# Heterogeneous anomaly detection in social diffusion with discriminative feature discovery

Siyuan Liu[a], Qiang Qu[b,c,*], Shuhui Wang[d]

[a] *Pennsylvania State University, USA*
[b] *Shenzhen Institutes of Advanced Technology, CAS, China*
[c] *MOE Key Laboratory of Machine Perception, Peking University, China*
[d] *Institute of Computing Technology, CAS, China*

## ABSTRACT

Social diffusion is a dynamic process of information propagation within social networks. In this paper, we study social diffusion from the perspective of *discriminative features*, a set of features differentiating the behaviors of social network users. We propose a new parameter-free framework based on modeling and interpreting of discriminative features that we have created, named HADISD. It utilizes a probability-distribution-based parameter-free method to identify the maximum vertex set with specified features. Using the maximum vertext set, a probability-distribution-based optimization approach is applied to find the minimum number of vertices in each feature category with the maximum discriminative information. HADISD includes an incremental algorithm to update the discriminative vertex set over time. The proposed model is capable of addressing anomaly detection in social diffusion, and the results can be leveraged for both spammer detection and influence maximization. The findings from our extensive experiments on four real-life datasets show the efficiency and effectiveness of the proposed scheme.

© 2018 Elsevier Inc. All rights reserved.

## 1. Introduction

The creation of massive amounts of data from both online (e.g., Facebook and Twitter) and offline (e.g., mobile communication networks) social networks is offering an unprecedented opportunity in studying human social behaviors. *Social diffusion* is among the most interesting and well-motivated problems, describing a dynamic process of information propagation in social networks [5,23,32,51,52]. Many previous studies propose to solve spammer detection and influence maximization based on social diffusion, which are of critical importance to a wide range of applications. Spammer detection usually utilizes any of the following: topology-based [50], content-classification-based [8], and social-behavior-based methods [35,53], while influence maximization often focuses on influence model [9,14,24], influence spread optimization [25], influence estimation [15], and retrieval efficiency speed-up [12,45]. These two problems are mostly evaluated independently and thus use application-specific solutions. In this paper, we study social diffusion from a new perspective based on "*discriminative features*" to find anomalies, and the results can be leveraged as a unified solution to both problems.

A dynamic social graph can be considered to have three components: Vertices, edges, and timestamps. At each timestamp, a vertex spreads information to a set of adjacent vertices via the edges. A large amount of vertices spreading informa-

---

tion within a sufficiently long period yields a social diffusion process. In the process, if a vertex diffuses an unsolicited bulk of messages (spam) to other vertices or creates fake vertices to steal private information from others, we call the vertex a spammer with spamming behaviors. The influence of a vertex in a process relies on its reachable neighbors, the influencing behaviors are determined by the network structures over the time. Influence maximization is the selection of vertices to achieve the maximum influencing behaviors on a network. Therefore, we can postulate: 1) Spamming and influencing behaviors are linked with diffusion processes, 2) both are heterogeneous with various social behaviors, and 3) if we consider spamming behaviors and maximum influencing behaviors as "abnormal" behaviors, both problems can be seen as anomaly detection in diffusion processes.

Recognizing and utilizing discriminative features are important to conduct heterogeneous anomaly detection in social diffusion processes. In this paper, we consider discriminative features as a set of features differentiating the target behaviors from the rest. Note that an anomaly detected from heterogeneous semantics and behaviors is called heterogeneous anomaly. To address the problem, there exist three challenges:

- **Behavior variety**. Since there may exist huge diversities with various features for both social activities and behaviors, anomalies can be difficult to be detected. For instance, spammers often create evasive accounts and behave like "real people". And these behaviors are often changing with various features.
- **Feature heterogeneity**. Social behavior features may be divergent and behavior-dependent. The selection of the most discriminative feature subset to well distinguish a behavior is thus hard. In online social networks, spammers often behave differently resulting the hardness to precisely define them. In offline social networks, people may also have various behaviors in different diffusion processes to influence others. Thus, the process of selecting a set of initial seeds to achieve the maximum influence is challenging. Hereby, in this paper, heterogeneity shall not only be considered as diverse data features and social structures, but also various diffusion processes.
- **Anomaly dynamics**. Abnormal social behaviors are dynamic in social networks. Finding the behaviors thus requires effective algorithms. To avoid being detected, spammers may constantly change their behavior patterns. To maximize influence for different set of vertices, we have to situate influence maximization in the context of dynamics not only in vertex features but also in the ways of influence.

Social diffusion related problems are challenging as they are hard to define and volatile to investigate. Despite of that, motivated by discriminative features that can distinguish anomalies, we propose a novel approach, HADISD, to tackle the problem of **H**eterogeneous **A**nomaly **D**etection **I**n **S**ocial **D**iffusion. Specifically, a storage structure is designed for low space complexity; a parameter-free searching algorithm is proposed to capture the maximum set of vertices satisfying specific features and adapting to data dynamics; and finally, we develop a solution to identify the minimum set of vertices having the discriminative features along the time. The contribution can be summarized as follows.

- We utilize discriminative features to investigate anomaly detection problem in dynamic social diffusion processes, and the results can be leveraged to solve two popular problems: spammer detection and influence maximization.
- We define a *heterogeneous anomaly detection* problem, which retrieves abnormal social behaviors in diffusion processes.
- We propose a novel solution HADISD. The proposed approach does not need to study the information spread in social networks, and most importantly, it has the advantages of being *parameter-free* in searching and updating with low storage consumption and high efficiency.
- We conduct four empirical studies on large-scale real-life datasets. Our approach outperforms the state of the art in terms of accuracy, efficiency, and scalability.

The remainder of the paper is organized as follows. Section 2 surveys the related work. In Section 3, we formally describe the problem. Section 4 presents the solution, HADISD followed by the experimental studies in Section 5. Finally, we conclude the work in Section 6 with potential directions for future work.

## 2. Related work

Social network analysis is of importance to numerous applications, which has attracted a large body of research. As follows, we survey the major latest work from three aspects: spammer detection, influence maximization, and heterogeneous anomaly detection, respectively.

### 2.1. Spammer detection

Spammer detection methods in message systems are hard to be adopted by social network applications due to various social activities and network structures [53]. Spammer detection in complex social networks is usually based on either user-generated content or social behaviors. We proceed to briefly discuss the representative approaches. Co-classification on a labeled training data has been used to detect Web spam on social media Web sites [8]. Page value metrics assessed by using the information of bidirectional links can be defined to filter out spam sites and identify reputable ones [50]. Based on deceptive spam profiles, a honeypot-based approach is designed to uncover social spammers in online social systems [31]. The methods based on social behaviors mostly employ spamming behavior modeling and analysis for spammer detection.