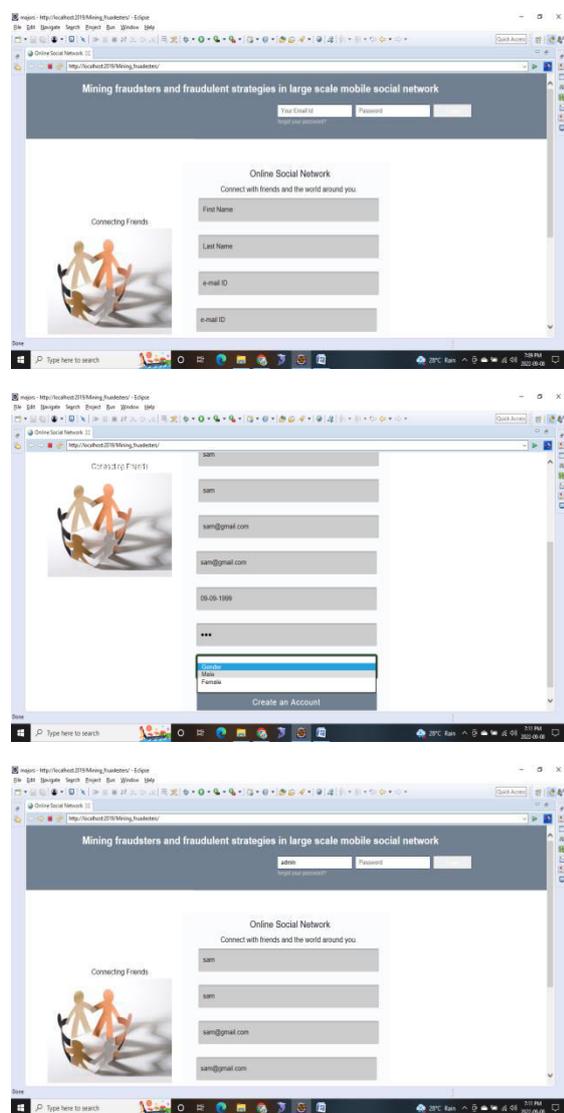
characteristics, the majority of which are call frequency and call duration. In the course of this effort, by analysing and differentiating con artists. We suggest a number of broad and efficient characteristics, all of which come from honest people. We are of the opinion that others may profit from our features.

Investigate and combat fraudulent activity. One more distinction between us and them work and those of others is the fact that, to the best of our knowledge, we are not aware of any are the ones who have initially begun researching fraudulent approach. Approaches that are based on decomposition. These kind of methods and strategies. Applying matrix decomposition will allow you to identify con artists. For instance, Akoglu and Faloutsos [29] describe an approach that is based on SVD decomposition and may discover anomalous nodes in a network. a graph that changes over time, each node of which is believed to be a pattern is considered an anomaly if it deviates considerably from the norm. pattern that came before. Ide and his colleague present a study that is comparable to this one. Kashima [30], both of which concentrate on the issue of monitoring a web-based system with several tiers. Sun et al. [31] suggest a methodology for the identification of anomalies in dynamic graphs, which the low-rank approximations are used as summaries of the sparse data.graph.

After that, the reconstruction error is used in the evaluation of the. degree of oddness or peculiarity. Rossi et al. [32] create a none negative iteratively extract values using a technique based on matrix factorization ode feature and identify the roles played by the nodes. Various other works like [33], [34], and [35] also employ a decomposition-based comparable to the one described. approaches to handle the fraud detection challenge. Outlier detection. Our work is also important to research on finding areas that deviate from the norm in terms of their local structure or characteristics diverge a great deal from the norms and practises of other members of a social group. For example, Gao et al. [36] and Perozzi et al. [37] suggest techniques that may identify communities and outliers at the same

time. A model that makes advantage of an outlier is presented by Muller et al. [38].strategy for rating the attributes in the graph. There are also a lot of them. works that centre on dynamic graphs as its primary theme. Take, for instance, Peel. Clauset [39] and utilise an extended version of the hierarchical random graphs (GHRG) to represent the architecture of communities inside a graph.[40] Sun and colleagues suggest using a mechanism they call Graph Scope, which is a kind of algorithm that does not need any parameters, and it monitors the Over time, the network's nodes will begin to divide.

IV RESULTS



V CONCLUSION

In this article, we investigate the challenge of mining fraudsters and fraudulent schemes in a massive mobile network by using data mining techniques.

By analysing a one-month sufficient amount of data of telecommunication metadata in Shanghai with 698 million call logs between 54 million users, we came to the conclusion that people who engage in fraudulent activity communicate with others in a manner that is distinct from that of people who do not engage in fraudulent activity. In addition, con artists chose their victims based on criteria such as the users' age and the amount of activity they have in their phone conversations. After doing an exploratory investigation, we come up with an original semi-supervised model that can differentiate between those who commit fraud and those who do not commit fraud. The findings of our experiments indicate that our model is able to significantly outperform a number of the baseline approaches that are currently considered to be state of the art.

Concerning the work that will be done in the future, it is interesting to consider how to identify a fraud group rather than a single fraudster, which is comprised of several fraudsters who each have their own distinct tasks and responsibilities. On the basis of this, the patterns of cooperation used by various fraud organisations might be uncovered. In addition, our work might be expanded by taking into further consideration the geographical information of users and by researching the offline activities of fraudsters, such as how they travel about the city.

The amount of data to which we have access restricts the scope of our job.

Even though China Telecom is a large service provider and Shanghai is an important global metropolis, the selection bias in our data may restrict the generalizability of our study. Shanghai is indeed an important worldwide city.

REFERENCES

- [1] B. Hooi, H. A. Song, A. Beutel, N. Shah, K. Shin, and C. Faloutsos, "Fraudar: Bounding graph fraud in the face of camouflage," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 895–904.
- [2] V. S. Tseng, J. Ying, C. Huang, Y. Kao, and K. Chen, "Fraudetector: A graph-mining-based framework for fraudulent phone call detection," in Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015, pp. 2157–2166.
- [3] J. Xu, A. H. Sung, and Q. Liu, "Behaviour mining for fraud detection." *Journal of Research and Practice in Information Technology*, 2007.
- [4] M. I. M. Yusoff, I. Mohamed, and M. R. A. Bakar, "Fraud detection in telecommunication industry using gaussian mixed model," in Research and Innovation in Information Systems (ICRIIS), 2013 International Conference on, 2013, pp. 27–32.
- [5] P. Chan, W. Fan, A. L. Prodromidis, and S. J. Stolfo, "Distributed data mining in credit card fraud detection," *IEEE Intelligent Systems & Their Applications*, vol. 14, no. 6, pp. 67–74, 1999.
- [6] T. Ormerod, N. Morley, L. Ball, C. Langley, and C. Spenser, "Using ethnography to design a mass detection tool (mdt) for the early discovery of insurance fraud," in CHI'03 Extended Abstracts on Human Factors in Computing Systems, 2003, pp. 650–651.
- [7] S. Aral, L. Muchnik, and A. Sundararajan, "Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks," *Proceedings of the National Academy of Sciences*, vol. 106, no. 51, pp. 21 544–21 549, 2009.
- [8] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," *Information Theory, IEEE Transactions on*, vol. 47, no. 2, pp. 498–519, 2001.
- [9] J. M. Hammersley and P. Clifford, "Markov fields on finite graphs and lattices," Unpublished manuscript, 1971. [10] K. P. Murphy, Y. Weiss, and M. I. Jordan, "Loopy belief propagation for approximate inference: An empirical study," in *UAI'99*, 1999, pp. 467–475.
- [11] J. Kleinberg, "Hubs, authorities, and communities," *ACM Computing Surveys*, vol. 31, p. 5, 1999.
- [12] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [13] A. Clauset, M. E. J. Newman, and C. Moore, "Finding community structure in very large networks," *Phys. Rev. E*, vol. 70, no. 06111, 2004.
- [14] R. A. Becker, C. Volinsky, and A. R. Wilks, "Fraud detection in telecommunications: History and lessons learned," *Technometrics*, vol. 52, no. 1, pp. 20–33, 2010.
- [15] P. Picard, "Economic analysis of insurance fraud," in *Handbook of insurance*, 2000, pp. 315–362.
- [16] S. B. E. Raj and A. A. Portia, "Analysis on credit card fraud detection methods," in *Computer, Communication and Electrical Technology (ICCCET)*, 2011 International Conference on, 2011, pp. 152–156.
- [17] J. Li, K.-Y. Huang, J. Jin, and J. Shi, "A survey on statistical methods for health care fraud detection," *Health care management science*, vol. 11, no. 3, pp. 275–287, 2008.
- [18] M. Weatherford, "Mining for fraud," *IEEE Intelligent Systems*, vol. 17, no. 4, pp. 4–6, 2002.
- [19] D. Olszewski, "A probabilistic approach to fraud detection in telecommunications," *Knowledge-Based Systems*, vol. 26, pp. 246–258, 2012.

- [20] Q. Cao, M. Sirivianos, X. Yang, and T. Pregueiro, "Aiding the detection of fake accounts in large scale social online services," in Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation, 2012, pp. 15–15.
- [21] S. Ghosh, B. Viswanath, F. Kooti, N. K. Sharma, G. Korlam, F. Benevenuto, N. Ganguly, and K. P. Gummadi, "Understanding and combating link farming in the twitter social network," in Proceedings of the 21st international conference on World Wide Web, 2012, pp. 61–70.
- [22] Z. Gyongyi, H. Garcia-Molina, and J. Pedersen, "Combating web spam with trustrank," in Proceedings of the Thirtieth international conference on Very large data bases-Volume 30. VLDB Endowment, 2004, pp. 576–587.
- [23] M. Jiang, P. Cui, A. Beutel, C. Faloutsos, and S. Yang, "Catchsync: catching synchronized behavior in large directed graphs," in Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, 2014, pp. 941–950.
- [24] B. Wu, V. Goel, and B. D. Davison, "Propagating trust and distrust to demote web spam." *MTW*, vol. 190, 2006.
- [25] M. Onderwater, "Detecting unusual user profiles with outlier detection techniques," VU University Amsterdam, 2010.
- [26] L. Akoglu, R. Chandy, and C. Faloutsos, "Opinion fraud detection in online reviews by network effects," in *ICWSM*, 2013.
- [27] X. Liu, X. Zhou, Z. Fu, F. Wei, and M. Zhou, "Exacting social events for tweets using a factor graph," in *Twenty-Sixth AAAI Conference on Artificial Intelligence*, 2012.
- [28] H.-A. Loeliger, J. Dauwels, J. Hu, S. Korl, L. Ping, and F. R. Kschischang, "The factor graph approach to model-based signal processing," *Proceedings of the IEEE*, vol. 95, no. 6, pp. 1295–1322, 2007. [29] L. Akoglu and C. Faloutsos, "Event detection in time series of mobile communication graphs," in *Army science conference*, 2010, pp. 77–79.
- [30] T. Ide and H. Kashima, "Eigenspace-based anomaly detection in computer systems," in Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2004, pp. 440–449.
- [31] J. Sun, Y. Xie, H. Zhang, and C. Faloutsos, "Less is more: Sparse graph mining with compact matrix decomposition," *Statistical Analysis and Data Mining: The ASA Data Science Journal*, vol. 1, no. 1, pp. 6–22, 2008.
- [32] R. A. Rossi, B. Gallagher, J. Neville, and K. Henderson, "Modeling dynamic behavior in large evolving graphs," in Proceedings of the sixth ACM international conference on Web search and data mining. ACM, 2013, pp. 667–676.
- [33] H. Bunke, P. J. Dickinson, M. Kraetzl, and W. D. Wallis, *A graphtheoretic approach to enterprise network dynamics*. Springer Science & Business Media, 2007, vol. 24.
- [34] M. A. Peabody, "Finding groups of graphs in databases," Ph.D. dissertation, Drexel University, 2002.
- [35] P. Shoubridge, M. Kraetzl, W. Wallis, and H. Bunke, "Detection of abnormal change in a time series of graphs," *Journal of Interconnection Networks*, vol. 3, no. 01n02, pp. 85–101, 2002.