# A Survey and Critique of Recent Advance in Online Spam Detection

Yingtong Dou

Department of Computer Science
University of Illinois at Chicago

March 9, 2019

# Road Map

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ●○○○ | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Background

## What is spam?

- Spam is the **fake** and **useless** information on the Internet.
- Spam is composed by **intentionally crafted** content.
- Spam is everywhere.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○●○○ | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Background

# Examples of Spams

Amazon

Twitter

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○●○ | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Background

## Effect of Spams

(UIC) COMPUTER SCIENCE

- ▶ Decay online experience.
- ▶ Bias users' choice/opinion.
- ▶ Mislead recommender system.
- ▶ Usually come with other security threats.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○○● | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Background

## Spam Detection Problem

**UIC COMPUTER SCIENCE**

- ▶ Spam detection is an **Anomaly Detection** task in data mining.
- ▶ We need to evaluate the **suspiciousness** of users, posts, reviews.
- ▶ Generally, spam detection is a **Supervised Learning** task.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 00 |
| ●00 | 00 | 000 | 0000 | 00000 |
| 0 | 00 | 00 | 0 | |

Methods

## Feature Extraction

UIC **COMPUTER SCIENCE**

- ▶ Semantic Features

| Feature Name | Description |
|---|---|
| RL | Average review length |
| ACS | Average content similarity |
| PCW | Percentage of all capital words |
| PC | Percentage of capital letters |
| $DL_b$ | Description length based on bigrams |
| PP1 | The ratio of 1st person pronouns |
| RES | The ratio of exclamation sentences |
| SW | The ratio of subjective words |
| OW | The ratio of objective words |
| F | The frequency of review |

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 00 |
| 0●0 | 00 | 000 | 0000 | 00000 |
| 0 | 00 | 00 | 0 | |

Methods

## Feature Extraction



- Behavioral Features

| Behavior | Description |
|---|---|
| MNR | Max. number of reviews posted in a day |
| PR | The ratio of positive reviews |
| NR | The ratio of negative reviews |
| WRD | Weighted Rating Deviation |
| BST | Burstiness |
| RD | Rating deviation of product's avg. rating |
| Rank | The rank order of the review |
| ETF | The early time frame of the reviewer |
| ISR | Is singleton? |
| DPW | Deceptive review count previous week |

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○○○ | ○○ | ○ | ○○ | ○○ |
| ○○● | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Methods

## Learning Methods

- ▶ **Traditional Model**
  Naive Bayes, Support Vector Machine, Random Forest etc.

- ▶ **Graph Model**

- ▶ **Deep Model**
  Relation Embedding, Convolutional Neural Network etc.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| oooo | oo | o | oo | oo |
| ooo | oo | ooo | oooo | ooooo |
| ● | oo | oo | o | |

Overview

## Overview of the Three Paper


UIC COMPUTER SCIENCE

| Paper | Problem | Model | Target | Dataset | Venue |
|---|---|---|---|---|---|
| Paper 1 | Optimizing Alg | Graph | Social Spammer | Twitter | ICDM2017 |
| Paper 2 | Cold Start | Deep Model | Spam Review | Yelp | ACL2017 |
| Paper 3 | Crowdsourcing | Graph&Deep | Crowd Worker | Amazon | WSDM2018 |

# GANG: Detecting Fraudulent Users in Online Social Networks via Guilt-by-Association on Directed Graphs

Binghui Wang, Neil Zhenqiang Gong
ECE Department, Iowa State University
{binghuiw, neilgong}@iastate.edu

Hao Fu
Microsoft Research Asia, China
fuha@microsoft.com

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | ●○ | ○ | ○○ | ○○ |
| 000 | ○○ | 000 | 0000 | 00000 |
| ○ | ○○ | ○○ | ○ | |

Directed Graph Model

## Graph Homophily Assumption

UIC COMPUTER SCIENCE

▶ Similar nodes are more likely to connect with each other than dissimilar ones.

▶ In the online social network, we represent **users as nodes**, their **friendship as edges**.

▶ **Suspicious users** tend to connect with each other; **regular users** tend to connect with each other.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○○○ | ○● | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Directed Graph Model

# Modeling Directed Edge Influence

**UIC COMPUTER SCIENCE**

**Bidirectional Edge**
$v_1(Benign) \Rightarrow u(Benign)$
$v_1(Suspicious) \Rightarrow u(Suspicious)$

**Unidirectional Incoming Edge**
$v_2(Benign) \Rightarrow u(Benign)$
$v_2(Suspicious) \Rightarrow u(?)$

**Unidirectional Outgoing Edge**
$v_3(Benign) \Rightarrow u(?)$
$v_3(Suspicious) \Rightarrow u(Suspicious)$

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 00 |
| 000 | ●0 | 000 | 0000 | 00000 |
| 0 | 00 | 00 | 0 | |

Belief Propagation

## Inference on MRF

UIC COMPUTER SCIENCE

- ▶ The social network could be modeled as a Markov Random Field.
- ▶ Belief propagation could infer the states of the nodes in MRF.



Message Passing Process of Belief Propagation

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
| --- | --- | --- | --- | --- |
| oooo | oo | o | oo | oo |
| ooo | o● | ooo | oooo | ooooo |
| o | oo | oo | o | |

Belief Propagation

# Optimizing Belief Propagation

▶ Eliminate message maintenance in BP.

▶ Approximate BP using matrix multiplication.



One Round of Approximated BP in Matrix Form

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 00 |
| 000 | 00 | 000 | 0000 | 00000 |
| 0 | ●0 | 00 | 0 | |

Critique

## Contributions of the Paper

**UIC COMPUTER SCIENCE**

- ▶ Models the influence of directed edges with a unified formulation.
- ▶ Proves the convergence condition of the proposed model.
- ▶ Strong theoretical guarantee.
- ▶ Algorithm is scalable .

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○○○ | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○○○ | ○○○○ | ○○○○○ |
| ○ | ○● | ○○ | ○ | |

Critique

## Drawbacks of the Paper

UIC COMPUTER SCIENCE

- ▶ Model is vulnerable to attacks.
- ▶ It is a trade off between model robustness and model performance

# Handling Cold-Start Problem in Review Spam Detection by Jointly Embedding Texts and Behaviors

**Xuepeng Wang**[1,2], **Kang Liu**[1], and **Jun Zhao**[1,2]

[1] National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, 100190, China

[2] University of Chinese Academy of Sciences, Beijing, 100049, China

{xpwang, kliu, jzhao}@nlpr.ia.ac.cn

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| oooo | oo | ● | oo | oo |
| ooo | oo | ooo | oooo | ooooo |
| o | oo | oo | o | |

Cold Start

## Cold Start Problem

UIC COMPUTER SCIENCE

- ▶ Cold Start refers to those new coming data items.
- ▶ Traditional features fail to model new users and reviews.
- ▶ The graph model is not useful for dealing with new users.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○○○ | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ●○○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Review Relation Embedding

## Translating Model

▶ Translating is a relation embedding model on knowledge graph.



$$\sum_{(h,\ell,t),(h',\ell,t'))\in T_{\text{batch}}} \nabla\big[\gamma + d(h+\ell,t) - d(h'+\ell,t')\big]_+$$

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| ○○○○ | ○○ | ○ | ○○ | ○○ |
| ○○○ | ○○ | ○●○ | ○○○○ | ○○○○○ |
| ○ | ○○ | ○○ | ○ | |

Review Relation Embedding

# Review Triplet

- Each review could be represented as a user-review-product triplet.



Relations in a Review Platform

# Review Relation Embedding

- Product as head, reviewers as relation, review embedding as tail.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 00 |
| 000 | 00 | 000 | 0000 | 00000 |
| 0 | 00 | ●0 | 0 | |

Critique

## Contributions

UIC COMPUTER SCIENCE

- ▶ First paper tackling cold start in spam detection.
- ▶ Encode latent relations with deep network.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
| --- | --- | --- | --- | --- |
| 0000 | 00 | 0 | 00 | 00 |
| 000 | 00 | 000 | 0000 | 00000 |
| 0 | 00 | 0● | 0 | |

Critique

## Drawbacks

UIC COMPUTER SCIENCE

- ▶ No explanation for review triplet setting.
- ▶ No comparison with other embedding models.
- ▶ No comparison with other dimension reduction models.

# Combating Crowdsourced Review Manipulators: A Neighborhood-Based Approach

Parisa Kaghazgaran
Texas A& M University
College Station, TX
kaghazgaran@tamu.edu

James Caverlee
Texas A& M University
College Station, TX
caverlee@tamu.edu

Anna Squicciarini
Pennsylvania State University
State College, PA
asquicciarini@ist.psu.edu

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | OO | O | ●O | OO |
| 000 | OO | 000 | 0000 | 00000 |
| O | OO | OO | O | |

Crowdsourcing Attack

## Crowdsourcing Spam

**UIC COMPUTER SCIENCE**

▶ Crowdsourcing is an activity that hires online freelancing workers to finish specific tasks.

| Available Tasks | Amount | Time |
|---|---|---|
| Youtube: Vote for this video | $0.10 | 1 min |
| Follow me on Twitter | $0.12 | 1 min |
| Insurance Form: Sign up | $1.50 | 5 min |
| Create Gmail account for me | $0.13 | 3 min |
| Online Game: Sign up | $0.20 | 3 min |
| Digg: Bookmark my page | $0.10 | 1 min |
| Upload 5 photos to this site | $0.39 | 4 min |

A Screenshot of Crowdsourcing Tasks Listed on a Website

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| oooo | oo | o | o● | oo |
| ooo | oo | ooo | oooo | ooooo |
| o | oo | oo | o | |

Crowdsourcing Attack

# Challenges in Defending Crowdsourcing Attack (UIC) COMPUTER SCIENCE

- ▶ Crowd workers look like regular users.
- ▶ It is difficult to acquire the ground truth.

Introduction  Paper1: Alg Optimizing  Paper2: Cold Start  Paper3: Crowdsourcing Attack  Discussion & Future Work
0000          oo                       o                  oo                              oo
ooo           oo                       ooo                ●ooo                            ooooo
o             oo                       oo                 o

TwoFace Framework

## Select Seed Users

UIC COMPUTER SCIENCE
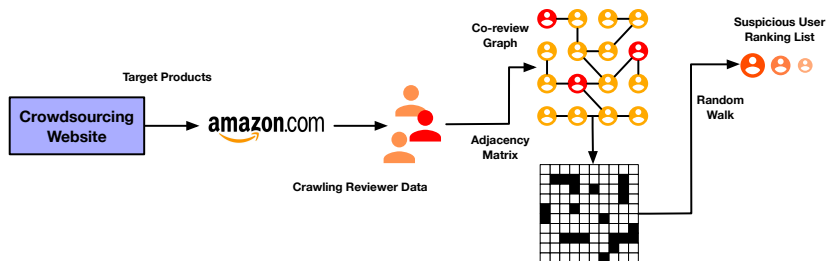
▶ Crawl all the Amazon products that have released tasks on a crowdsourcing website.

▶ Crawl all the reviewers having reviewed those products.

▶ Select seed users.



Target Products

Crowdsourcing Website → amazon.com →

Crawling Reviewer Data

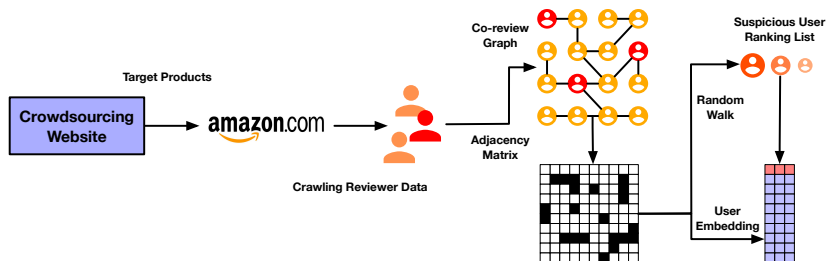| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
| 0000 | 00 | 0 | 00 | 00 |
| 000 | 00 | 000 | 0●00 | 00000 |
| 0 | 00 | 00 | 0 | |

TwoFace Framework

# Discover Local Similar Users

**UIC COMPUTER SCIENCE**

- Construct a co-review graph where reviewers are nodes. Edges connect users who both have reviewed the same product.
- Set the suspicious score of seed users to 1.
- Using random walk to propagate the suspiciousness.

# Discover Distant Similar Users

UIC COMPUTER SCIENCE
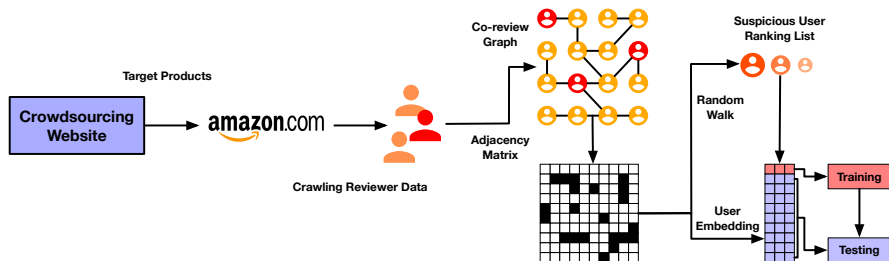
▶ Random walk only discovers local similar users. It cannot pass suspiciousness to users having no connection with seed users.

▶ Use **node2vec** model to learn the node embedding.

▶ Use node embedding to find structural similar users.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 00 |
| 000 | 00 | 000 | 000● | 00000 |
| 0 | 00 | 00 | 0 | |

TwoFace Framework

# TwoFace Framework

▶ Select seed users; use random walk to generate suspicious user ranking list; learn node embeddings of all users.

▶ Train traditional classifiers with node embeddings.

▶ Validate the model with holdout data.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| oooo | oo | o | oo | oo |
| ooo | oo | ooo | oooo | ooooo |
| o | oo | oo | ● | |

Model Limitations

## Model Limitations

**UIC COMPUTER SCIENCE**

- ▶ No side information of users.
- ▶ No comparison with other embedding models.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| oooo | oo | o | oo | ●o |
| ooo | oo | ooo | oooo | ooooo |
| o | oo | oo | o | |

Discussion

## Summary of Three Papers

**UIC COMPUTER SCIENCE**

- ▶ Traditional features and models have limitations.
- ▶ Graph models have strong theoretical guarantee.
- ▶ Deep models have weak interpretability.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| 0000 | 00 | 0 | 00 | 0● |
| 000 | 00 | 000 | 0000 | 00000 |
| 0 | 00 | 00 | 0 | |

Discussion

# Challenges in Spam Detection Research

**UIC COMPUTER SCIENCE**

- ▶ Vulnerability to attacks.
- ▶ Reproducibility of deep model.
- ▶ Lack of theoretical guarantee.
- ▶ Quality of benchmark datasets.
- ▶ Practical performance of detectors.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
| :--- | :--- | :--- | :--- | :--- |
| 0000 | 00 | 0 | 00 | 00 |
| 000 | 00 | 000 | 0000 | ●0000 |
| 0 | 00 | 00 | 0 | |

Future Work

# Adapt More Deep Models

UIC COMPUTER SCIENCE

- ▶ Graph Convolutional Network.
- ▶ Long Short Term Memory Network.
- ▶ Auto-Encoder Framework.

| Introduction | Paper1: Alg Optimizing | Paper2: Cold Start | Paper3: Crowdsourcing Attack | Discussion & Future Work |
|---|---|---|---|---|
| oooo | oo | o | oo | oo |
| ooo | oo | ooo | oooo | o●ooo |
| o | oo | oo | o | |

Future Work

## Other Promising Research Directions

- ▶ Adversarial machine learning.

- ▶ Dynamic detection model.

- ▶ Heterogeneous information network.

- ▶ New problems:
    - ▶ Poisoning reviews
    - ▶ Fake news
    - ▶ Multi-intention reviews
    - ▶ Machine-generated content

Formulation of the proposed model in Paper 1.

$$\begin{cases} \mathbf{A}_i'^{(t-1)} = I\left(\mathbf{A}_i \circ \mathbf{P}^{(t-1)^T}\right) \\ \mathbf{A}_o'^{(t-1)} = I\left(-\mathbf{A}_o \circ \mathbf{P}^{(t-1)^T}\right) \\ \mathbf{p}^{(t)} = \mathbf{q} + 2 \cdot w \cdot \left(\mathbf{A}_b + \mathbf{A}_i'^{(t-1)} + \mathbf{A}_o'^{(t-1)}\right) \cdot \mathbf{p}^{(t-1)} \end{cases}$$

$A_i$: incoming edge adjacency matrix  $\qquad$  $P$: node posterior belief

$A_o$: outgoing edge adjacency matrix  $\qquad$  $q$: node prior belief

$A_b$: bidirectional edge adjacency matrix  $\qquad$  $w$: coupling strength

---

**Algorithm 1** Learning TransE

---

**input** Training set $S = \{(h, \ell, t)\}$, entities and rel. sets $E$ and $L$, margin $\gamma$, embeddings dim. $k$.

1: **initialize** $\boldsymbol{\ell} \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$ for each $\ell \in L$

2: $\qquad \boldsymbol{\ell} \leftarrow \boldsymbol{\ell}/\left\|\boldsymbol{\ell}\right\|$ for each $\ell \in L$

3: $\qquad \mathbf{e} \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$ for each entity $e \in E$

4: **loop**

5: $\qquad \mathbf{e} \leftarrow \mathbf{e}/\left\|\mathbf{e}\right\|$ for each entity $e \in E$

6: $\qquad S_{batch} \leftarrow \text{sample}(S, b)$ // sample a minibatch of size $b$

7: $\qquad T_{batch} \leftarrow \emptyset$ // initialize the set of pairs of triplets

8: $\qquad$ **for** $(h, \ell, t) \in S_{batch}$ **do**

9: $\qquad\qquad (h', \ell, t') \leftarrow \text{sample}(S'_{(h,\ell,t)})$ // sample a corrupted triplet

10: $\qquad\qquad T_{batch} \leftarrow T_{batch} \cup \left\{\left((h, \ell, t), (h', \ell, t')\right)\right\}$

11: $\qquad$ **end for**

12: $\qquad$ Update embeddings w.r.t. $\displaystyle\sum_{\left((h,\ell,t),(h',\ell,t')\right)\in T_{batch}} \nabla\left[\gamma + d(\boldsymbol{h} + \boldsymbol{\ell}, \boldsymbol{t}) - d(\boldsymbol{h'} + \boldsymbol{\ell}, \boldsymbol{t'})\right]_{+}$

13: **end loop**

---

# Reference

1 Wang, Binghui, Neil Zhenqiang Gong, and Hao Fu. "GANG: Detecting fraudulent users in online social networks via guilt-by-association on directed graphs." In 2017 IEEE International Conference on Data Mining (ICDM), pp. 465-474. IEEE, 2017.

2 Wang, Xuepeng, Kang Liu, and Jun Zhao. "Handling cold-start problem in review spam detection by jointly embedding texts and behaviors." In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), vol. 1, pp. 366-376. 2017.

3 Kaghazgaran, Parisa, James Caverlee, and Anna Squicciarini. "Combating crowdsourced review manipulators: A neighborhood-based approach." In Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, pp. 306-314. ACM, 2018.

4 Rayana, Shebuti, and Leman Akoglu. "Collective opinion spam detection: Bridging review networks and metadata." In Proceedings of the 21th acm sigkdd international conference on knowledge discovery and data mining, pp. 985-994. ACM, 2015.

5 Mukherjee, Arjun, Vivek Venkataraman, Bing Liu, and Natalie Glance. "What yelp fake review filter might be doing?." In Seventh international AAAI conference on weblogs and social media. 2013.