# CoReRank: Ranking to Detect Users Involved in Blackmarket-Based Collusive Retweeting Activities

Aditya Chetan*, Brihi Joshi*, Hridoy Sankar Dutta*, Tanmoy Chakraborty
IIIT-Delhi, India
{aditya16217,brihi16142,hridoyd,tanmoy}@iiitd.ac.in

## ABSTRACT

Twitter's popularity has fostered the emergence of various illegal user activities – one such activity is to artificially bolster visibility of tweets by gaining large number of retweets within a short time span. The natural way to gain visibility is time-consuming. Therefore, users who want their tweets to get quick visibility try to explore shortcuts – one such shortcut is to approach the blackmarket services, and gain retweets for their own tweets by retweeting other customers' tweets. Thus the users intrinsically become a part of a collusive ecosystem controlled by these services.

In this paper, we propose CoReRank, an unsupervised framework to detect *collusive users* (who are involved in producing artificial retweets), and *suspicious tweets* (which are submitted to the blackmarket services) simultaneously. CoReRank leverages the retweeting (or quoting) patterns of users, and measures two scores – the 'credibility' of a user and the 'merit' of a tweet. We propose a set of axioms to derive the interdependency between these two scores, and update them in a recursive manner. The formulation is further extended to handle the cold start problem. CoReRank is guaranteed to converge in a finite number of iterations and has linear time complexity. We also propose a semi-supervised version of CoReRank (called CoReRank+) which leverages a partial ground-truth labeling of users and tweets. Extensive experiments are conducted to show the superiority of CoReRank compared to six baselines on a novel dataset we collected and annotated. CoReRank beats the best unsupervised baseline method by 269% (20%) (relative) average precision and 300% (22.22%) (relative) average recall in detecting collusive (genuine) users. CoReRank+ beats the best supervised baseline method by 33.18% AUC. CoReRank also detects suspicious tweets with 0.85 (0.60) average precision (recall). To our knowledge, **CoReRank is the first** *unsupervised method* **to detect collusive users and suspicious tweets** *simultaneously* **with** *theoretical guarantees.*

## CCS CONCEPTS

• **Information systems** → **Social networks**; • **Security and privacy**;

---

*Equal contribution. Arranged in the alphabetical order of the first name.

---

## KEYWORDS

Retweets, collusion, Twitter, blackmarket, Online Social Networks

## 1 INTRODUCTION

With the engagement of a large number of users, Online Social Networks (OSNs) have become a major platform to achieve social reputation. The organic way to gain significant social reputation is a time consuming process. Therefore, users who join social media for promotion, big-product launches, etc. may not want to wait for such a long time to gain reputation; rather they may want to explore some shortcuts. This has led to the increasing trend of gaining rapid boost with the help of blackmarket services. These services provide support on all major social networks – Instagram, Facebook, Twitter, YouTube, Vimeo, Pinterest, SoundCloud, Vine etc., and assure customers that the services (followers/retweets/views/likes) they provide would appear to be genuine and extremely difficult to spot. Blackmarket services are categorized into two types based on the mode of service [28] – *premium* and *freemium*. Premium blackmarkets provide services upon deposit of money. On the other hand, freemium services provide an additional option of unpaid services where customers themselves become a part of these services, participate in fake activities (following, retweeting other, etc.) and gain (virtual) credits. Hence, they become a part of the collusive echo-chamber controlled by these services. Here, we focus our attention on users involved in 'freemium retweeting services' and call them **"collusive users"** - *users who retweet/quote tweets submitted to the blackmarket services and earn credits in return.* We address the following problem – *how can we design an efficient system to simultaneously detect users (based on their unusual retweeting pattern) and tweets (based on the credibility of the users who retweet them) involved in collusive blackmarket services?*

**Background and Motivation:** Current state-of-the-art algorithms mostly focus on detecting bots [5, 33] and users involved in synchronous fraudulent activities [11]. The only work close to the problem we address here is a very recent effort by Dutta et al. [7]. They empirically showed that existing fake detection methods fail to identify collusive users as these users are not bots and their behavior is not synchronous. They conducted experiments only on 743 collusive retweeters collected from various freemium blackmarket services, and used standard supervised methods to

classify collusive retweeters and genuine users. They mentioned that detection of collusive users is challenging for two reasons – (i) These users neither resemble bots nor fake users as they express *a mix of organic and inorganic behavior*. Therefore, bot detection and fake user detection algorithms can not detect them. (ii) Collecting large scale labeled data of collusion users is extremely challenging. This necessitates the design of an unsupervised approach to detect collusive users.

**Proposed Approach:** In this paper, we propose CoReRank, an unsupervised approach to simultaneously detect collusive users and suspicious tweets. We model the interactions between users and tweets in terms of retweets/quotes (collectively called as 'support') using a directed bipartite graph. We capture the interdependency between the *credibility of users* and the *merit of tweets* via a set of axioms. We further design recurrence formulations combining the graph, behavioral information and topical diversity of supported tweets to obtain the final score of the 'credibility' of users and the 'merit' of tweets. The cold start problem is also handled using Laplace smoothing. We further design CoReRank+ which extends CoReRank to a semi-supervised setting (when few labeled data are available). We theoretically show that CoReRank is guaranteed to converge in a finite number of iterations, and it scales linearly with the number of edges present in the graph.

**Experimental Results:** We started by collecting users and tweets from two blackmarket services, and used them as seeds to expand the data to $10K$ users and 2.5 million tweets. Human annotators were employed to label unknown users as collusive/genuine. **This, to our knowledge, is the first labeled dataset of collusive/genuine users based on their retweeting activities**. We compare CoReRank with SCoRe [7], the only available system to detect collusive users involved in artifical retweeting activity. We also compare CoReRank with five other baselines designed to address similar type of problems (fake/bot detection, etc.). CoReRank outperforms the best unsupervised baseline by 269% and 300% for collusive user detection, and 20% and 22.22% for genuine user detection in terms of average precision and average recall (relative) respectively. CoReRank+ also beats the best supervised method by 33.18% (relative) AUC. The added benefit of CoReRank is that it also identifies suspicious tweets with 0.85 precision and 0.60 recall. Empirical results further show that – (i) graph-based interdependency is the most effective component of our models, (ii) the running time of CoReRank is linear in the number of edges and much faster than any baseline.

**Contribution:** The overall contribution is six-fold:

- **Problem definition:** This paper is the second attempt after [7] which addresses the problem of detecting blackmarket users involved in collusive retweeting activities.
- **Algorithm:** CoReRank is the first *unsupervised approach* to simultaneously detect collusive users as well as suspicious tweets (submitted to the blackmarket services to gain retweets).
- **Theoretical guarantee:** CoReRank is guaranteed to converge in a finite number of iterations.
- **Effectiveness:** CoReRank outperforms six state-of-the-art methods by a significant margin.
- **Scalability:** CoReRank is fast, and scales linearly with the number of edges present in the constructed graph.

- **Dataset:** We collected and annotated a novel dataset of collusive users, which is the first labeled dataset of this kind.

**Reproducibility:** The code of CoReRank and the dataset are available at https://github.com/LCS2-IIITD/CoReRank-WSDM-2019.

## 2 RELATED WORK

We organize our related work into two parts: (i) general study on fraud detection in OSNs, and (ii) study on understanding blackmarket activities.

**Fraud Detection in OSNs:** A series of works investigated fraud detection from various aspects. [19] proposed a honeypot-based approach to detect spammers on Twitter and MySpace. They used profile-based features of the spammers to design supervised classifiers. [2] also used supervised classifiers with features based on tweet content and user social behavior. [36] detected criminal profiles on Twitter. A large number of papers designed methods to detect bots on Twitter. [3] analyzed the behavior of humans, bots, and cyborgs. They reported that humans have complex timing behavior while bots and cyborgs have a periodic behavior. [5] and [33] designed unsupervised and supervised approaches respectively to detect bots. [18, 34] focused on the detection of spam tweets based on tweet-related features. [17] studied fraud favoritism on retweets. Some other studies of fraud detection based on tweet content are [1, 9, 10, 23]. Multiple papers focused on detecting fraud using URLs embedded in tweets [19, 35] and blacklisted URLs [10, 13, 14]. [11] proposed NDSync to detect synchronized fraud activities (fake retweets) on Twitter. [12] studied the influence of fraudulent and genuine retweet threads and discovered patterns ('Triangles' and 'Homogeneity') followed by fraudulent users.

**Study of Blackmarket Services:** Compared to the study on general fraud detection, the effort so far has been limited to investigate fraudulent activities by blackmarket services. A detailed analysis of blackmarkets with the impact on multiple social networks was studied by [6]. [28] studied the multifaceted behavior of the blackmarket agencies. [30] studied the market size and market price of multiple Twitter follower markets. [31] used supervised approach to detect accounts of blackmarket services. [29] studied the characteristics of Twitter follower merchant markets. [20] proposed a method to detect followers who provide voluntary following services to make the profit. [4] used supervised approach for fake blackmarket follower detection.

Few studies attempted to analyze blackmarkets in other platforms such as live video-streaming, online recruitment, etc. [24] analyzed six different underground forums where users are involved in selling goods and services. [27] proposed an unsupervised approach to detect bot-generated broadcasts and views for famous broadcasting platforms (YouTube and Twitch). [32] used machine learning approach and textual analysis on a dataset containing real-life job ads to detect whether given employment or job advertisement is legitimate or fraudulent. [7] studied blackmarket services which provide fake retweets and designed supervised approach to detect collusive retweeters.

**Major Differences with Existing Approaches:** CoReRank combines both *graph structure* and *behavioral properties* in an *unsupervised manner* and returns *ranked lists of users and tweets simultaneously* based on *collusive retweeting activities*. CoReRank is also

**Table 1: Comparison of `CoReRank` and other baseline methods w.r.t different dimensions of an algorithm.**

|  | [5] | [33] | [9] | [11] | [7] | [15] | Our |
|---|---|---|---|---|---|---|---|
| Address collusion phenomenon |  |  |  |  | ✓ |  | ✓ |
| Consider graph information |  | ✓ |  |  |  | ✓ | ✓ |
| Consider topic information |  |  |  |  |  |  | ✓ |
| Unsupervised approach | ✓ |  |  | ✓ |  | ✓ | ✓ |
| Return ranked list of users |  |  |  | ✓ |  |  | ✓ |
| Detect both collusive users and tweets |  |  |  |  |  |  | ✓ |
| Theoretical guarantees |  |  |  |  |  | ✓ | ✓ |

*guaranteed to converge* in a finite number of iterations with linear time complexity. Table 1 summarizes a comparison of `CoReRank` with other existing approaches. `CoReRank` is the only one which matches all specifications.

## 3 DATA DESCRIPTION

We solicited tweets from various freemium blackmarket services. We searched for these services by querying search engines with keywords such as 'Free Retweets', 'Retweet my Tweet', etc. We collected data from the following services (after taking proper IRB approval) - YouLikeHits[1] and Like4Like[2]. These services provide an 'earning area' where tweets submitted by the customers of the service are displayed so that other customers can gain credits by retweeting the tweets. We designed a web-scraper to extract tweets ($T_b$) as well as the customers ($U_c$) who submitted these tweets. We created three user sets: ground-truth genuine user set $\mathcal{S}_g$, ground-truth collusive user set $\mathcal{S}_c$, and unknown user set $\mathcal{S}_u$. We collected users ($U_b$) who retweeted/quoted the extracted tweets. Out of these users, we found (i) 4 verified Twitter users ($U_b^v$), and (ii) 329 users who were also a part of $U_c$. We added the former to $\mathcal{S}_g$ and latter to $\mathcal{S}_c$. The remaining 7451 users were added to $\mathcal{S}_u$. We further increased $\mathcal{S}_g$ by adding the followees of the verified users ($U_f^v$) with the assumption that verified users are more likely to follow genuine users. We also collected the tweets ($T_f^v$) that were retweeted/quoted by $U_f^v$. Finally, we collected all retweets and quotes (max. 3200) of users from their timeline. At this stage, the size of the ground-truth sets is as follows: $|\mathcal{S}_c| = 329$, $|\mathcal{S}_g| = 2667$ and $|\mathcal{S}_u| = 7451$. Users present in $\mathcal{S}_u$ were further annotated by human experts as collusive/genuine (see Section 5.2). The graph constructed from the dataset is described in Section 4.1.2.

## 4 PROPOSED METHODOLOGY

In this section, we explain our efforts in formulating `CoReRank`. (which is motivated by [7, 16]).

### 4.1 CoReRank Preliminaries

We start by constructing a graph comprising users and tweets as nodes. We hypothesize that users and tweets have intrinsic traits that often demonstrate their collusive nature and their credibility. Thus, users and tweets can be allotted scores that define these traits. The reason to operate on the graph and not individual users/tweets (as done in [7]) is that these scores are interdependent on each other and the graph as a whole.

*Definition 4.1.* [**User Support**] A tweet $t$ is considered to be *supported* (by retweeting or quoting) by a user $u$ if $u$ either retweeted or quoted $t$. This is given by $S(u, t)$ as follows -

$$S(u, t) = \begin{cases} w_q & \text{if } u \text{ quoted } t \\ w_r & \text{if } u \text{ retweeted } t \\ 0 & \text{otherwise} \end{cases}$$

Here, $w_r$ ($w_q$) denote the weight of the edge when $u$ retweeted (quoted) $t$. The relation between $w_r$ and $w_q$ can be defined by: $0 < w_r \leq w_q < 1$ as a quote is essentially a retweet augmented with further text or media, allowing the tweet to have more importance or weight than a simple retweet. In our model, we set $w_r$ and $w_q$ to 0.5 and 0.75 respectively. However, we show in Supplementary [26] how the performance of our algorithm varies by changing the values of $w_r$ and $w_q$.

*Definition 4.2.* [**Support Graph**] A bipartite support graph $G = (U, T, E)$ is a directed bipartite graph where $U$ indicating the set of users forms the left partition, and $T$ indicating the tweets supported by $U$ forms the right partition. Edge $E_{(u, t)}$ connecting user $u \in U$ and tweet $t \in T$ indicates that $u$ supported $t$ with the edge weight $S(u, t)$ denoting the kind of support $u$ extends to $t$.

*4.1.1 Graph Construction.* The graph construction is divided into four distinct steps as follows (we use the same notations mentioned in Section 3).

(i) The users $U_b$ who supported tweets $T_b$ (submitted to Blackmarkets) form the left partition, and the submitted tweets $T_b$ form the right partition of $G$. $U_b$ and $T_b$ are then connected by directed edges $E_{(U_b, T_b)}$; $S(u, t)$ denotes the edge weight connecting $u \in U_b$ and $t \in T_b$. We call this graph as the upper half $G_u$ of the final graph $G$.

(ii) The left partition of the graph is further augmented with $U_b^v$, the verified users in $U_b$ and $U_f^v$, the followees of $U_b^v$. The right partition is populated by augmenting $T_f^v$, the tweets supported by $U_f^v$. $U_f^v$ and $T_f^v$ are connected by directed edges $E_{(U_f^v, T_f^v)}$ with weights given by $S(., .)$. We call this as the lower half $G_l$ of the final graph $G$.

(iii) The next step is to look for possible connections $E_{(U_b, T_f^v)}$ between users in $G_u$ and tweets in $G_l$. If a user $u \in U_b$ supported a tweet $t \in T_f^v$, a directed edge is added from $u$ and $t$ with weight $S(u, t)$.

(iv) Similarly, we search for possible connections $E_{(U_f^v, T_b)}$ between users in $G_l$ and tweets in $G_u$. If a user $u \in U_f^v$ supported a tweet $t \in T_b$, a directed edge is added from $u$ and $t$ with weight $S(u, t)$.

After performing the above operation, we obtain the final support graph $G = (U, T, E)$, where $U = \{U_b, U_f^v\}$, $T = \{T_b, T_f^v\}$ and $E = \{E_{(U_b, T_b)}, E_{(U_f^v, T_f^v)}, E_{(U_b, T_f^v)}, E_{(U_f^v, T_b)}\}$. $G$ turns out to be a connected graph. Table 2 shows the statistics of $G$.

*4.1.2 Intrinsic Traits of Users and Tweets.* Users and Tweets have intrinsic traits on the basis of which they can be given a score to determine the 'credibility' (for users) and 'merit' (for tweets). Let $Out(u)$ be the set of tweets user $u$ supported, and $In(t)$ be the set

**Table 2: Statistics of the bipartite support graph $G$.**

|  | Left partition | Right partition | Edges |
|---|---|---|---|
| $G_u$ | $|U_b| = 7784$ | $|T_b| = 1001$ | $|E_{(U_b, T_b)}| = 55382$ |
| $G_l$ | $|U_f^v| = 2667$ | $|T_f^v| = 2439319$ | $|E_{(U_f^v, T_f^v)}| = 2862793$ |
| $G_u$ to $G_l$ | # nodes in $U_b = 294$ | # nodes in $T_f^v = 10512$ | $|E_{(U_b, T_f^v)}| = 44466$ |
| $G_l$ to $G_u$ | # nodes in $U_f^v = 10$ | # nodes in $T_b = 14$ | $|E_{(U_f^v, T_b)}| = 96$ |
| $G$ | $|U| = 10451$ | $|T| = 2440320$ | $|E| = 2962737$ |

of users that support $t$. Below we provide a detailed discussion on the interpretation and relevance of these intrinsic properties.

**Credibility of Users:** The credibility of a user is an indication of how likely they are to support a tweet based on their genuine agreement with the content of the tweet. A highly credible user would only support those tweets which are about the *topics* that they frequently support. In contrast, a user with low credibility, who might be involved in collusive activities would support any tweet that they come across on a blackmarket service. This would highly diversify their timeline w.r.t the topics of supported tweets. Figure 1(b) supports this hypothesis. The credibility score of a user $u$ is given by $C(u)$, where $C(u)$ ranges from 0 (very high collusive trait) to 1 (very high credibility trait).

**Merit of Tweets:** The merit of a tweet is an indication of the genuine organic support of users. A meritorious tweet would be supported by more credible users, indicating that it has a genuine support in the graph and is not earning support because of its submission to a blackmarket service. In contrast, a tweet with lower merit, even though it has a higher user support overall, would still lack the support of users with high credibility. The merit of a tweet $t$ is given by $M(t)$, where $M(t)$ ranges from 0 (highly suspicious) to 1 (highly genuine).

## 4.2 CoReRank Properties

The credibility of user $u$ depends majorly on the merit of the tweets $Out(u)$ that $u$ supported. Similarly, the merit of tweet $t$ depends majorly on the credibility of the users $In(t)$ who supported $t$. The axioms mentioned below capture this interdependency.

*Definition 4.3.* [**Inter-support time**] Given a user $u$, let $t_1$ and $t_2$ be two consecutive tweets that $u$ supported at $T_1$ and $T_2$, respectively. The time difference between $t_1$ and $t_2$ is thus given by $T_1 - T_2$. We construct a set of such time differences for *all* consecutive tweets that $u$ supported, and form a set of inter-support times for $u$. This set is denoted by $IST(u)$.

AXIOM 1 (**Collusive users have very less inter-support times compared to genuine users**). *Freemium blackmarket services require users to get 'points' to gain retweets and other privileges. In order to keep scoring, these users keep on retweeting random tweets, which leads to very less inter-support times. Hence, the average inter-support time of a collusive user, $avg(IST(u_1))$ is very less than that of a genuine user. Formally,*

$$\exists u_1, u_2 \in U, C(u_1) < C(u_2) \implies avg(IST(u_1)) < avg(IST(u_2))$$

Figure 1(a) validates Axiom 1. Collusive users support tweets in a succession, within a minute; whereas genuine users take days or months to support tweets.

*Definition 4.4.* [**Identically collusive support of tweets**] Two tweets, $t_1$ and $t_2$, are said to have identically collusive support if,
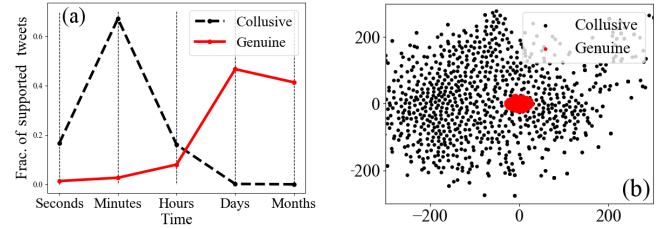


**Figure 1: (a) Inter-support time for collusive and genuine users. Collusive users are very fast in supporting tweets. (b) Projection of tweets supported by one collusive user and one genuine user. We use t-SNE plot [21] to visualize the tweet space obtained from GloVe embedding (see Section 4.3.2 for details of embedding). Tweets supported by the genuine (collusive) user are on same (different) topic(s).**

$|In(t_1)| = |In(t_2)|$ and there exists a bijection $f$ from $In(t_1)$ to $In(t_2)$ such that, $f(u) = u' \implies C(u) = C(u') \ \forall u \in In(t_1), u' \in In(t_2)$.

AXIOM 2 (**Among tweets with identically collusive support, a highly meritorious tweet receives higher support**). *For two tweets, $t_1$ and $t_2$ with identically collusive support, if $S(u, t_1) \geq S(u', t_2)$, where $u \in In(t_1) \land u' \in In(t_2)$ such that $C(u) = C(u')$, then $M(t_1) \geq M(t_2)$.*

AXIOM 3 (**A collusive user associated with blackmarket services demonstrates immense topical diversity**). *As per Axiom 1, a collusive user tends to retweet every tweet that can help in receiving more credits. Thus, the tweets supported by that user tend to be very diverse in terms of the topic.*

Figure 1(b) supports Axiom 2. We show that the topics of the tweets supported by a collusive user are highly diverse (in the embedding space, points are scattered around the space), whereas the same for a genuine user is uniform (points are clustered).

## 4.3 CoReRank Formulation

CoReRank is an unsupervised approach which considers a directed bipartite graph of users and tweets. We propose that users have unknown intrinsic scores that quantify how trustworthy they are, and tweets have unknown intrinsic scores that quantify its natural merit of being retweeted organically. Naturally, these scores are interdependent and unknown apriori. Here we start by describing how to obtain seed score which quantifies the intrinsic score of users and tweets. We also define topical diversity score of a user by calculating the inter-support similarity of the tweets s/he retweeted. Finally, we show how one can combine seed score, topical diversity and behavioral activities to jointly obtain the credibility of users and the merit of tweets.

*4.3.1 Seed Score.* In order to compute the seed scores both for users and tweets that help CoReRank propagate further, we use Birdnest [15] (as suggested by [16]). Birdnest takes the following scores as inputs and uses Bayesian estimates to formulate how much the properties of a user or a tweet deviates from the rest.

- **User-specific Score**, is calculated by providing a vector that contains the inter-support times of the users. This array can be of variable length for each user and is provided as an input to Birdnest.

- **Tweet-specific Score**, is calculated with the aid of the length (word count) of the tweet. The length of a retweet (quote) is the length of the original tweet that was retweeted (the sum of the lengths of the original tweet and the quoted text). Each tweet is associated with a vector, whose each entry indicates the length of its retweet/quote. This vector is passed as an input to Birdnest. We hypothesize that the text length of a supported tweet is a measure of its behavioural activity [16].

The final outputs of Birdnest are the suspicion scores for users $S_U(u)$ and tweets $S_U(T)$. Finally, the seed scores, $\pi_U(u)$ and $\pi_T(t)$ for a user and a tweet, respectively are given by,

$$\pi_U(u) = 1 - S_U(u) \; \forall u \in U$$

$$\pi_T(t) = 1 - S_T(t) \; \forall t \in T$$

In some cases, it might happen that the user (tweet) provided (received) only a single support. In such cases, Birdnest does not provide us with a seed score. For these users and tweets, we have assigned the highest possible seed score of 1, giving them the benefit of doubt.

*4.3.2 Inter-support Similarity.* To calculate the topical diversity score for a user, we first extract information from all the tweets present in our graph $G$. For each tweet of a user, we split it into words and derive their word representation in embedding space using GloVe embedding trained specifically with Tweets [25]. For this work, we use 100 dimensional pre-trained GloVe vectors. The final embedding of a tweet is obtained by combining the embedding of each word present in the tweet. We do not take into account the words for which the pre-trained word vectors are not found in GloVe. Finally, we measure $\tau_U(u)$, the inter-support similarity for user $u$ by calculating the average cosine similarity of the embeddings corresponding to the tweets supported by $u$.

*4.3.3 Recurrence Formulation.* We here propose a systematic approach to combine the following three signals into a recurrence framework – (i) graph-based interdependency between users and tweets, (ii) behavioral activities of users and tweet, and (iii) topical diversity of tweets supported by a user. We also present a way of handling cold start problem.

**Graph-based Interdependency of Users and Tweets:** As presented in the axioms mentioned previously, the credibility of a user is influenced by the merit of all the tweets that the user supported. Similarly, the merit of a tweet depends on the credibility of the users who supported it. We present below two mathematical formulations that incorporate interdependency among credibility and merit of users and tweets respectively.

For tweets,

$$M(t) = \frac{\sum_{u \in \text{In}(t)} \gamma_{1t} \cdot C(u) \cdot S(u, t)}{\gamma_{1t} + |\text{In(t)}|} \tag{1}$$

Similarly for users,

$$C(u) = \frac{\sum_{t \in \text{Out}(u)} \gamma_{1u} \cdot M(t) \cdot S(u, t)}{\gamma_{1u} + |\text{Out(u)}|} \tag{2}$$

Here, $\gamma_{1u}, \gamma_{1t}$ are constants for users and tweets, respectively. Their values will be learned by parameter sweeping (Section 4.4).

**Behavioral Activities:** Apart from graph properties, behavioral properties are also included while calculating the credibility and merit of users and tweets. We modify Equations 1 and 2 to incorporate the seed scores explained in Section 4.3.1 of the users and tweets in the formulation of credibility and merit respectively.

$$M(t) = \frac{\sum_{u \in \text{In}(t)} \gamma_{1t} \cdot C(u) \cdot S(u, t) + \gamma_{2t} \cdot \pi_T(t)}{\gamma_{1t} + \gamma_{2t} + |\text{In(t)}|} \tag{3}$$

$$C(u) = \frac{\sum_{t \in \text{Out}(u)} \gamma_{1u} \cdot M(t) \cdot S(u, t) + \gamma_{2u} \cdot \pi_U(u)}{\gamma_{1u} + \gamma_{2u} + |\text{Out(u)}|} \tag{4}$$

Here, $\gamma_{2t}$ and $\gamma_{2u}$ are also constants for users and tweets respectively and will be learned by parameter sweeping (Section 4.4).

**Topic-based Inter-support Similarity:** Finally, we add $\tau_U(u)$, the inter-support similarity of user $u$ into the credibility formulation. The lower the inter-support similarity, the higher the probability that the user is collusive. The modified version of Equation 3 is as follows:

$$C(u) = \frac{\sum_{t \in \text{Out}(u)} \gamma_{1u} \cdot M(t) \cdot S(u, t) + \gamma_{2u} \cdot \pi_U(u) + \gamma_{3u} \cdot \tau_U(u)}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + |\text{Out(u)}|} \tag{5}$$

Here, $\gamma_{3u}$ is a parameter for users. We learn this parameter by performing parameter sweeping as well (Section 4.4).

**Handling Cold Start:** For some users as well as tweets, multiple supports may not be available (as discussed in Section 4.3.1). Hence, getting seed scores for them is a problem. Also, such users and tweets can lead to biased rankings with high credibility or merit scores. We solve this problem using Laplace smoothing by assigning a default score $\mu_T$ to tweets and $\mu_U$ to users, weighed by parameters $\gamma_{3t}$ and $\gamma_{4u}$ respectively. The value of $\gamma_{3t}$ and $\gamma_{4u}$ decides how much importance is given to the default scores – high value implies less influence of the graph structure in the scores of users and tweets. Thus, the modified formulation of merit and credibility is given below.

$$M(t) = \frac{\sum_{u \in \text{In}(t)} \gamma_{1t} \cdot C(u) \cdot S(u, t) + \gamma_{2t} \cdot \pi_T(t) + \gamma_{3t} \cdot \mu_T}{\gamma_{1t} + \gamma_{2t} + \gamma_{3t} + |\text{In(t)}|} \tag{6}$$

$$C(u) = \frac{\sum_{t \in \text{Out}(u)} \gamma_{1u} \cdot M(t) \cdot S(u, t) + \gamma_{2u} \cdot \pi_U(u) + +\gamma_{3u} \cdot \tau_U(u) + \gamma_{4u} \cdot \mu_U}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + \gamma_{4u} + |\text{Out(u)}|} \tag{7}$$

Here, $\gamma_{3t}$ and $\gamma_{4u}$ are parameters for tweets and users respectively, and will be learned by parameter sweeping (Section 4.4). The default scores are set as the average of the initial scores of all users and tweets.

## 4.4 The CoReRank Algorithm

We now briefly describe the CoReRank algorithm (see Algorithm 1 for the pseudo-code). Given $G(U, T, E)$, the support graph of users and tweets, CoReRank takes all constants – $\gamma_{1t}$, $\gamma_{2t}$, $\gamma_{3t}$, $\gamma_{1u}$, $\gamma_{2u}$, $\gamma_{3u}$, $\gamma_{4u}$ as parameters. At first we calculate the seed scores for all the users and tweets, using the Birdnest algorithm. For users and tweets for which the seed score could not be found, it is initialized to the highest value, i.e., 1. In the first iteration, the scores for all the users and tweets are initialized to their seed scores. Next, we calculate the cold start constants for users and tweets, i.e., $\mu_U$ and $\mu_T$ respectively. At the beginning of each iteration, we normalise

**Algorithm 1:** CoReRank Algorithm

**Input** : $G(U, T, E), \gamma_{1t}, \gamma_{2t}, \gamma_{3t}, \gamma_{1u}, \gamma_{2u}, \gamma_{3u}, \gamma_{4u}$

**Output** : Credibility and Merit scores for all users and tweets

1 Calculate $\pi_U(u) \ \forall u \in U$ and $\pi_T(t) \ \forall t \in T$

2 Initialize $C(u)^0 = \pi_U(u)$ and $M(t)^0 = \pi_T(t) \ \forall u \in U, \forall t \in T$

3 Initialize $\mu_U = \frac{\sum_{u \in U} C(u)^0}{|U|}$ and $\mu_T = \frac{\sum_{t \in T} M(t)^0}{|T|}$

4 $k = 0$

5 error = maximum possible integer value

6 **while** error $> \epsilon$ **do**

7      $k = k + 1$

8      $\tilde{C}^{k-1}(u) = \text{norm}(C^{k-1}(u)) \forall u \in U$ such that
         $\tilde{C}^{k-1}(u) \in [0, 1]$

9      Update the merit of tweets using Equation 6: $\forall t \in T$,

10      $M^k(t) = \frac{\sum_{u \in \text{In}(t)} \gamma_{1t} \cdot \tilde{C}^{k-1}(u) \cdot S(u,t) + \gamma_{2t} \cdot \pi_T(t) + \gamma_{3t} \cdot \mu_T}{\gamma_{1t} + \gamma_{2t} + \gamma_{3t} + |\text{In}(t)|}$

11      Update the credibility of users using Equation 7: $\forall u \in U$,

12      $C^k(u) =$
         $\frac{\sum_{t \in \text{Out}(u)} \gamma_{1u} \cdot M^k(t) \cdot S(u,t) + \gamma_{2u} \cdot \pi_U(u) + \gamma_{3u} \cdot \tau_U(u) + \gamma_{4u} \cdot \mu_U}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + \gamma_{4u} + |\text{Out}(u)|}$

13      error =
         $\max(\max_{u \in U} |C^k(u) - C^{k-1}(u)|, \max_{t \in T} |M^k(t) - M^{k-1}(t)|)$

14 **end**

15 **return** $C^k(u) \ \forall u \in U$ and $M^k(t) \ \forall t \in T$

the credibility scores for users using *Min-Max Normalization*. Following this, we keep computing scores using Equations 6 and 7. Convergence is achieved when the maximum of the maximum of differences of scores between $(t + 1)^{th}$ and $t^{th}$ iterations, i.e., *error* is less than a very small value, $\epsilon$ (set as $10^{-6}$).

The convergence factor $\epsilon$ is chosen using a grid search in $\{10^{-6}, 10^{-4}, 0.01, 0.1, 0.5, 1\}$. We use *parameter sweep* to tune other parameters − $\gamma_{1t}, \gamma_{2t}, \gamma_{3t}, \gamma_{1u}, \gamma_{2u}, \gamma_{3u}$ and $\gamma_{4u}$. In this method, possible combinations of all the parameters from a given list of values are considered. The possible values for $\gamma_{1t}, \gamma_{2t}, \gamma_{3t}, \gamma_{1u}, \gamma_{2u}, \gamma_{4u}$ are taken in the range $[0, 1]$ with a step of 0.3, while $\gamma_{3u}$ is chosen from $\{0, 1, 2, 3\}$. Maximum accuracy is achieved with the following parameter setting: $\gamma_{1t} = 0.6, \gamma_{2t} = 0.6, \gamma_{3t} = 0.3, \gamma_{1u} = 0.6, \gamma_{2u} = 0.6, \gamma_{3u} = 3, \gamma_{4u} = 0.3$ (see Supplementary [26] for more details).

> **Output:** The outputs of CoReRank are ranked lists of users and tweets based on $C(u)$ and $M(t)$ respectively. Ranking in descending (ascending) order of $C(u)$ will place the genuine (collusive) users at the top of the ranking (same for the tweets based on $M(t)$).

## 4.5 CoReRank+: A Semi-supervised Version

Oftentimes, we have partial knowledge about the labels of some users (verified, blackmarket customers, etc.) and tweets. We can leverage such prior information and incorporate them to our formulation in a semi-supervised fashion. We provide each user $u$ with a label score, $\alpha_U(u)$ which can be defined as follows -

$$\alpha_U(u) = \begin{cases} \alpha_U^c & \text{if } u \text{ is a collusive user} \\ \alpha_U^g & \text{if } u \text{ is a genuine user} \\ 0 & \text{if } u \text{ has no pre-defined label} \end{cases}$$

Similarly, the label score for a tweet $t$, $\alpha_T(t)$ can be defined as follows -

$$\alpha_T(t) = \begin{cases} \alpha_T^c & \text{if } t \text{ is a labeled as suspicious} \\ 0 & \text{if } t \text{ has no pre-defined label} \end{cases}$$

We set the values of these constants in such a manner that the genuine users are awarded with a high positive label score and the customers of blackmarket services are given a high negative label score. In our experiment, we set these constants as follows: $\alpha_U^c = -100, \alpha_U^g = 100$ and $\alpha_T^c = -100$.

The final equations of credibility and merit after incorporating the labels are given below -

$$M(t) = \frac{\sum_{u \in \text{In}(t)} \gamma_{1t} \cdot C(u) \cdot S(u, t) + \gamma_{2t} \cdot \pi_T(t) + \gamma_{3t} \cdot \mu_T + \alpha_T(t)}{\gamma_{1t} + \gamma_{2t} + \gamma_{3t} + |\text{In}(t)|} \tag{8}$$

$$C(u) = \frac{\sum_{t \in \text{Out}(u)} \gamma_{1u} \cdot M(t) \cdot S(u, t) + \gamma_{2u} \cdot \pi_U(u) + + \gamma_{3u} \cdot \tau_U(u) + \gamma_{4u} \cdot \mu_U + \alpha_U(u)}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + \gamma_{4u} + |\text{Out}(u)|} \tag{9}$$

We design a modified algorithm, called CoReRank+ by replacing the equations mentioned in Lines 10 and 12 of Algorithm 1 by Equations 8 and 9.

## 4.6 Theoretical Guarantee

In this section, we provide the theoretical guarantee to show that CoReRank will converge for a given set of inputs. Before we begin, let $C^\infty(u)$ and $M^\infty(t)$ be the final scores of user $u$ and tweet $t$ respectively. Then we have the following results (see Supplementary [26] for the detailed proof):

LEMMA 4.5 (**Lemma 1**). *For any given tweet, $t$, the difference between their final score and their score after the first iteration of* CoReRank *cannot exceed* $\frac{3}{4}$. *Formally, it means that* $|M^\infty(t) - M^1(t)| \leq \frac{3}{4}$. *Similarly, for users the upper bound is* $\frac{3}{4}$, *i.e.,* $|C^\infty(u) - C^1(u)| \leq \frac{3}{4}$

THEOREM 4.6 (**Theorem of convergence**). *Between successive iterations, the difference in the scores of the users and tweets is bounded. For any user $u \in U$, $|C^\infty(u) - C^k(u)| \leq \left(\frac{3}{4}\right)^k$. Thus, as the algorithm proceeds through more and more iterations, the value of $k$ keeps on increasing and the difference in score from the final score keeps on decreasing. Similarly for a tweet $t \in T$, $|M^\infty(t) - M^k(t)| \leq \left(\frac{3}{4}\right)^{k-1}$.*

**Proof Sketch:** We will prove this using induction on $k$.

**Base Cases ($k = 1$):** This is vacuously true from Lemma 4.5.

**Induction Hypothesis:** Assume that for any $n \leq k$, we have, $|C^\infty(u) - C^k(u)| \leq \left(\frac{3}{4}\right)^k$ and $|M^\infty(t) - M^k(t)| \leq \left(\frac{3}{4}\right)^{k-1}$, $\forall u \in U, \forall t \in T$.

**Induction Step ($k = n + 1$):** For $k = n + 1$, we have,

$$|M^\infty(t) - M^{n+1}(t)| \leq \gamma_{1t} \frac{\sum_{u \in \text{In}(t)} |(\tilde{C}^\infty(u) - \tilde{C}^n(u))| \cdot |S(u, t)|}{\gamma_{1t} + \gamma_{2t} + \gamma_{3t} + |\text{In}(t)|}$$

As $|(\tilde{C}^\infty(u) - \tilde{C}^n(u))| \leq |C^\infty(u) - C^n(u)| \leq \left(\frac{3}{4}\right)^n$,

$$\implies |M^\infty(t) - M^{n+1}(t)| \leq \frac{\gamma_{1t} \cdot \left(\frac{3}{4}\right)^n \cdot |\text{In}(t)| \cdot |S(u, t)|}{\gamma_{1t} + \gamma_{2t} + \gamma_{3t} + |\text{In}(t)|} \leq \left(\frac{3}{4}\right)^n$$
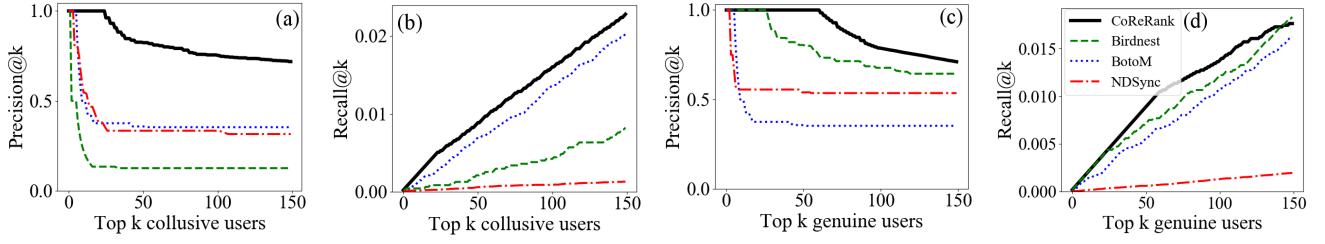
**Figure 2: Change in performance of the competing unsupervised methods with the increase of $k$ (the number of results returned) for detecting both (a-b) collusive and (c-d) genuine users.**

Also, $|C^\infty(u) - C^{n+1}(u)| \leq \frac{\sum_{t \in \text{Out(u)}} |\gamma_{1u}| \cdot |(M^\infty(t) - M^{n+1}(t))| \cdot |S(u,t)|}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + \gamma_{4u} + |\text{Out(u)}|}$

$\leq \frac{|\gamma_{1u}| \cdot \left(\frac{3}{4}\right)^n \cdot |S(u,t)|}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + \gamma_{4u} + |\text{Out(u)}|}$. As $S(u,t) \leq \frac{3}{4}$, $|C^\infty(u) - C^{n+1}(u)|$

$\leq \frac{|\gamma_{1u}| \cdot \left(\frac{3}{4}\right)^{n+1}}{\gamma_{1u} + \gamma_{2u} + \gamma_{3u} + \gamma_{4u} + |\text{Out(u)}|} \leq \left(\frac{3}{4}\right)^{n+1}$

THEOREM 4.7 (**Bound on iterations**). *There exists a bound on the number of iterations until* CoReRank *converges. This bound is governed by the precision to which the score is calculated before convergence is declared, i.e., $\epsilon$. The number of iterations till convergence is at most $2 + \lceil \frac{\log\left(\frac{\epsilon}{2}\right)}{\log\left(\frac{3}{4}\right)} \rceil$*

**Proof:** Let $k = \lceil \frac{\log\left(\frac{\epsilon}{2}\right)}{\log\left(\frac{3}{4}\right)} \rceil$. By Theorem 4.6, after $k + 1$ iterations,

$\forall t \in \mathcal{T}$, $|M^\infty(t) - M^{k+1}(t)| \leq \frac{3}{4}^k \leq \frac{3}{4}^{\log_{\frac{3}{4}}\left(\frac{\epsilon}{2}\right)} = \frac{\epsilon}{2}$. Similarly, $M^\infty(t) - M^{k+2}(t)| \leq \frac{\epsilon}{2} \cdot \frac{3}{4} \leq \frac{\epsilon}{2}$. Thus,

$$|M^{k+1}(t) - M^{k+2}(t)| = |M^{k+1}(t) - M^\infty(t) + M^\infty(t) - M^{k+2}(t)|$$

As $|x + y| \leq |x| + |y|$,

$\implies |M^{k+1}(t) - M^{k+2}(t)| \leq |M^{k+1}(t) - M^\infty(t)| + |M^\infty(t) - M^{k+2}(t)| \leq 2 \cdot \frac{\epsilon}{2} = \epsilon$

Similarly for credibility, $|C^{k+1}(u) - C^{k+2}(u)| \leq \epsilon \forall u \in \mathcal{U}$.

Thus, by Line 6 of Algorithm 1 it will take $k + 2$ iterations to converge. Hence, a low value of $\epsilon$ would require a larger number of iterations for the algorithm to converge.

**Time Complexity of CoReRank:** In each iteration, CoReRank updates the scores of users and tweets in constant time. Thus, the complexity of each iteration is $O(|E| + |V|)$. Since $|E| >> |V|$ in $G(U, T)$, $O(|E| + |V|) \approx O(|E|)$. As explained in Theorem 4.7, CoReRank converges in a constant number of iterations. Let $n$ be the product of the number of iterations till convergence and the number of runs of CoReRank. Thus, the time complexity of algorithm is $O(n|E|)$, which is linear in the number of edges of $G(U, T)$ (as supported empirically in Figure 3(d)).

## 5  EXPERIMENTAL RESULTS

This section starts by briefly explaining the baselines and human annotation process, followed by the detailed performance analysis of the competing methods.

### 5.1  Baseline Methods

As mentioned before, CoReRank is the second algorithm after SCoRe which addresses the problem of detecting collusive users involved in fake retweeting activities. We choose six state-of-the-art (supervised and unsupervised) methods as baselines, which are similar to the problem we address here. (i) **Baseline I** (BotoM): [5] measured bot score for each user and produced a ranked list of users. (ii) **Baseline II** (NDSync): [11] detected synchronous retweet fraudulent activities by calculating a user-level suspiciousness score which combines the suspiciousness score for each retweet thread by projecting them into a multi-dimensional feature space. (iii) **Baseline III** (Birdnest): [15] detected fraudulent reviewers by combining the rating and temporal information and generating a likelihood-based suspiciousness metric for users and reviews. We adopted Birdnest to rank users and tweets.

The **supervised approaches** are as follows. (i) **Baseline IV** (SpamBot): [33] detected spam bots using a set of content-based and graph-based features. The proposed set of features is used in Naive Bayes classifier to classify a user into collusive or genuine. (ii) **Baseline V** (FakeAcc): [9] used a classification method to detect fake accounts on Twitter based on minimum weighted feature set calculated over 22 features. The feature set is then applied on several classifiers among which SVM performed the best. (iii) **Baseline VI** (SCoRe): [7] is the closest baseline of our method which identified blackmarket customers by running SVM using a set of 64 features.

### 5.2  Annotating Unknown Users

We asked three human annotators[3] to label $7,451$ unknown users $\mathcal{S}_u$ (mentioned in Section 3) into *collusive* or *genuine*. Annotators were given the definition of collusive users and Twitter Terms of Service. They were also given complete freedom to search for more information related to the (collusive) users on the web and apply their intuitions. The following guidelines were also given to them:

(1) Collusive users are members of blackmarket services who retweet/quote the tweets submitted to blackmarket services to earn credits in return.

(2) To earn credits from the blackmarket services, collusive users tend to show more aggressive behavior in terms of retweeting activity. Moreover, collusive users are less likely to retweet content of their friends.

(3) Generally, tweets submitted to the blackmarket services are related to promotions of tweets, accounts or services [8].

(4) Users created by the blackmarket services will only be involved in retweeting other tweets rather than publishing their own tweets. The annotators were asked to check the retweet statistics such as number of retweets, inter-retweet times, etc.

---

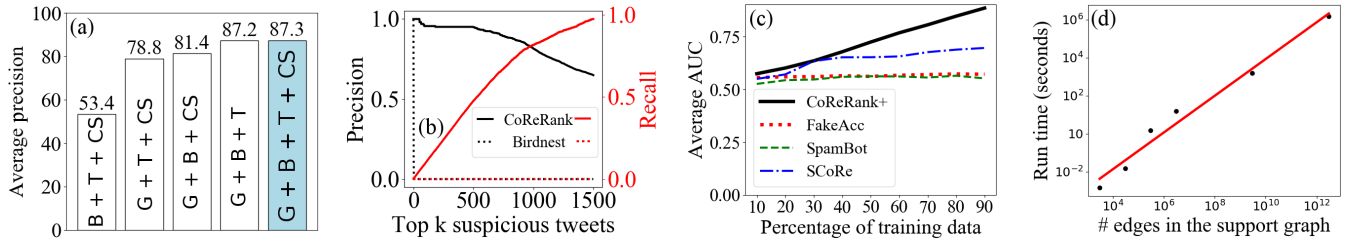[3]Annotators were experts in social media, and their age ranged between 25-35.

Figure 3: (a) Importance of different components of the recurrence formulation – graph (G), behavior (B), topic (T) and cold start (CS). (b) Precision and recall of `CoReRank` and `Birdnest` in detecting suspicious tweets. (c) Average AUC with different percentage of training data. (d) Scalability analysis of `CoReRank`.

We found high inter-annotator agreement (Fleiss' kappa coefficient of 0.75). We finally considered 4732 users as collusive and 2719 users as genuine (for which at least two annotators agreed). The newly labeled sets were augmented with the sets of ground-truth users $S_c$ and $S_g$ (as described in Section 3).

## 5.3 Comparative Evaluation

We observe the performance of the competing unsupervised methods with the increase of $k$, number of results returned. Figure 2 shows that till $k = 30$ ($k = 60$), precision of `CoReRank` is almost 1 for collusive (genuine) user detection; even after this, the decrease in performance is significantly less for `CoReRank` compared to other methods. The recall curve corresponding to `CoReRank` also increases steadily with $k$ compared to other methods.

Table 3 shows the performance of unsupervised and (semi-) supervised methods separately. We observe that `CoReRank` beats the best baseline (`Birdnest`) with 269% high average precision and 300% high average recall for collusive user detection. For genuine user detection, its improvement is 20% and 22.22% higher than `Birdnest` based on average precision and average recall respectively. The performance of (semi-) supervised methods is reported separately after averaging over 10-fold cross validation. For `CoReRank+`, we choose 0.202 as the threshold of $C(u)$ [4]; users whose corresponding $C(u)$ values are higher than (lower than or equal to) the threshold are considered as genuine (collusive). We observe that `CoReRank+` beats the best baseline (SCoRe) by 33.18% in terms of AUC (similar pattern is observed for other evaluation metrics).

**Robustness Analysis:** We also report the performance of the competing (semi-)supervised methods with the increase of training size to show the robustness of the methods w.r.t the size of the training set. We vary the training data from 10% to 90%. Figure 3(c) shows the AUC on test sets, averaged over 50 random iterations of training data. We observe that `CoReRank+` always outperforms others for all training percentage. This shows the efficiency of our framework even when a small amount of training data is available.

## 5.4 Importance of Different Components

In order to understand the importance of different components in the recurrence formulation of `CoReRank`, we drop each component in isolation (set its corresponding coefficient as 0) and measure

Table 3: Performance of the competing methods. We show the accuracy separately for unsupervised (in terms of Average Precision (AP) and Average Recall (AR)) and (semi-) supervised (in terms of ROC-AUC, Precision (P) and Recall (R), averaged over 10-fold cross validation) methods.

| Metric | BotoM | NDSync | Birdnest | CoReRank |
|---|---|---|---|---|
| AP (Collusive) | 0.212 | 0.218 | 0.221 | **0.817** |
| AR (Collusive) | 0.006 | 0.010 | 0.003 | **0.012** |
| AP (Genuine) | 0.394 | 0.532 | 0.727 | **0.875** |
| AR (Genuine) | 0.005 | 0.009 | 0.009 | **0.011** |
| Metric | SpamBot | FakeAcc | SCoRe | CoReRank+ |
| AUC | 0.573 | 0.578 | 0.696 | **0.927** |
| P (Collusive) | 0.611 | 0.589 | 0.724 | **0.933** |
| R (Collusive) | 0.336 | 0.362 | 0.727 | **0.887** |
| P (Genuine) | 0.547 | 0.559 | 0.603 | **0.910** |
| R (Genuine) | 0.790 | 0.762 | 0.599 | **0.947** |

the accuracy. Figure 3(a) shows maximum decrease in accuracy (38%) after removing graph (G) component, followed by behavior (9%), topic (6.75%) and cold start (0.11%). However, removing none of the components increases the accuracy, indicating that all the components are important.

## 5.5 Suspicious Tweet Detection

One of the advantages of `CoReRank` is that apart from ranking users based on their collusive activities, it can also rank tweets based on the merit score, which none of the baselines (except `Birdnest`) can. We take 1001 tweets which we collected from the blackmarket services as ground-truth suspicious tweets [5]. Figure 3(b) shows the precision and recall of `CoReRank` and `Birdnest` with the number of tweets returned ($k$). `CoReRank` achieves 0.85 average precision and 0.60 average recall, whereas both values for `Birdnest` is 0. In the ranked list that `Birdnest` returns, the first suspicious tweet appears at $3485^{th}$ position.

## 5.6 Scalability Analysis

Theorem 4.7 already showed that the running time of `CoReRank` scales linearly in the number of edges in the bipartite support graph. To show this empirically, we take the complete bipartite graph and keep adding edges (within the range of $10^4 - 10^{12}$) randomly without changing the number of nodes. Figure 3(d) shows that the running

---

[4]We plot the non-cumulative distribution of $C(u)$ and choose the threshold based on sudden change in the slope of the curve; see Supplementary [26].

[5]We did not show the results for genuine tweet detection as we were not sure about the ground-truth genuine tweets.

time increases linearly with the number of edges. On the dataset we collected, the average running time over 50 iterations is 170 seconds, which is much faster than BotoM (14400 seconds), NDSync (3000 seconds) and Birdnest (1200 seconds)[6]. We implemented CoReRank in Python using the Pandas library [22] for efficiency. All experiments were executed on a 1.8 GHz Intel Core i5 Macbook Air, 8 GB DDR3 RAM, running macOS High Sierra v10.13.3.

## 6 CONCLUSION

We presented CoReRank, the first unsupervised framework to simultaneously detect collusive users and suspicious tweets, controlled by blackmarket retweeting services. We showed the effectiveness of our framework over six other baselines on our own (manually curated and annotated) dataset. We also provided theoretical guarantees to show the convergence of CoReRank. Its runtime scales linearly with the number of edges, and it is faster than other competing methods. The dataset we collected is also the first dataset of this kind. We also made the code and dataset available for the purpose of reproducibility.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Faraz Ahmed and Muhammad Abulaish. 2013. A generic statistical approach for spam detection in online social networks. *Computer Communications* 36, 10-11 (2013), 1120–1129.
[2] Fabricio Benevenuto, Gabriel Magno, Tiago Rodrigues, and Virgilio Almeida. 2010. Detecting spammers on twitter. In *Collaboration, electronic messaging, anti-abuse and spam conference (CEAS)*, Vol. 6. 12.
[3] Zi Chu, Steven Gianvecchio, Haining Wang, and Sushil Jajodia. 2010. Who is tweeting on Twitter: human, bot, or cyborg?. In *Proceedings of the 26th annual computer security applications conference*. ACM, 21–30.
[4] Stefano Cresci, Roberto Di Pietro, Marinella Petrocchi, Angelo Spognardi, and Maurizio Tesconi. 2015. Fame for sale: efficient detection of fake Twitter followers. *Decision Support Systems* 80 (2015), 56–71.
[5] Clayton Allen Davis, Onur Varol, Emilio Ferrara, Alessandro Flammini, and Filippo Menczer. 2016. Botornot: A system to evaluate social bots. In *Proceedings of the 25th International Conference Companion on World Wide Web*. International World Wide Web Conferences Steering Committee, 273–274.
[6] Carlo De Micheli and Andrea Stroppa. 2013. Twitter and the underground market. In *11th Nexa Lunch Seminar*, Vol. 22.
[7] Hridoy Dutta, Aditya Chetan, Brihi Joshi, and Tanmoy Chakraborty. 2018. Retweet Us, we will Retweet You: Spotting Collusive Retweeters Involved in Blackmarket Services. In *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*. 242–249.
[8] Milad Eftekhar and Nick Koudas. 2013. Some Research Opportunities on Twitter Advertising. *IEEE Data Eng. Bull.* 36, 3 (2013), 77–82.
[9] Ahmed ElAzab. 2016. Fake Accounts Detection in Twitter based on Minimum Weighted Feature. *World* (2016).
[10] Hongyu Gao, Yan Chen, Kathy Lee, Diana Palsetia, and Alok N Choudhary. 2012. Towards Online Spam Filtering in Social Networks.. In *NDSS*, Vol. 12. 1–16.
[11] Maria Giatsoglou, Despoina Chatzakou, Neil Shah, Alex Beutel, Christos Faloutsos, and Athena Vakali. 2015. Nd-sync: Detecting synchronized fraud activities. In *PAKDD*. 201–214.
[12] Maria Giatsoglou, Despoina Chatzakou, Neil Shah, Christos Faloutsos, and Athena Vakali. 2015. Retweeting activity on twitter: Signs of deception. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 122–134.
[13] Chris Grier, Kurt Thomas, Vern Paxson, and Michael Zhang. 2010. @ spam: the underground on 140 characters or less. In *Proceedings of the 17th ACM conference on Computer and communications security*. ACM, 27–37.
[14] Srishti Gupta, Abhinav Khattar, Arpit Gogia, Ponnurangam Kumaraguru, and Tanmoy Chakraborty. 2018. Collective Classification of Spam Campaigners on Twitter: A Hierarchical Meta-Path Based Approach. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*. 529–538.
[15] Bryan Hooi, Neil Shah, Alex Beutel, Stephan Günnemann, Leman Akoglu, Mohit Kumar, Disha Makhija, and Christos Faloutsos. 2016. Birdnest: Bayesian inference for ratings-fraud detection. In *Proceedings of the 2016 SIAM International Conference on Data Mining*. SIAM, 495–503.
[16] Srijan Kumar, Bryan Hooi, Disha Makhija, Mohit Kumar, Christos Faloutsos, and V.S. Subrahmanian. 2018. REV2: Fraudulent User Prediction in Rating Platforms. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (WSDM '18)*. ACM, New York, NY, USA, 333–341. DOI:http://dx.doi.org/10.1145/3159652.3159729
[17] Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. 2010. What is Twitter, a social network or a news media?. In *Proceedings of the 19th international conference on World wide web*. ACM, 591–600.
[18] Eric Lancaster, Tanmoy Chakraborty, and VS Subrahmanian. 2018. MALTP: Parallel Prediction of Malicious Tweets. *IEEE Transactions on Computational Social Systems* (2018).
[19] Kyumin Lee, James Caverlee, and Steve Webb. 2010. Uncovering social spammers: social honeypots+ machine learning. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 435–442.
[20] Yuli Liu, Yiqun Liu, Min Zhang, and Shaoping Ma. 2016. Pay Me and I'll Follow You: Detection of Crowdturfing Following Activities in Microblog Environment.. In *IJCAI*. 3789–3796.
[21] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *JMLR* 9, Nov (2008), 2579–2605.
[22] Wes McKinney. 2011. pandas: a foundational Python library for data analysis and statistics. *Python for High Performance and Scientific Computing* (2011), 1–9.
[23] Zachary Miller, Brian Dickinson, William Deitrick, Wei Hu, and Alex Hai Wang. 2014. Twitter spammer detection using data stream clustering. *Information Sciences* 260 (2014), 64–73.
[24] Marti Motoyama, Damon McCoy, Kirill Levchenko, Stefan Savage, and Geoffrey M Voelker. 2011. An analysis of underground forums. In *Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference*. ACM, 71–80.
[25] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
[26] CoReRank Online repository: code & data & Supplementary. 2018. https://github.com/LCS2-IIITD/CoReRank-WSDM-2019. (2018).
[27] Neil Shah. 2017. FLOCK: Combating Astroturfing on Livestreaming Platforms. In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 1083–1091.
[28] Neil Shah, Hemank Lamba, Alex Beutel, and Christos Faloutsos. 2017. The Many Faces of Link Fraud. In *Data Mining (ICDM), 2017 IEEE International Conference on*. IEEE, 1069–1074.
[29] Monika Singh, Divya Bansal, and Sanjeev Sofat. 2016. Followers or fradulents? An analysis and classification of Twitter followers market merchants. *Cybernetics and Systems* 47, 8 (2016), 674–689.
[30] Gianluca Stringhini, Gang Wang, Manuel Egele, Christopher Kruegel, Giovanni Vigna, Haitao Zheng, and Ben Y Zhao. 2013. Follow the green: growth and dynamics in twitter follower markets. In *Proceedings of the 2013 conference on Internet measurement conference*. ACM, 163–176.
[31] Kurt Thomas, Damon McCoy, Chris Grier, Alek Kolcz, and Vern Paxson. 2013. Trafficking Fraudulent Accounts: The Role of the Underground Market in Twitter Spam and Abuse.. In *USENIX Security Symposium*. 195–210.
[32] Sokratis Vidros, Constantinos Kolias, Georgios Kambourakis, and Leman Akoglu. 2017. Automatic detection of online recruitment frauds: Characteristics, methods, and a public dataset. *Future Internet* 9, 1 (2017), 6.
[33] Alex Hai Wang. 2010. Detecting spam bots in online social networking sites: a machine learning approach. In *IFIP Annual Conference on Data and Applications Security and Privacy*. Springer, 335–342.
[34] Bo Wang, Arkaitz Zubiaga, Maria Liakata, and Rob Procter. 2015. Making the most of tweet-inherent features for social spam detection on twitter. *arXiv preprint arXiv:1503.07405* (2015).
[35] De Wang, Shamkant B Navathe, Ling Liu, Danesh Irani, Acar Tamersoy, and Calton Pu. 2013. Click traffic analysis of short url spam on twitter. In *Collaborative Computing: Networking, Applications and Worksharing (Collaboratecom), 2013 9th International Conference on*. IEEE, 250–259.
[36] Chao Yang, Robert Harkreader, Jialong Zhang, Seungwon Shin, and Guofei Gu. 2012. Analyzing spammers' social networks for fun and profit: a case study of cyber criminal ecosystem on twitter. In *Proceedings of the 21st international conference on World Wide Web*. ACM, 71–80.

---

[6]We did not report the running time of supervised methods as it may not be appropriate to compare two different classes of methods.