

Opinion Spam Detection: A New Approach Using Machine Learning and Network-Based Algorithms

Kiril Danilchenko*, Michael Segal and Dan Vilenchik

School of Electrical and Computer Engineering
 Ben-Gurion University of the Negev
 Beer-Sheva, Israel 8410501
 vilenchi@bgu.ac.il

Abstract

E-commerce is the fastest-growing segment of the economy. Online reviews play a crucial role in helping consumers evaluate and compare products and services. As a result, fake reviews (opinion spam) are becoming more prevalent and negatively impacting customers and service providers. There are many reasons why it is hard to identify opinion spammers automatically, including the absence of reliable labeled data. This limitation precludes an off-the-shelf application of a machine learning pipeline. We propose a new method for classifying reviewers as spammers or benign, combining machine learning with a message-passing algorithm that capitalizes on the users' graph structure to compensate for the possible scarcity of labeled data. We devise a new way of sampling the labels for the training step (active learning), replacing the typical uniform sampling. Experiments on three large real-world datasets from Yelp.com show that our method outperforms state-of-the-art active learning approaches and also machine learning methods that use a much larger set of labeled data for training.

Introduction

In the era of e-commerce, consumers typically buy products or services based on reviews. Therefore reviews are increasingly valuable for sellers and service providers due to the benefits of positive reviews or the damage from negative ones. In light of this, fake reviews are flourishing and pose a real threat to the proper conduct of e-commerce platforms. A group of reviewers (we will call them opinion spammers) post fake reviews to promote their products or demote their competitors' products.

Fake reviews are often written by experienced professionals who are paid to write high-quality, believable reviews. Detecting opinion fraud is a non-trivial and challenging problem that was extensively studied in the literature using various approaches. Some approaches are based solely on the review content [Jindal and Liu 2008; Li, Cardie, and Li 2013; Ott et al. 2011], reviewer behavior [Lim et al. 2010; Mukherjee et al. 2013b] and the tripartite relationships between reviewers, reviews, and products [Mukherjee, Liu, and Glance 2012; Rayana and Akoglu 2015, 2016; Wang

et al. 2011, 2018]. While each paper presented a method that is useful to some extent for detecting certain kinds of spamming activities, there is no one-size-fits-all solution. This is because spammers keep changing their strategies, many times in an adaptive manner to the spam detection policies. Therefore, there is a need to study and incorporate as many approaches as possible.

Another challenge is the fact that most datasets are imbalanced, as the three datasets that we used for evaluation, where only about 20% of the users are spammers (Table 1).

Despite the vast commercial impact of opinion spam detection, most machine-learning-based solutions for this problem do not achieve very high performance due to insufficient labeled data to properly train an ML model. In addition, standard ML models often treat each sample separately, disregarding the underlying graph structure of the spammer group. Indeed, this graph is often latent and should be derived from the data in an ad-hoc manner. One can use deep graph learning methodology to automatically embed the users, and thus infer the underlying graph structure, but such an approach requires a large training set, which is often not available or costly to obtain.

The setting where only part of the data is labeled is often called "semi-supervised learning". When very few examples are labeled, this setting is also termed "few-shot learning". When in addition one is given the possibility to choose which set of users will be labeled (given a budget, one can decide for which users to invest the budget in order to acquire their label), this is called the active-learning setting. In this paper we study the few-shot active learning setting for the opinion spam detection problem.

The few-shot active learning setting is very reasonable for the opinion-spam detection problem because labeled data never comes labeled for free, and operators of e-commerce websites choose de-facto which users they want to check manually and label, and typically only a small fraction of the users are labelled.

Our Contribution

We propose a classification algorithm, *Clique Reviewer Spammer Detection Network*, *CRSDnet*, for detecting fake reviews in the few-shot active learning setting. *CRSDnet* harnesses the power of both machine learning algorithms and classical graphical models algorithms such as

*Supported by ISF grant number 1388/16.
 Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

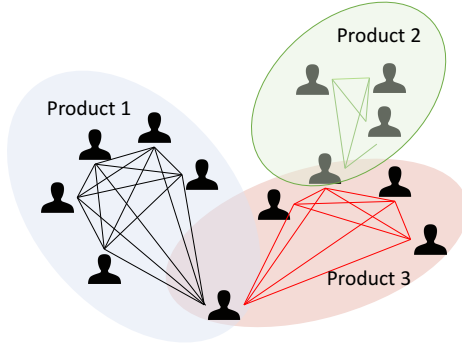


Figure 1: Illustration of the user-user graph, inducing an interconnected set of cliques, one for each product.

Belief Propagation. We show that this combination yields a better performance than each approach separately.

We evaluate our algorithm on a golden-standard dataset for the spammer detection task, the Yelp Challenge Data [Yelp 2014]. The performance of our algorithm is better than all previous work in almost every metric. We also outperform other methods that use the graph structure, such as graph embedding, and neural graph networks algorithms [Liu et al. 2020; Wang et al. 2019a; Shehneepoor et al. 2021]. These algorithms use much more labeled data for training (30% and over, compared to at most 2.5% that we use).

We show that using both machine learning and the graph structure (via label propagation algorithms) improves over stand-alone machine learning by at least 10% in AUC measure.

We attribute the success of our algorithm to the following key innovations in our approach:

- We derive from the raw data a user-user graph, where two users share an edge if they wrote a review for the same item. Figure 1 illustrates this graph, which is a clique graph by definition. In previous work a tripartite user-review-product graph was used.
- The user-user graph may be much denser than the tripartite user-review-product graph. To overcome computational issues that such density entails, we design a careful edge sparsification procedure to speed up the algorithm without compromising performance much. The sparsification is guided by the rule that each node will end up having just enough edges connecting it to nodes both from his own class (spammer or not) and from the opposite class.
- We run a label-propagation algorithm (concretely, Belief Propagation), but some parameters of that algorithm (the node and edge potentials) are determined using a machine learning model. This is the first time that such a combination of approaches is undertaken, and its usefulness demonstrated.
- We propose a new way of choosing the set of users whose

labels will be obtained (active learning). Instead of randomly choosing a set of users up to the allowed budget, we choose random users from the largest clique of the user-user graph. The intuition behind this rule comes from the work of Wang et. al. [Wang et al. 2020] where it was shown that collusive spamming (or, co-spamming) is a useful lens to identify spammers.

Related Work

Opinion spam detection has different nuances such as fake review detection [Jindal and Liu 2008; Ott et al. 2011; Xie et al. 2012; Xu et al. 2007], fake reviewer detection [Lim et al. 2010; Wang et al. 2011; Xu et al. 2013] and spammer group detection [Mukherjee, Liu, and Glance 2012; Ye and Akoglu 2015; Wang et al. 2020, 2018; Wang, Gu, and Xu 2018]. Two survey papers, [Crawford et al. 2015; Viviani and Pasi 2017], provide a broad perspective on the field.

Our method is part of the graph-based models, which take into account the relationships among reviewers, comments, and products. The key algorithm in this approach is Belief Propagation (BP) [Pearl 2014] which is applied to a carefully designed graph and the Markov Random Field (MRF) associated with it. The first to use this approach were Akoglu et al. [Akoglu, Chandy, and Faloutsos 2013] who suggested FraudEagle, a BP-based algorithm that runs on the bipartite reviewer-product graph, where the edge potentials are based on the sentiment in the review. In later work, Rayana et al. [Rayana and Akoglu 2015] introduced SPEagle, where node and edge potentials are derived from a richer set of meta-data features, improving significantly over the performance of FraudEagle. Wang et al. [Wang et al. 2011] consider the tripartite user-review-product network and define scores for trustiness of users, honesty of reviews, and reliability of products. They use an ad-hoc iterative procedure to compute the scores, rather than BP. In [Shehneepoor et al. 2017] an algorithm called NetSpam was introduced which utilizes spam features for modeling review datasets as heterogeneous information networks.

A graph-based approach was suggested in [Fei et al. 2013] but this time the graph contains edges for reviews that were written within a certain time difference from each other (a “burst”). The authors use a different dataset to evaluate their method, reviews from Amazon.com.

The authors of [Wang et al. 2020] present ColluEagle, a graph-based algorithm to detect both review spam campaigns and *collusive* review spammers. To measure the collusiveness of pairs of reviewers, they identify reviewers that review the same product in a similar way.

A different approach to the problem of spammer detection is via Graph Neural Networks (GNNs) and Graph Convolutional Networks (GCNs). These are deep learning architectures for graph-structured data. The core idea is to learn node representations through local neighborhoods [Kipf and Welling 2016]. In [Liu et al. 2020], the authors design a new GNN framework, GraphConsis, to tackle the fraud detection task. The authors evaluated the method on four data sets, where one of them is used by us too. GraphConsis is benchmarked on different training set sizes, from 40% to 80% of the data. GraphConsis with 80% training set size achieves

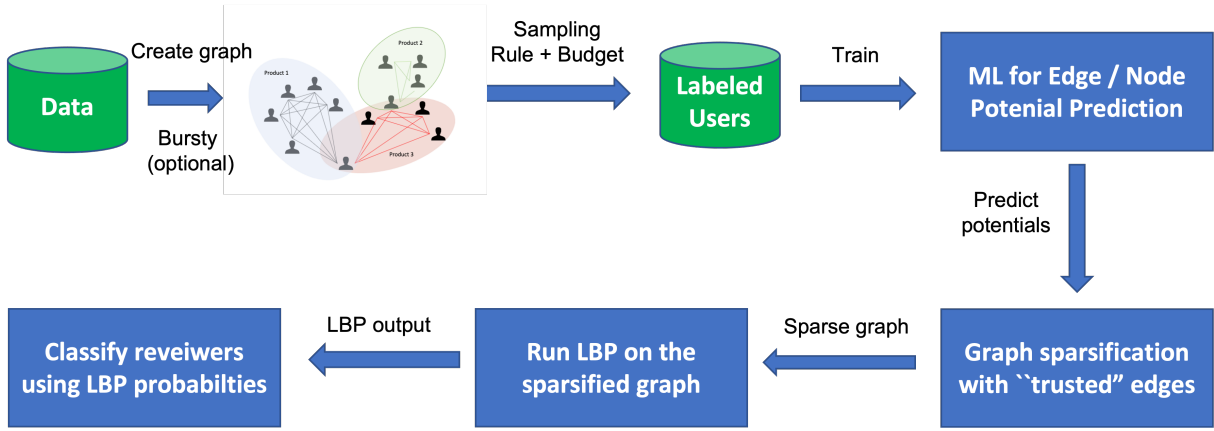


Figure 2: The flow chart of our pipeline, from raw data, through sampling users for labeling, using them to train an ML algorithm for predicting edge and node potentials, sparsification of the graph using trusted edges, running LBP, and completing the classification task.

AUC of 0.742 on the Yelp Chicago dataset, and our method achieves 0.754 with only 2.5%. Note that GraphConsis uses all the metadata that we use as well (the graph structure, and the reviews).

A GCN-based algorithm was designed in [Wang et al. 2019b], and tested on reviews from Tencent Inc. The algorithm outperformed four baseline algorithms (Logistic Regression, Random Forest, DeepWalk [Perozzi, Al-Rfou, and Skiena 2014], LINE [Tang et al. 2015]).

Our work departs from previous work in several ways. Compared to the works where label-propagation algorithms were used, we use machine learning to predict the edge and node potentials rather than hand-crafted threshold functions. Second, we consider the user-user graph and not a bi/tri-partite user-review-product graph. To overcome the computational challenge incurred by the density of the user-user graph, we apply a new rule for edge sparsification, based on the ML prediction. Finally, in the active learning setting, we introduce a new sampling rule. All these modifications have led to an improvement over FraudEagle [Akoglu, Chandy, and Faloutsos 2013] and SPEagle [Rayana and Akoglu 2015].

Active Learning: The active learning approach aims to achieve high accuracy by using few queries, and therefore the “most informative” points are natural candidates for label acquisition. Various heuristics were proposed to determine the “most informative” nodes, e.g., uncertainty sampling [Lewis and Catlett 1994; Culotta and McCallum 2005; Settles and Craven 2008] and variance reduction [Flaherty, Arkin, and Jordan 2006; Schein and Ungar 2007]. In our setting, we chose a rule that is native to the problem itself – sampling from the largest clique, following the take-home message from the work of [Wang et al. 2020] about collusive spamming.

Many works on active learning choose the train set using adaptive rules, point by point. This however is infeasible in our case as re-running the entire pipeline for every new example is computationally prohibitive for datasets as large as

ours. Therefore we choose all users for labelling in bulk.

Methodology

In this section we describe our pipeline, end to end. The flow chart is depicted in Figure 2.

We formulate the spam detection problem as a classification task on the user network. The dataset consists of n reviewers who write reviews on m products from the set P . The vertex set of the graph $G = (V, E)$ is the set of users (reviewers); user i and j share an undirected edge if there exists some product $p \in P$ such that user i and j wrote a review for p . The resulting graph consists of interconnected cliques, each corresponds to a different product. Figure 1 illustrates such a network. Each node $i \in V$ has in addition a vector of features F_i associated with it, and a binary class variable $v_i \in \{1, -1\}$, for spammer (1) or benign (-1).

The classification task is given the graph G (with the nodes’ features), and possibly a set of labeled users $\{i_1, \dots, i_k\}$ (the “few shots” training set), predict the value of v_i for the remaining nodes (the test set).

Ideally, to solve the classification task, we would find an assignment $s : \{-1, 1\}^n \rightarrow V$ that maximizes

$$Pr [v_1 = s_1, \dots, v_n = s_n | v_{i_1} = s_{i_1}, \dots, v_{i_k} = s_{i_k}, G]. \quad (1)$$

This Maximum Likelihood Estimation (MLE) task is in general NP-hard. However, in practice, a useful solution s (perhaps not the maximizer) may be obtained by using a Markov Random Field modelling for the probability space.

Markov Random Field

Markov Random Field (MRF) is often used to model a set of random variables having a Markov property described by an undirected dependency graph. A pairwise-MRF (pMRF) is an MRF satisfying the pairwise Markov property: a random variable depends only on its neighbors and is independent of all other variables.

A pMRF model involves two types of potentials, node potentials, and edge potentials. Our node potentials, $\phi_i(v_i)$, stand for the probability that reviewer i belongs to either class (spam/benign):

$$\phi_i(v_i) = \begin{cases} a_i & , v_i = 1 \text{ (} i \text{ is a spammer)} \\ 1 - a_i & , v_i = -1 \text{ (} i \text{ is a benign)}. \end{cases} \quad (2)$$

The edge potential $\psi_{ij}(v_i, v_j)$ signifies the affinity of reviewer i and j , namely, the probability p_{ij} that both belong to the same class. Formally,

$$\psi_{ij}(v_i, v_j) = \begin{cases} p_{ij} & , v_i = v_j \\ 1 - p_{ij} & , v_i \neq v_j. \end{cases} \quad (3)$$

The parameters a_i and p_{ij} satisfy $a_i, p_{ij} \in [0, 1]$ for all i, j . To determine the values of these parameters we use machine learning applied to features that are extracted from the meta-data of the reviewers dataset.

The pMRF model is used to approximate the expression for $Pr[s|G]$ in Eq. (1):

$$\Pr(s) = \frac{1}{Z} \prod_{v_i \in V} \phi_i(v_i) \prod_{(i,j) \in E} \psi_{ij}(v_i, v_j), \quad (4)$$

where Z is a normalization factor, the sum over the energies of all possible $2^{|V|}$ assignments s .

Finding the assignment s that maximizes the probability in Eq. (4) is still an intractable problem; LBP (loopy belief propagation) is the go-to heuristic for approximating the intractable maximization problem.

The LBP algorithm [Pearl 2014] is based on an iterative message passing along the edges of the graph. Messages are initialized according to some user-defined rule. At iteration t , a message $m_{ij}^{(t)}$ is sent from node i to each neighboring node j . The message represents the belief of i about the label of j . If G is a tree, then BP is guaranteed to converge; if G contains cycles then convergence is not guaranteed (hence the name loopy), but in practice, a cap on the number of iterations is set. We use the standard LBP messages, omitted for brevity, and can be found in [Pearl 2014].

Each iteration of LBP takes $O(|V| + |E|)$ time, hence the number of edges, which may be quadratic in $|V|$, plays a key role in the computational complexity. The more iterations one can perform for the same time budget, the better the performance. In the next section, we describe how to address the computational aspect using graph sparsification.

Running LBP on a Sparse Graph

Recall that our graph is defined over users, and not as a tripartite user-product-review graph. This may result in a rather dense graph, which poses a computational impediment even on LBP, when the number of nodes is large. For example, the graph created from the Yelp Chicago dataset is very dense (average degree 1193). Therefore our first step is to sparsify the graph by choosing a linear number of “useful” edges (linear in the number of nodes).

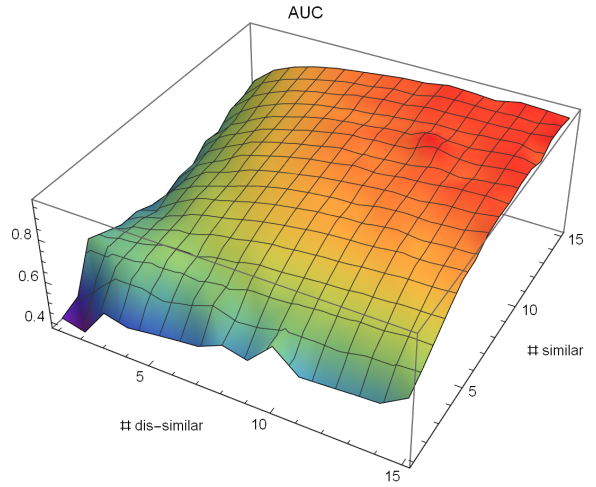


Figure 3: The AUC of the LBP algorithm for the Yelp Chicago dataset. LBP is run on a subgraph in which each node has k_1 neighbors from the similar class and k_2 neighbors from the other (dis-similar) class. Edge potentials are set according to the ground truth.

To gain intuition into a useful way of sparsification, we conducted the following experiment using the Chicago Yelp data [Rayana and Akoglu 2015]. The initial graph contains 38063 nodes, and, $2.4 \cdot 10^7$, edges. The sparsification procedure is parameterized with two numbers k_1 and k_2 . For every node i we choose k_1 neighbors from i ’s class (spammer or benign) and k_2 neighbors from the opposite class, and color these $k_1 + k_2$ edges red. We then remove all edges that were not colored red. The resulting graph has an average degree of at most $k_1 + k_2$ (multiple edges are merged). We set all node potentials to $\phi(v_i) \leftarrow \{0.5, 0.5\}$; in other words, we don’t provide any prior knowledge about the class of the node v_i . We set all edge potentials ψ_{ij} as follows: $p_{ij} = \epsilon$ if $v_i \neq v_j$ and $p_{ij} = 1 - \epsilon$ if $v_i = v_j$ (that is, according to the true agreement relationship between users i and j). We fix $\epsilon = 0.001$.

We run LBP on the resulting graph, and label each node v_i as a spammer if the probability that LBP assigns it is larger than a pre-defined threshold τ .

Figure 3 depicts the AUC of the LBP classification when varying k_1 and k_2 between 0 and 15. We see that the AUC approaches 1 when each node has “enough” neighbors from each class.

Sparsification and Edge Potentials

In practice, we have the true labels of a small set of nodes, and we use it to learn and predict the edge potentials between the remaining unlabelled nodes. We shall use this prediction for a sparsification procedure similar to what was just described. Our sparsification proceeds as follows:

- (1) The first step is to train a machine learning algorithm on the set of users whose labels we know with the objective of predicting p_{ij} , the probability that a pair of reviewers i, j belongs to the same class (either both spammers or both

benign). The exact choice of ML algorithm, along with the parameters is explained in experimental setting section. (2) Compute p_{ij} using the ML model for all the remaining edges of the graph (edges that do not connect two users from the training set). (3) Choose all edges for which $p_{ij} \in [0.95, 1]$ or $p_{ij} \in [0, 0.05]$ and set the potential in Eq. (3) accordingly. We call these edges the “trusted” edges. LBP will be run on the graph containing the trusted edges.

The input to the machine learning algorithm in steps (1) and (2) is a set of features that is extracted from the metadata of both the users and the reviews. Typical metadata includes the text of the review, the rating that the reviewer gave the product, the total number of reviews that the reviewer wrote, etc. The Yelp set of features is described in Tables 2 and 3.

Additional sparsification of the graph can be obtained by removing edges that connect users whose reviews of the same product were written in faraway times. Namely, two users i and j share an edge only if they wrote a review for the same product, and these reviews were written within a period of T days (following previous work, we fixed $T = 7$). Such a graph with time-dependent edges is called a *bursty* graph and was introduced in [Fei et al. 2013]. We tested our pipeline with and without the bursty variant. The bursty sparsification, if applied, is done before the trusted edges are selected.

Node Potential

In the experiment just described, all node potentials were set to $\{0.5, 0.5\}$, and only edge potentials played a role. However, there may be a gain in setting the node potentials according to the metadata features rather than ignoring it.

Similar to the way we set the edge potential, we use machine learning to predict a_i in Eq. (2). The machine learning algorithm is trained on the set of users chosen for labeling (active learning setting) with the objective of predicting spam or benign. The value of a_i is predicted for all the remaining users using the trained model (which gives the “probability” of being a spammer or benign, alongside the discrete label).

Active learning: Sampling Users

The final component in our methodology is the way we choose the set of users for training. In this work we explore two sampling rules with the same budget of k users:

1. *Random Sampling*: Pick k reviewers uniformly at random from V .
2. *Sampling from largest clique*: In this strategy, we sample k users that belong to the largest clique. The largest clique corresponds to the product on which the largest number of reviews were written. If the budget is not consumed, we sample the remainder from the second-largest clique, and so on.

Data Description

To evaluate our methodology we use three datasets that contain reviews from Yelp.com, summary statistics of which are presented in Table 8. The datasets contain reviews of

Dataset	#Reviews (fake %)	#Users (spammer %)	#Products
Y’Chi	67,395 (13.23 %)	38,063 (20.33%)	201
Y’NYC	359,052 (10.27%)	160,225 (17.79%)	923
Y’Zip	608,598 (13.22%)	260,277 (23.91%)	5,044

Table 1: Summary statistics of the three Yelp datasets [Mukherjee et al. 2013b; Rayana and Akoglu 2015].

restaurants and hotels and were collected by [Mukherjee et al. 2013b; Rayana and Akoglu 2015]. YelpChi covers the Chicago area, YelpNYC covers NYC and YelpZip is the largest, and it includes ratings and reviews for restaurants in a continuous region including NJ, VT, CT, PA, and NY. They differ in size (YelpChi is the smallest and YelpZip is the largest), as well as in the percentage of spammers out of the total number of users. Yelp has a filtering algorithm that identifies fake/suspicious reviews. The three datasets contain these labels. We partition the users into spammers: authors of at least one filtered review and benign: authors with no filtered reviews.

Alongside the text of the reviews, the dataset contains additional metadata such as ratings, timestamps. From the text and the additional data, various features are extracted, which were used in previous work that studied these datasets [Rayana and Akoglu 2015; Mukherjee et al. 2013b; Lim et al. 2010]. Tables 2 and 3 include brief descriptions of these features. Most of them are self-explanatory, and hence we omit detailed explanations for brevity. Note that we used exactly the same set of features as [Rayana and Akoglu 2015] to allow a fair comparison.

The features are used to compute both the node and the edge potentials as explained in the Methodology section.

Evaluation

In this section, we describe the results of the experiments we ran on the three Yelp datasets. We report our results and results obtained by previous work on the same datasets.

Evaluation Metrics

We evaluated the performance of *CRSDnet* using four popular metrics, which were used by previous work as well. The Average Precision (AP), which is the area under the precision-recall curve, the ROC AUC, and the precision@k. To compute precision@k we rank the reviewers according to the probability that LBP assigned each one to be a spammer. We compute the fraction of real spammers among the top k places. We compute precision@k for $k = 100, 200, \dots, 1000$.

Finally, we use the Discounted Cumulative Gain (DCG@k) which provides a weighted score that favors correct spammer predictions at the top indices. Formally, $DCG@k = \sum_{i=1}^k \frac{2^{l_i} - 1}{\log_2(i+1)}$ where $l_i = 1$ if the user

User Features	
MNR	Max. number of reviews written in a day [Mukherjee et al. 2013a,b]
PR	Ratio of positive reviews (4-5 star) [Mukherjee et al. 2013b]
NR	Ratio of negative reviews (1-2 star) [Mukherjee et al. 2013b]
avgRD	Avg. rating deviation of user’s reviews [Mukherjee et al. 2013b; Lim et al. 2010; Fei et al. 2013]
WRD	Weighted rating deviation [Lim et al. 2010]
BST	Burstiness [Mukherjee et al. 2013b; Fei et al. 2013] (spammers are often short-term members of the site).
RL	Avg. review length in number of words [Mukherjee et al. 2013b]
ACS	Avg. content similarity—pairwise cosine similarity among user’s (product’s) reviews, where a review is represented as a bag-of-bigrams [Lim et al. 2010; Fei et al. 2013]
MCS	Max. content similarity—maximum cosine similarity among all review pairs [Mukherjee et al. 2013a]

Table 2: User Features

at the i^{th} place is a correctly identified spammer, and 0 otherwise. For compatibility with other works, we actually report the normalized DCG , which is obtained by dividing $DCG@k$ by the ideal DCG which is the $DCG@k$ where all $l_i = 1$ (all top k are indeed spammers in the ideal ranking, for the k that we choose).

The Experimental Setting

There are four choices that effect the performance of *CRSDnet*: (a) the way the node potentials are computed, using ML or using a threshold function as in [Rayana and Akoglu 2015]; (b) the way the edge potentials are computed, again using ML or using a threshold function [Rayana and Akoglu 2015]; (c) the active learning sampling rule which specifies how to choose the users for which the label is revealed; (d) with the time-dependent bursty sparsification or without.

In Table 4 we summarize the seven different configurations with which we tested *CRSDnet*. Each configuration was tested with a labeled set of users that is of size 0.25%, 0.5%, 1%, 2.5% of the entire set of users. In total we have $7 \times 4 = 28$ experiments, each was ran 10 times with fresh randomness.

The machine learning algorithm that we used to compute the edge and node potentials was random forest, written in the Wolfram Language. The code is available on Github ¹. We chose this implementation as Wolfram has a good support for graph structures, on which LBP can be easily run. All the parameters of the random forest are the default ones

¹<https://github.com/users/KirilDan/projects/1>

Review Features	
Rank	Rank order among all the reviews of product [Jindal and Liu 2008]
RD	Absolute rating deviation from product’s average rating [Li et al. 2011]
EXT	Extremity of rating [Mukherjee et al. 2013a]
DEV	Thresholded rating deviation of review [Mukherjee et al. 2013a]
ETF	Early time frame [Mukherjee et al. 2013a] (spammers often review early to increase impact)
ISR	If review is user’s sole review, then $x_{ISR} = 1$, otherwise 0 [Rayana and Akoglu 2015]
PCW	Percentage of ALL-capitals words [Li et al. 2011; Jindal and Liu 2008]
PC	Percentage of capital letters [Li et al. 2011]
L	Review length in words [Li et al. 2011]
PPI	Ratio of 1st person pronouns (‘I’, ‘my’, etc.) [Li et al. 2011]
RES	Ratio of exclamation sentences containing ‘!’ [Li et al. 2011]

Table 3: Review Features

Setting	Nodes	Edges	Sampling	Bursty
1	ML	None	Random	No
2	ML	Threshold	Random	No
3	Threshold	ML	Random	No
4	ML	ML	Random	No
5	ML	ML	Clique	No
6	ML	ML	Random	Yes
7	ML	ML	Clique	Yes

Table 4: Various configurations of running *CRSDnet*. The first two columns say how potentials were computed: using ML or the threshold method of [Rayana and Akoglu 2015]. The sampling rule corresponds to the two options mentioned at the end of the Methodology section. Bursty refers to time-dependent sparsification. Setting #1 consists of only using ML to predict nodes class, without LBP.

besides the following: 950 trees, a maximum tree depth of 16, and a maximum of 0.65-fraction of the features are considered per split. The features that we used are the ones in Tables 2 and 3.

To measure the extent to which each new component in our pipeline is responsible for the improvement over previous results, we ran our pipeline also with the way that edge and node potentials were computed in [Rayana and Akoglu 2015]. For completeness, we describe this method briefly. This method is completely unsupervised. A set of features F_1, \dots, F_r is computed for every user and review. Let $f_{u,i}$ be the value of feature i for user u . For every feature F_i , the probability $p_{u,i} = Pr[F_i < f_{u,i}]$ is estimated from the data.

Method	Y*Chi				Y*NYC				Y*Zip			
	0.25%	0.5%	1%	2.5%	0.25%	0.5%	1%	2.5%	0.25%	0.5%	1%	2.5%
SpEagle			0.691				0.657				0.671	
SpEagle+			0.708				0.683				0.691	
NetSPAM							0.650	0.650				
Set. 1	0.519	0.602	0.620	0.632	0.561	0.578	0.558	0.585	0.566	0.593	0.672	0.693
Set. 2	0.669	0.701	0.711	0.729	0.664	0.663	0.687	0.692	0.685	0.700	0.784	0.794
Set. 3	0.688	0.699	0.702	0.702	0.659	0.677	0.691	0.692	0.562	0.708	0.779	0.831
Set. 4	0.689	0.712	0.723	0.731	0.673	0.681	0.685	0.665	0.684	0.707	0.706	0.828
Set. 5	0.718	0.724	0.735	0.754	0.669	0.688	0.720	0.766	0.703	0.729	0.790	0.848
Set. 6	0.673	0.730	0.719	0.741	0.668	0.671	0.682	0.696	0.628	0.707	0.783	0.835
Set. 7	0.672	0.694	0.662	0.668	0.666	0.666	0.645	0.645	0.631	0.726	0.763	0.824

Table 5: AUC performance of compared methods on all three Yelp datasets. Best results are in bold. Empty cells stand for results that were not reported in the paper or were not computed by us.

Method	Y*Chi				Y*NYC				Y*Zip			
	0.25%	0.5%	1%	2.5%	0.25%	0.5%	1%	2.5%	0.25%	0.5%	1%	2.5%
SpEagle+			0.396				0.348				0.424	
NetSPAM							0.300	0.28				
Set. 1	0.874	0.802	0.896	0.852	0.902	0.913	0.917	0.916	0.885	0.896	0.901	0.886
Set. 2	0.901	0.879	0.882	0.906	0.915	0.912	0.921	0.912	0.782	0.805	0.859	0.883
Set. 3	0.901	0.793	0.897	0.883	0.912	0.919	0.924	0.924	0.794	0.845	0.923	0.941
Set. 4	0.890	0.825	0.896	0.901	0.917	0.922	0.968	0.927	0.859	0.870	0.869	0.875
Set. 5	0.909	0.906	0.900	0.914	0.875	0.727	0.925	0.926	0.903	0.906	0.935	0.942
Set. 6	0.907	0.707	0.886	0.913	0.914	0.920	0.926	0.921	0.864	0.891	0.927	0.946
Set. 7	0.885	0.873	0.886	0.883	0.914	0.920	0.948	0.914	0.899	0.833	0.847	0.870

Table 6: AP performance of compared methods on all three datasets. The best results are in bold. Empty cells stand for results that were not reported in the paper or were not computed by us.

Finally, a “spam score” S_u is computed for every user u via

$$S_u = 1 - \sqrt{\frac{\sum_{i=1}^r p_{u,i}}{r}} \quad (5)$$

The potential of reviewer v_i is set to $\phi(v_i) \leftarrow \{1 - S_u, S_u\}$. A similar procedure is carried out to determine edge potentials.

Results

Tables 5,6 and 7 present the results of running *CRSDnet* according to the aforementioned experimental setting, reporting the different evaluation metrics.

We compared the performance of *CRSDnet* to other algorithms that were evaluated on the same dataset: SpEagle+ [Rayana and Akoglu 2015], FraudEagle [Akoglu, Chandy, and Faloutsos 2013], NetSpam [Shehnepoor et al. 2017], Wang et al. [Wang et al. 2011] and ColluEgale [Wang et al. 2020]. The results that we report are taken from the relevant papers and are not reproductions that we carried out.

Different papers report different metrics and for different budgets. Hence some cells in the tables are left empty. Some algorithms are completely missing from a table/plot if the paper did not report that metric at all.

Table 5 provides a comparison using the AUC measure. As evident from the table, our method is superior to all pre-

vious work. The most interesting comparison is with SpEagle+ [Rayana and Akoglu 2015]. That work reports only results for a budget of 1% of the users. For all three datasets, we already obtain a better result than SpEagle+ when using only 0.25% of the users (6% better for the Chicago dataset, 10% better for the NYC dataset and 17% for the ZIP).

Table 6 reports the AP measure. Here the difference is even more dramatic. Our results are between 2 to 3 times better than SpEagle+ on all three datasets.

Table 7 reports the precision@k measure when using 1% of the users and our best configuration, #5. For the Chicago and ZIP dataset, our algorithm has the upper hand; for NYC, SpEagle+ outperforms *CRSDnet* for $k = 200$ up to $k = 700$, but the difference is very small. When our algorithm outperforms the other competitors, it is by a very large margin in most cases.

Tables 5 and 6 suggest that the best way to run *CRSDnet* is according to configuration #5 in Table 4. Namely, both edge and node potentials are set using ML, and the budget is spent on users from the largest clique. Comparing settings 2,3 vs 4 we see that using ML for both nodes and edges (setting 4) is preferable to using ML only on one of them (settings 2,3). Settings 6 and 7 show that adding the time-dependent aspect, bursty edges, only harms the performance. Configuration #1, using only ML applied to the user’s features, gave the worse performance in the AUC measure, with a big gap.

k	Y'Zip				Y'NYC				Y'Chi			
	FraudEagle	Wang	SpEagle+	CRSDnet	FraudEagle	Wang	SpEagle+	CRSDnet	FraudEagle	Wang	SpEagle+	CRSDnet
100	0.30	0.21	0.93	1	0.21	0.15	0.96	0.98	0.55	0.18	0.90	0.99
200	0.30	0.19	0.81	1	0.19	0.19	0.96	0.91	0.52	0.18	0.91	0.99
300	0.38	0.21	0.69	0.93	0.17	0.18	0.95	0.86	0.48	0.20	0.91	0.99
400	0.33	0.26	0.61	0.80	0.21	0.17	0.95	0.86	0.49	0.20	0.92	0.99
500	0.29	0.27	0.57	0.75	0.22	0.17	0.95	0.88	0.48	0.20	0.92	0.93
600	0.28	0.27	0.56	0.74	0.27	0.17	0.96	0.89	0.47	0.21	0.92	0.90
700	0.27	0.29	0.54	0.76	0.37	0.16	0.95	0.90	0.47	0.21	0.92	0.91
800	0.26	0.30	0.51	0.73	0.45	0.16	0.90	0.91	0.49	0.22	0.91	0.92
900	0.26	0.30	0.50	0.69	0.5	0.15	0.85	0.92	0.48	0.22	0.91	0.92
1000	0.28	0.32	0.49	0.67	0.45	0.16	0.82	0.92	0.47	0.22	0.90	0.93

Table 7: Precision@k of compared methods when using 1% of the users for training. The best results are in bold. *CRSDnet* runs in setting 5.

Dataset	Method	AUC
Y'Zip	DFraud (80%) [Shehnepoor et al. 2021]	0.733
Y'Zip	RF (30%)	0.740
Y'Zip	<i>CRSDnet</i> (2.5%)	0.847
Y'Chi	GraphConsis [Liu et al. 2020] (80%)	0.742
Y'Chi	RF (30%)	0.735
Y'Chi	<i>CRSDnet</i> (2.5%)	0.754

Table 8: comparison with NN-based algorithms and a RF baseline. The percentage of data used for training appears in parenthesis.

Figures 4,5 and 6 plot the NDCG@k measure for k between 0 and 1000. Here, all five competing algorithms are represented. Again, *CRSDnet* outperforms all algorithms for most values of k , in all three datasets.

Table 8 shows a comparison with two GNN-based approaches, [Shehnepoor et al. 2021] and [Liu et al. 2020]. These approaches use much more data for training (80%). As an additional baseline, we trained our Random Forest classifier this time on 30% of the data, and predicted the reviewers' class. As evident from the table, more data is not a guarantee for better performance. One possible reason for the relative poor performance of NN-based methods is that much more data (in absolute value) is needed for successfully training the NN.

Conclusion

In this work, we proposed a new holistic framework called *CRSDnet* for detecting review spammers. Our method combines both machine learning and more classical algorithmic approaches (Belief Propagation) to better exploit the relational data (user-review-product) and metadata (behavioral

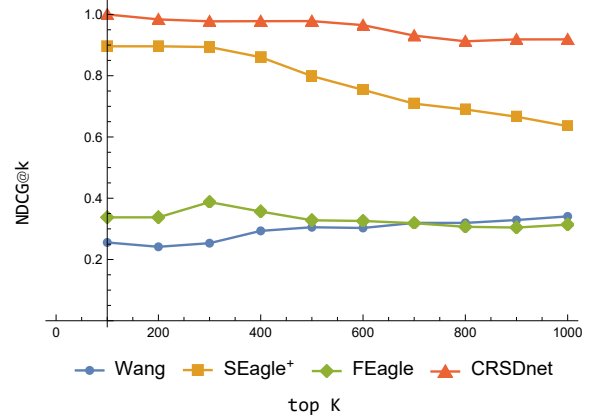


Figure 4: NDCG@k of compared methods on YelpChi

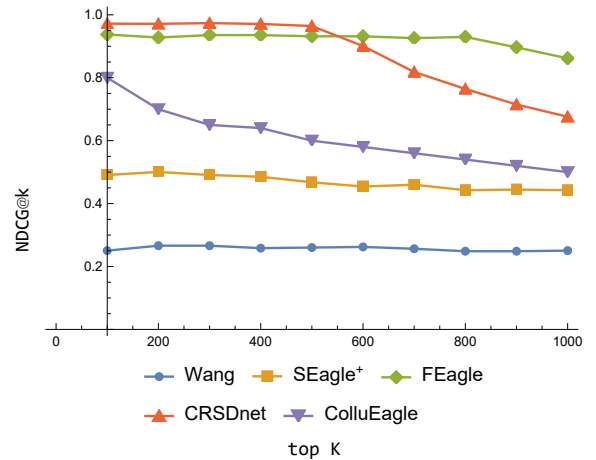


Figure 5: NDCG@k of compared methods on YelpNYC

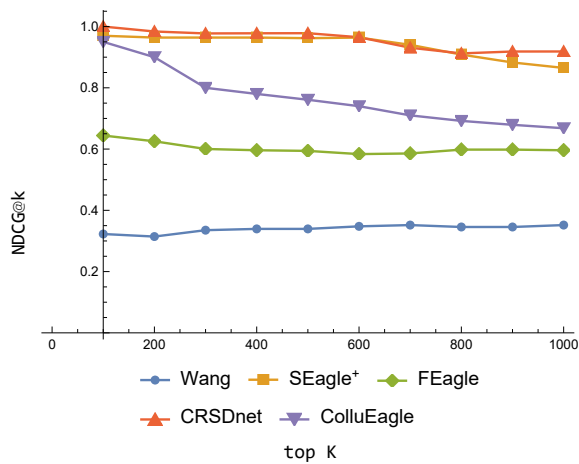


Figure 6: NDCG@k of compared methods on YelpZip

and text data) to detect such users. Adding to previous work in this line of research, we come up with two new components: using machine learning to predict the edge and node potentials, and a new sampling rule in the active learning setting – sample users from the largest clique.

Our results suggest that the two components improve performance one on top of the other, and when combined, give the best result obtained so far for the Yelp dataset.

Another point that our work highlights is that while in many settings, NN-based methods give the best results, this is highly contingent upon having sufficient data for training. The spammer detection problem is exactly one of those problems where obtaining a lot of labeled data is expensive and non-trivial. Fake reviews are many times written by professionals, and it takes an experienced person to identify them. Hence platforms like Amazon Turk may not provide an easy solution to the shortage of data problem. In such cases, old-school algorithmic ideas become relevant again (Belief Propagation), and as we demonstrate in this paper, their performance may be boosted by incorporating ML in a suitable manner (computing potentials in our case) alongside domain expertise (sampling from the largest clique, following the insight about collusive spamming [Wang et al. 2020]).

One limitation of our work is the fact that we only tested on one platform, Yelp. Future work should run our pipeline on other datasets, once they become publicly available (as far as we know there is only one more dataset in English from Amazon which is publicly available). Also, we only considered the problem of user classification. It will be interesting to extend our method for the task of review classification.

References

Akoglu, L.; Chandy, R.; and Faloutsos, C. 2013. Opinion fraud detection in online reviews by network effects. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 7.

Crawford, M.; Khoshgoftaar, T.; Prusa, J. D.; Richter, A. N.;

and Najada, H. A. 2015. Survey of review spam detection using machine learning techniques. *Journal of Big Data*, 2: 1–24.

Culotta, A.; and McCallum, A. 2005. Reducing labeling effort for structured prediction tasks. In *AAAI*, volume 5, 746–751.

Fei, G.; Mukherjee, A.; Liu, B.; Hsu, M.; Castellanos, M.; and Ghosh, R. 2013. Exploiting burstiness in reviews for review spammer detection. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 7.

Flaherty, P.; Arkin, A.; and Jordan, M. I. 2006. Robust design of biological experiments. In *Advances in neural information processing systems*, 363–370.

Jindal, N.; and Liu, B. 2008. Opinion spam and analysis. In *Proceedings of the 2008 international conference on web search and data mining*, 219–230.

Kipf, T. N.; and Welling, M. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.

Lewis, D. D.; and Catlett, J. 1994. Heterogeneous uncertainty sampling for supervised learning. In *Machine Learning Proceedings 1994*, 148–156. Elsevier.

Li, F. H.; Huang, M.; Yang, Y.; and Zhu, X. 2011. Learning to identify review spam. In *Twenty-second international joint conference on artificial intelligence*.

Li, J.; Cardie, C.; and Li, S. 2013. TopicSpam: a Topic-Model based approach for spam detection. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 217–221. Sofia, Bulgaria: Association for Computational Linguistics.

Lim, E.-P.; Nguyen, V.-A.; Jindal, N.; Liu, B.; and Lauw, H. W. 2010. Detecting product review spammers using rating behaviors. In *Proceedings of the 19th ACM international conference on Information and knowledge management*, 939–948.

Liu, Z.; Dou, Y.; Yu, P. S.; Deng, Y.; and Peng, H. 2020. Alleviating the inconsistency problem of applying graph neural network to fraud detection. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1569–1572.

Mukherjee, A.; Kumar, A.; Liu, B.; Wang, J.; Hsu, M.; Castellanos, M.; and Ghosh, R. 2013a. Spotting opinion spammers using behavioral footprints. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 632–640.

Mukherjee, A.; Liu, B.; and Gance, N. 2012. Spotting Fake Reviewer Groups in Consumer Reviews. In *Proceedings of the 21st International Conference on World Wide Web, WWW '12*, 191–200. New York, NY, USA: Association for Computing Machinery. ISBN 9781450312295.

Mukherjee, A.; Venkataraman, V.; Liu, B.; and Gance, N. 2013b. What yelp fake review filter might be doing? In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 7.

- Ott, M.; Choi, Y.; Cardie, C.; and Hancock, J. T. 2011. Finding Deceptive Opinion Spam by Any Stretch of the Imagination. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, 309–319. Portland, Oregon, USA: Association for Computational Linguistics.
- Pearl, J. 2014. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann.
- Perozzi, B.; Al-Rfou, R.; and Skiena, S. 2014. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 701–710.
- Rayana, S.; and Akoglu, L. 2015. Collective opinion spam detection: Bridging review networks and metadata. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 985–994. ACM.
- Rayana, S.; and Akoglu, L. 2016. Collective opinion spam detection using active inference. In *Proceedings of the 2016 SIAM International Conference on Data Mining*, 630–638. SIAM.
- Schein, A. I.; and Ungar, L. H. 2007. Active learning for logistic regression: an evaluation. *Machine Learning*, 68(3): 235–265.
- Settles, B.; and Craven, M. 2008. An analysis of active learning strategies for sequence labeling tasks. In *Proceedings of the conference on empirical methods in natural language processing*, 1070–1079. Association for Computational Linguistics.
- Shehnpoor, S.; Salehi, M.; Farahbakhsh, R.; and Crespi, N. 2017. NetSpam: A Network-Based Spam Detection Framework for Reviews in Online Social Media. *IEEE Transactions on Information Forensics and Security*, 12(7): 1585–1595.
- Shehnpoor, S.; Togneri, R.; Liu, W.; and Bennamoun, M. 2021. DFraud3: Multi-Component Fraud Detection free of Cold-start. *IEEE Transactions on Information Forensics and Security*, 1–1.
- Tang, J.; Qu, M.; Wang, M.; Zhang, M.; Yan, J.; and Mei, Q. 2015. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on world wide web*, 1067–1077.
- Viviani, M.; and Pasi, G. 2017. Credibility in social media: opinions, news, and health information - A survey. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 7: e1209.
- Wang, G.; Xie, S.; Liu, B.; and Philip, S. Y. 2011. Review graph based online store review spammer detection. In *2011 IEEE 11th international conference on data mining*, 1242–1247. IEEE.
- Wang, J.; Wen, R.; Wu, C.; Huang, Y.; and Xion, J. 2019a. Fdgars: Fraudster detection via graph convolutional networks in online app review system. In *Companion Proceedings of The 2019 World Wide Web Conference*, 310–316.
- Wang, J.; Wen, R.; Wu, C.; Huang, Y.; and Xion, J. 2019b. FdGars: Fraudster Detection via Graph Convolutional Networks in Online App Review System. In *Companion Proceedings of The 2019 World Wide Web Conference*, WWW '19, 310–316. New York, NY, USA: Association for Computing Machinery. ISBN 9781450366755.
- Wang, Z.; Gu, S.; and Xu, X. 2018. GSLDA: LDA-Based Group Spamming Detection in Product Reviews. *Applied Intelligence*, 48(9): 3094–3107.
- Wang, Z.; Gu, S.; Zhao, X.; and Xu, X. 2018. Graph-Based Review Spammer Group Detection. *Knowl. Inf. Syst.*, 55(3): 571–597.
- Wang, Z.; Hu, R.; Chen, Q.; Gao, P.; and Xu, X. 2020. ColuEagle: collusive review spammer detection using Markov random fields. *Data Mining and Knowledge Discovery*, 34: 1621–1641.
- Xie, S.; Wang, G.; Lin, S.; and Yu, P. S. 2012. Review Spam Detection via Time Series Pattern Discovery. In *Proceedings of the 21st International Conference on World Wide Web*, WWW '12 Companion, 635–636. New York, NY, USA: Association for Computing Machinery. ISBN 9781450312301.
- Xu, C.; Zhang, J.; Chang, K.; and Long, C. 2013. Uncovering Collusive Spammers in Chinese Review Websites. In *Proceedings of the 22nd ACM International Conference on Information; Knowledge Management, CIKM '13*, 979–988. New York, NY, USA: Association for Computing Machinery. ISBN 9781450322638.
- Xu, X.; Yuruk, N.; Feng, Z.; and Schweiger, T. A. J. 2007. SCAN: A Structural Clustering Algorithm for Networks. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '07, 824–833. New York, NY, USA: Association for Computing Machinery. ISBN 9781595936097.
- Ye, J.; and Akoglu, L. 2015. Discovering Opinion Spammer Groups by Network Footprints. In Appice, A.; Rodrigues, P. P.; Santos Costa, V.; Soares, C.; Gama, J.; and Jorge, A., eds., *Machine Learning and Knowledge Discovery in Databases*, 267–282. Cham: Springer International Publishing.
- Yelp. 2014. Yelp Challenge Data. <https://www.yelp.com/dataset>. Accessed: 2020-01-01.