

# Net Spam: Opinion Spam Detection in Online Review Communities

Abdul Qadir Ansari<sup>1</sup> Akshit Bansal<sup>2</sup> Anuradha Sharma<sup>3</sup> Kumar Sanu<sup>4</sup> Sanjoy Das<sup>5</sup>

<sup>1,2,3,4,5</sup>Research Scholar, India

**Abstract**—These days, a majority of peoples depends on accessible substance in online networking in their choices (e.g. surveys and criticism on a theme or item). The likelihood that anyone can leave a survey give a brilliant chance to spammers to compose spam reviews about items and administrations for various interests. Distinguishing these spammers and the spam content is an intriguing issue of research and despite the fact that a significant number of studies have been done as of late toward this end, yet so far the approaches set forth still scarcely identify spam surveys, and none of them demonstrate the significance of each separated element compose. In this investigation, we propose a novel structure, named NetSpam, which uses spam highlights for demonstrating survey datasets as heterogeneous data systems to delineate identification technique into a classification issue in such systems. Utilizing the significance of spam highlights help us to get better outcomes as far as various measurements probed true review datasets from Yelp and Amazon sites. The outcomes demonstrate that NetSpam beats the current strategies and among four classes of highlights; including review behavioral, client behavioral, and survey linguistic, user-linguistic, the first type of features performs better than alternate classifications.

**Keywords**—Social Media, Social Network, Spammer, Spam Review, Fake Review, Heterogeneous Information Networks

## I. INTRODUCTION

Online Social Media entryways assume an influential part in data spread which is considered as an essential hotspot for makers in their promoting efforts and additionally for clients in choosing items and administrations. In the previous years, individuals depend on a considerable measure on the composed surveys in their basic leadership procedures, and positive/negative review empowering/disheartening them in their choice of items and administrations. Furthermore, composed review likewise help specialist organizations to upgrade the nature of their items and administrations. These surveys hence have turned into a vital factor in the achievement of a business while positive review can bring benefits for an organization, negative review can conceivably affect believability and cause monetary misfortunes. The way that anybody with any character can leave remarks as review, gives an enticing chance to spammers to compose counterfeit surveys intended to delude clients' assessment. These deceptive surveys are then increased by the sharing capacity of online networking and engendering over the web. The review are written to change clients' impression of how great an item or an administration are considered as spam [11], and are regularly composed in return for cash. As appeared in [1], 20% of the review on the Yelp site are really spam surveys.

Then again, a lot of writing has been distributed on the systems used to recognize spam and spammers and also extraordinary kind of examination on this point [30], [31].

These systems can be classified into various classifications; some utilizing etymological examples in content [2], [3], [4], which are generally in light of bigram, and unigram, others depend on behavioral examples that depend on highlights extricated from designs in clients' conduct which are for the most part metadata based [34], [6], [7], [8], [9], and even a few strategies utilizing charts and diagram based calculations and classifiers [10], [11], [12].

In spite of this awesome arrangement of endeavors, numerous viewpoints have been missed or stayed unsolved. One of them is a classifier that can compute highlight weights that demonstrate each element's level of significance in deciding spam review. The general idea of our proposed system is to show a given survey dataset as a Heterogeneous Information Network (HIN) [19] and to outline issue of spam discovery into a HIN classification issue. Specifically, we show survey dataset as a HIN in which review are associated with various hub composes, (for example, highlights, and clients). A weighting calculation is then utilized to compute each element's significance (or weight). These weights are used to ascertain the final marks for review utilizing both unsupervised and regulated methodologies.

To assess the proposed arrangement, we utilized two-example survey datasets from Yelp and Amazon sites. In light of our perceptions, defining two perspectives for highlights (survey client and behavioral-etymological), the classified includes as review behavioral have more weights and yield better execution on spotting spam surveys in both semi-regulated and unsupervised methodologies. What's more, we show that utilizing diverse supervisions, for example, 1%, 2.5% and 5% or utilizing an unsupervised approach, make no discernible variety in the execution of our approach. We watched that element weights can be included or evacuated for marking and thus time many-sided quality can be scaled to a specific level of exactness. As the consequence of this weighting step, we can utilize fewer highlights with more weights to acquire better precision with less time many-sided quality. What's more, sorting highlights in four noteworthy classifications (survey behavioral, client behavioral, review semantic, client etymological), causes us to see how much every class of highlights is added to spam identification. In outline, our principal commitments are as per the following:

- 1) We propose Net Spam system that is a novel system based approach which models survey organizes as heterogeneous data systems. The classification step employments distinctive Meta way composes which are imaginative in the spam discovery area.
- 2) Another weighting technique for spam highlights is proposed to decide the relative significance of each component and shows how powerful every one of highlights is in distinguishing spams from ordinary surveys. Past works [12], [20] additionally expected to address the significance of highlights predominantly in term of got precision, however not as an inherent

capacity in their system (i.e., their approach is needy to ground truth for deciding each component significance). As we clarify in our unsupervised approach, Net Spam can find highlights significance even without ground truth, and just by depending on Meta way definition and in light of qualities ascertained for each review.

- 3) Net Spam enhances the exactness contrasted with the best in class regarding time many-sided quality, which very depends to the quantity of highlights used to distinguish a spam survey; subsequently, utilizing highlights with more weights will bring about recognizing counterfeit review less demanding with less time many-sided quality.

## II. PRELIMINARIES

As specified before, we show the issue as a heterogeneous system where hubs are either genuine segments in a dataset, (for example, reviews, clients, and items) or spam highlights. To better comprehend the proposed structure we first exhibit an outline of a portion of the ideas and definitions in heterogeneous data systems [23], [22], [24].

### A. Definitions

- Definition 1 (Heterogeneous Information Network). Suppose we have  $r(> 1)$  types of nodes and  $s(> 1)$  types of relation links between the nodes, then a heterogeneous information network is defined as a graph  $G = (V, E)$  where each node  $v \in V$  and each link  $e \in E$  belongs to one particular node type and link type respectively. If two links belong to the same type, the types of starting node and ending node of those links are the same.
- Definition 2 (Network Schema). Given a heterogeneous information network  $G = (V, E)$ , a network schema  $T = (A, R)$  is a meta path with the object type mapping  $\tau : V \rightarrow A$  and link mapping  $\phi : E \rightarrow R$ , which is a graph defined over object type  $A$ , with links as relations from  $R$ . The schema describes the meta structure of a given network (i.e., how many node types there are and where the possible links exist).
- Definition 3 (Meta path). As mentioned above, there are no edges between two nodes of the same type, but there are paths. Given a heterogeneous information network  $G = (V, E)$ , a meta path  $P$  is defined by a sequence of relations in the network schema  $T = (A, R)$ , denoted in the form  $A_1(R_1)A_2(R_2)...(R_{l-1})A_l$ , which defines a composite relation  $P = R_1 \circ R_2 \circ ... \circ R_{l-1}$  between two nodes, where  $\circ$  is the composition operator on relations. For convenience, a meta path can be represented by a sequence of node types when there is no ambiguity, i.e.,  $P = A_1A_2...A_l$ . The meta path extends the concept of link types to path types and describes the different relations among node types through indirect links, i.e. paths, and also implies diverse semantics.
- Definition 4 (Classification problem in heterogeneous information networks). Given a heterogeneous information network  $G = (V, E)$ , suppose  $V^*$  is a subset of  $V$  that contains nodes of the target type (i.e., the type of nodes to be classified).  $k$  denotes the number of the class, and for each class, say  $C_1...C_k$ , we have some pre-labeled nodes in  $V^*$  associated with a single user.

The classification task is to predict the labels for all the unlabeled nodes in  $V^*$ .

### B. Feature Types

In this paper, we use an extended definition of the Meta path concept as follows. A Meta path is defined as a path between two nodes, which indicates the connection of two nodes through their shared features. When we talk about metadata, we refer to its general definition, which is data about data. In our case, the data is the written review, and by metadata, we mean data about the reviews, including the user who wrote the review, the business that the review is written for, the rating value of the review, date of written review and finally its label as spam or genuine review.

Specifically, in this work highlights for clients and surveys fall into the classes as takes after (appeared in Table I):

Review-Behavioral (RB) based highlights. This component write depends on metadata and not simply the review content. The RB classification contains two highlights; early time allotment (ETF) and Threshold rating deviation of review (DEV) [16].

- Review Linguistic (RL) based highlights. Highlights in this class depend on the survey itself and separated specifically from the content of the review. In this work we utilize two fundamental highlights in RL classification; the Ratio of first Personal Pronouns (PP1) and the Ratio of outcry sentences containing '!' (RES) [6].
- User Behavioral (UB) based highlights. These highlights are specific to every individual client and they are figured per client, so we can utilize these highlights, to sum up, the greater part of the surveys composed by that specific client. This classification has two principle includes; the Burstiness of surveys composed by a solitary client [7], and the normal of a clients' negative proportion given to various organizations [20].
- User Linguistic (UL) based highlights. These highlights are separated from the clients' dialect and show how clients are depicting their inclination or supposition about what they've encountered as a client of a specific business. We utilize this kind of highlights to see how a spammer imparts as far as wording. There are two highlights connected with for our structure in this class; Average Content Similarity (ACS) and Maximum Content Similarity (MCS). These two highlights demonstrate how much two reviews composed of two distinct clients are like each other, as spammers have a tendency to compose fundamentally the same as surveys by utilizing format pre-composed content [11].

## III. NET SPAM: THE PROPOSED SOLUTION

In this section, we provide details of the proposed solution which is shown in Algorithm III.1.

### A. Prior Knowledge

The first step is computing prior knowledge, i.e. the initial probability of review  $u$  being spam which denoted as  $y_u$ . The proposed framework works in two versions; semi-supervised learning and unsupervised learning. In the semi-supervised method,  $y_u = 1$  if review  $u$  is labeled as spam in



the pre-labeled reviews, otherwise  $y_u = 0$ . If the label of this review is unknown due the amount of supervision, we consider  $y_u = 0$  (i.e., we assume  $u$  as a non-spam review). In the unsupervised method, our prior knowledge is realized by using  $y_u = (1/L) \sum_{l=1}^L f(x_{lu})$  where  $f(x_{lu})$  is the probability of review  $u$  being spam according to feature  $l$  and  $L$  is the number of all the used features (for details, refer to [12]).

### B. Network Schema Definition

The accompanying stage is defining framework piece in perspective of a given once-over of spam features which chooses the features involved with spam disclosure. This Schema is general definitions of Meta way and shows when all is said in done how remarkable framework parts are related. For example, if the once-over of features joins NR, ACS, PP1, and ETF, the yield diagram is as presented in Fig. 1.

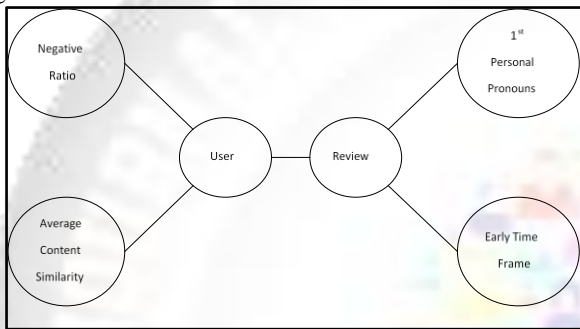


Fig. 1: An Example for a Network Schema Generated Based on a Given Spam Features list; NR, ACS, PP1 and ETF

### C. Metapath Definition & Creation

As mentioned in Section II-A, a met path is defined by a sequence of relations in the network schema. Table II shows all the met paths used in the proposed framework. As shown, the length of user-based met paths is 4 and the length of review-based met paths is 2.

For met path creation, we define an extended version of the met path concept considering different levels of spam certainty. In particular, two reviews are connected to each other if they share same value. Hassanzadeh et al. [25] propose a fuzzy-based framework and indicate for spam detection, it is better to use fuzzy logic for determining a review's label as a spam or non-spam. To be sure, there are diverse levels of spam conviction. We utilize a stage capacity to decide these levels. Specifically, given a survey  $u$ , the levels of spam conviction for met path de 1) is figured as where  $s$  means the quantity of levels. In the wake of registering  $m_u^{p1}$  for all reviews and met paths, two surveys  $u$  and  $v$  with the same met path esteems (i.e.,  $m_u^{p1} = m_v^{p1}$ ) for met path  $p_1$  are associated with each other through that met path and make one connection of review organize. The met path esteem between them indicated as  $m_u^{p1} = \frac{s \times f(x_{lu})}{s}$

Using  $s$  with a higher value will increase the number of each feature's Metapaths and hence fewer reviews would be connected to each other through these features. Conversely, using the lower value for  $s$  leads us to have bipolar values (which means reviews take value 0 or 1). Since we need enough spam and non-spam reviews for each step, with a fewer number of reviews connected to each other for every step, the spam probability of reviews take a

uniform distribution, but with a lower value of  $s$ , we have enough reviews to calculate final spamicity for each review. Therefore, accuracy for lower levels of  $s$  decreases because of the bipolar problem, and it decodes for higher values of  $s$ , because they take uniform distribution. In the proposed framework, we considered  $s = 20$ , i.e.  $m_u^{p1} \in \{0, 0.05, 0.10, \dots, 0.85, 0.90, 0.95\}$ .

### D. Classification

The classification part of NetSpam includes two steps; (i) weight calculation which determines the importance of each spam feature in spotting spam reviews, (ii) Labeling which calculates the final probability of each review being spam. Next, we describe them in detail.

#### 1) Weight Calculation

This progression registers the heaviness of each metapath. We expect that hubs' classification is done in view of their relations to different hubs in the survey organize; connected hubs may have a high likelihood of taking similar marks. The relations in a heterogeneous data arrange to incorporate the immediate connection as well as the way that can be estimated by utilizing the metapath idea. In this manner, we have to use the metapaths defined in the past advance, which speak to heterogeneous relations among hubs. In addition, this progression will have the capacity to process the heaviness of every connection way (i.e., the significance of the metapath), which will be utilized as a part of the subsequent stage

(Labeling) to gauge the mark of each unlabeled review.

Spam Feature	User Based	Review Based
Behavioral based Features	<p>Burstiness [20]: Spammers, usually write their spam reviews in short period of time for two reasons: first, because they want to impact readers and other users, and second because they are temporal users, they have to write as much as reviews they can in short time.</p> $x_{BST}(i) = \begin{cases} 0 & (L_i - F_i) \notin \tau \\ 1 - \frac{L_i - F_i}{\tau} & (L_i - F_i) \in \tau \end{cases}$ <p>where <math>L_i - F_i</math> describes days between last and first review for <math>\tau = 28</math>. Users with calculated value greater than 0.5 take value 1 and others take 0.</p> <p>Negative Ratio [20]: Spammers tend to write reviews which defame businesses which are competitor with the ones they have contract with, this can be done with</p>	<p>Early Time Frame [16]: Spammers try to write their reviews asap, in order to keep their review in the top reviews which other users visit them sooner.</p> $x_{ETF}(i) = \begin{cases} 0 & (T_i - F_i) \notin \delta \\ 1 - \frac{T_i - F_i}{\delta} & (T_i - F_i) \in \delta \end{cases}$ <p>denotes days specified written review and first written review for a specific business. We have also <math>\delta = 7</math>. Users with calculated value greater than 0.5 takes value 1 and others take 0.</p> <p>Rate Deviation using threshold [16]: Spammers, also tend to promote businesses they have contract with, so they rate these businesses with high scores. In result, there is high diversity in their given scores to different businesses</p>

	destructive reviews, or with rating those businesses with low scores. Hence, ratio of their scores tends to be low. Users with average rate equal to 2 or 1 take 1 and others take 0	which is the reason they have high variance and deviation. $x_{DEV}(i) = \begin{cases} 0 & \text{otherwise} \\ 1 - \frac{r_{ij} - \text{avg}_{i \in \mathcal{U}}}{4} & \end{cases}$ <p>where <math>\beta 1</math> is some threshold determined by recursive minimal entropy partitioning. Reviews are close to each other based on their calculated value, take same values (in [0,1)).</p>
Linguistic based Features	Average Content Similarity [7], Maximum Content Similarity [16]: Spammers, often write their reviews with same template and they prefer not to waste their time to write an original review. In result, they have similar reviews. Users have close calculated values take same values (in [0,1)).	Number of first Person Pronouns, Ratio of Exclamation Sentences containing '!' [6]: First, studies show that 'spammers use second personal pronouns much more than first personal pronouns. In addition, spammers put '!' in their sentences as much as they can to increase impression on users and highlight their reviews among other ones. Reviews are close to each other based on their calculated value, take same values (in [0,1)).

Table 1: Features for Users & Reviews in Four Defined Categories (The Calculated Values Are Based On Table 2 In [12])

The weights of the metapaths will answer a basic request; which metapath (i.e., spam feature) is better at situating spam reviews? Moreover, the weights enable us to understand the advancement to an instrument of a spam review. In like manner, since some of these spam features may get critical computational costs (for example, figuring phonetic based features through NLP methods in a gigantic review dataset), picking the more essential features in the spam distinguishing proof framework prompts better execution at whatever point the estimated cost is an issue.

To compute the weight of metapath  $p_i$ , for  $i = 1, \dots, L$  where  $L$  is the number of metapaths, we propose following equation:

$$W_{p_i} = \frac{\sum_{r=1}^n \sum_{s=1}^n mp_{r,s}^{p_i} \times y_r \times y_s}{\sum_{r=1}^n \sum_{s=1}^n mp_{r,s}^{p_i}}$$

where  $n$  denotes the number of reviews and  $mp_{r,s}^{p_i}$  is a metapath value between reviews  $r$  and  $s$  if there is a path between them through metapath  $p_i$ , otherwise  $mp_{r,s}^{p_i} = 0$ . Moreover,  $y_r(y_s)$  is 1 if review  $r(s)$  is labeled as spam in the pre-labeled reviews, otherwise 0.

## 2) Labeling

Let  $Pr_{u,v}$  be the likelihood of unlabeled review  $u$  being spam by considering its association with spam survey  $v$ . To assess  $Pr_u$ , the likelihood of unlabeled review  $u$  being spam, we propose the accompanying conditions:

$$Pr_{u,v} = 1 - \prod_{i=1}^L 1 - mp_{u,v}^{p_i} \times W_{p_i}$$

$$Pr_u = \text{avg}(Pr_{u,1}, Pr_{u,2}, \dots, Pr_{u,n})$$

Where  $n$  denotes number of reviews connected to review  $u$ . Fig. 2 shows an example of a review network and different steps of proposed framework.

It is worth to note that in creating the HIN, as much as the number of links between a review and other reviews increases, its probability to have a label similar to them increase too because it assumes that a node relation to other nodes shows their similarity. In particular, more links between a node and other non-spam reviews, more probability for a review to be non-spam and vice versa. It means that it shares features with other reviews with low spamicity and hence its probability to be a non-spam review increases.

## IV. EXPERIMENTAL EVALUATION

Row	Notation	Type	MetaPath	Semantic
1	R-DEV-R	RB	Review-Threshold Rate Deviation-Review	Reviews with same Rate Deviation from average Item rate (based on recursive minimal entropy partitioning)
2	R-U-NR-U-R	UB	Review-User-Negative Ratio-User-Review	Reviews written by different Users with same Negative Ratio
3	R-ETF-R	RB	Review-Early Time Frame-Review	Reviews with same released date related to Item
4	R-U-BST-U-R	UB	Review-User-Burstiness-User-Review	Reviews written by different users in same Burst
5	R-RES-R	RL	Review-Ratio of Exclamation Sentences containing '!'-Review	Reviews with same number of Exclamation Sentences containing '!'
6	R-PP1-R	RL	Review-first Person Pronouns-Review	Reviews with same number of first Person Pronouns
7	R-U-ACS-U-R	UL	Review-User-Average Content Similarity-User-Review	Reviews written by different Users with same Average Content Similarity using cosine similarity score



8	R-U- MCS-U- R	UL	Review-User- Maximum Content Similarity- User-Review	Reviews written by different Users with same Maximum Content Similarity using cosine similarity score
---	---------------------	----	--	--

Table 2: Meta Paths Used in the Net Spam Framework

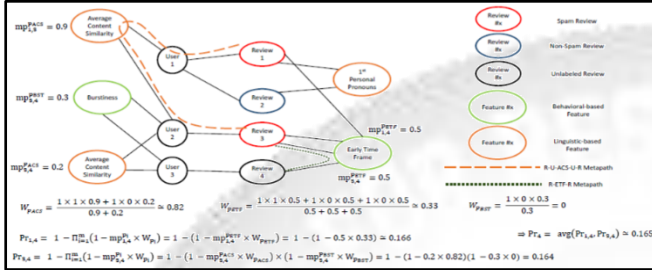
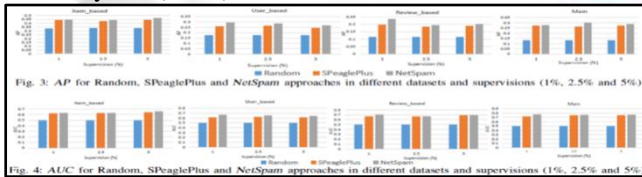


Fig. 2: An Example of Review Network & Different Steps of Proposed Framework

### 1) Datasets

Table III incorporates an outline of the datasets and their qualities. We utilized a dataset from Yelp, presented in [12], which incorporates right around 608,598 reviews composed by clients of eateries and lodgings in NYC. The dataset incorporates the commentators' impressions and remarks about the quality, and different viewpoints identified with eateries (or inns). The dataset additionally contains marked surveys as ground truth (purported close ground-truth [12]), which demonstrates whether a review is a spam or not. Howl dataset was marked utilizing filtering calculation drew in by the Yelp recommender, and albeit none of the recommenders are impeccable, as indicated by [36] it produces trustable outcomes. It discloses enlisting somebody to compose diverse phony reviews on various online networking destinations, it is the cry calculation that can spot spam surveys and rank one specific spammer at the highest point of spammers. Different properties in the dataset are the rate of analysts, the date of the composed review, and date of genuine visit, and also the client's and the eatery's id (name).



We created three other datasets from this main dataset as follow:

Dataset	Reviews (spam%)	Users	Business (Resto. & hotels)
Main	608,598 (13%)	260,277	5,044
Review-based	62,990 (13%)	48,121	3,278
Item-based	66,841 (34%)	52,453	4,588
User-based	183,963 (19%)	150,278	4,568
Amazon	8,160 (-)	7685	243

Table 3: Review Datasets used in This Work

A review-based dataset includes 10% of the reviews from the Main dataset, randomly selected using uniform distribution.

- An item-based dataset composes of 10% of the randomly selected reviews of each item, also based on uniform distribution (as with Review-based dataset)
- A user-based dataset includes randomly selected reviews using uniform distribution in which one review is selected from every 10 reviews of a single user and if the number of reviews was less than 10, uniform distribution has been changed in order to at least one review from every user get selected.
- In addition to the presented dataset, we also used another real-world set of data from Amazon [34] to evaluate our work on unsupervised mode. There is no credible label in the Amazon dataset (as mentioned in [35]), but we used this dataset to show how much our idea is viable on other datasets beyond Yelp and results for this dataset is presented in Sec. IV-C3.

### A. Evaluation Metrics

We have used Average Precision (AP) and Area Under the Curve (AUC) as two metrics in our evaluation. AUC measures accuracy of our ranking based on False Positive Ratio (FPR as y-axis) against True Positive Ratio (TPR as x-axis) and integrate values based on these two measured values. The estimation of this metric increments as the proposed strategy performs well in positioning, and tight clamp versa. Let A be the rundown of arranged spam reviews with the goal that A(i) signifies a survey arranged on the i<sup>th</sup> list in A. In the event that the quantity of spam (non-spam) surveys before review in the i<sup>th</sup> record is equivalent to n<sub>i</sub> and the aggregate number of spam (non-spam) surveys is equivalent to f, at that point TPR (FPR) for the i<sup>th</sup> is registered as n<sub>i</sub>/f. To figure the AUC, we set TPR esteems as the x-hub and FPR esteems on the y-pivot and after that incorporate the region under the bend for the bend that uses their qualities. We get an incentive for the AUC utilizing:

$$AUC = \sum_{i=2}^n (FPR(i) - FPR(i-1)) * (TPR(i))$$

Where n indicates a number of reviews. For AP we first need to figure file of best-arranged reviews with spam names. Let records of arranged spam surveys in list A with spam marks in ground truth resemble list I, at that point for AP we have:

$$AP = \sum_{i=1}^n \frac{i}{I(i)}$$

As the first step, two metrics are rank-based which means we can rank the final probabilities. Next, we calculate the AP and AUC values based on the reviews' ranking in the final list.

In the most optimum situation, all of the spam reviews are ranked on top of the sorted list; In other words, when we sort spam probabilities for reviews, all of the reviews with spam labels are located on top of the list and ranked as the first reviews. With this assumption, we can calculate the AP and AUC values. They are both highly dependent on the number of features. For the learning process, we use different supervisions and we train a set for

weight calculation. We also engage these supervisions as fundamental labels for reviews which are chosen as a training set.

## B. Main Results

In this area, we assess NetSpam from the alternate point of view and contrast it and two different methodologies, Random approach and SPeaglePlus [12]. To contrast and the first one, we have built up a system in which reviews are associated with each other haphazardly. The second approach utilizes an outstanding chart based calculation called as "LBP" to figure final names. Our perceptions indicate NetSpam, outflanks these current techniques. At that point investigation of our perception is performed and finally we will inspect our structure in unsupervised mode. Finally, we examine time the



TABLE IV: Weights of all features (with 5% data as train set); features are ranked based on their overall average weights.

Dataset - Weights	DEV	NR	ETF	BST	RES	PP1	ACS	MCS
Main	0.0029	0.0032	0.0015	0.0029	0.0010	0.0011	0.0003	0.0002
Review-based	0.0023	0.0017	0.0017	0.0015	0.0010	0.0009	0.0004	0.0003
Item-based	0.0010	0.0012	0.0009	0.0009	0.0010	0.0010	0.0004	0.0003
User-based	0.0017	0.0014	0.0014	0.0010	0.0010	0.0009	0.0005	0.0004

Many-sided quality of the proposed system and the effect of camouflage methodology on its execution.

### 1) Accuracy

Figures 3 and 4 show the execution as far as the AP and AUC. As it's appeared in the majority of the four datasets NetSpam beats SPeaglePlus extraordinarily when the quantity of highlights increments. What's more, extraordinary supervisions have no impressive impact on the metric esteems neither on NetSpam nor SPeaglePlus. Results additionally demonstrate the datasets with higher level of spam reviews have better execution since when portion of spam surveys in a specific dataset expands, likelihood for anreview to be a spam survey increments and subsequently more spam surveys will be named as spam reviews and in the aftereffect of AP measure which is exceptionally reliant on spam rate in a dataset. Then again, AUC measure does not fluctuate excessively, in light of the fact that this metric isn't subject to spam reviews rate in a dataset, however on the final arranged rundown which is ascertained in light of the final spam likelihood.

### 2) Feature Weights Analysis

Next we discuss features weights and their involvement to determine spamicity. First, we inspect how much AP and AUC are dependent on a variable number of features. Then we show these metrics are different for the four feature types explained before (RB, UB, RL, and UL). To show how much our work on weights calculation is effective, first we have simulated framework on several runs with whole features and used most weighted features to find out best combination which gives us the best results. Finally, we found which category is most effective category among those listed in Table I.

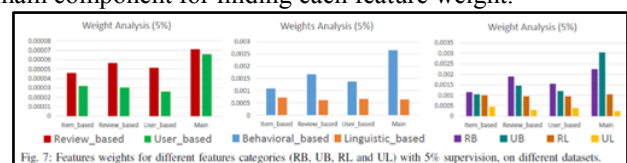
Dataset Impression on Spam Detection: As we clarified beforehand, unique datasets yield diverse outcomes in light of their substance. For all datasets and most weighted highlights, there is a sure grouping of highlights weights. As is appeared in fig.5forfourdatasets, in every one of them, highlights for the Main data set have more weights and features for Review-based dataset remain in the second

position. The third position has a place with a User-based dataset and finally Item-based dataset has the base weights (for at any rate the four highlights with general weights).

Features Weights Importance: As shown in Table IV, there are a couple of features which are more weighted than others. Combination of these features can be a good hint for obtaining better performance. The results of the Main dataset show all the four behavioral features are ranked as first features in the final overall weights. In addition, as shown in the Review based as well as other two datasets, DEV is the most weighted feature. This is also same for our second most weighted feature, NR. From the third feature to the last feature there is a different order of the mentioned features. The third feature for both datasets User-based and Review-based is same, ETF, while for the other dataset, Item-based, PP1 is at rank 3. Going further, we see in the Review-based data set all four most weighted features are behavioral-based features which shows how much this type of features are important in detecting spams as acknowledged by other works as well [12], [20].

AswecanseeinFig.6, there is a strong correlation between features weights and the accuracy. For the Main dataset, we can see this correlation is much more obvious and also applicable. Calculating weights using NetSpam help us to understand how much a feature is effective in detecting spam reviews; since as much as their weights increase two metrics including AP and AUC also increase respectively and therefore our framework can be helpful in detecting spam reviews based on feature importance.

The observations indicate larger datasets yield better correlation between features weights and also its accuracy in term of AP. Since we need to know each feature rank and importance we use Spearman's rank correlation for our work. In this experience our main dataset has correlation value equal to 0.838 (p-value=0.009), while this value for our next dataset, User-based one, is equal to 0.715 (p-value = 0.046). As much as the size of dataset gets smaller in the experiment, this value drops. This problem is more obvious in Item and Review-based datasets. For Item-based dataset, a correlation value is 0.458 which is low, because sampling Item-based datasetneeds Item-based features. The features are identical to each item and are similar to user-based features. Finally, the obtained results for our smallest dataset is satisfying, because final results considering AP show a correlation near to 0.683 between weights and accuracy (similar results for SPeaglePlus as well). Weights and accuracy (in terms of AP) are completely correlated. We observed values 0.958 (pvalue=0.0001), 0.764 (p=0.0274), 0.711 (p=0.0481) and 0.874 (p=0.0045) for the Main, User-based, Item-based and Review based datasets, respectively. This result shows using weight calculation method and considering metapath concept can be effective in determining the importance of features. The similar result for SPeaglePlus also shows our weights calculation method can be generalized to other frameworks and can be used as a main component for finding each feature weight.





Our results also indicate feature weights are completely dependent on datasets, considering this fact two most important features in all datasets are same features. This means except the first two features, other features weights are highly variable regarding dataset used for extracting weights of features.

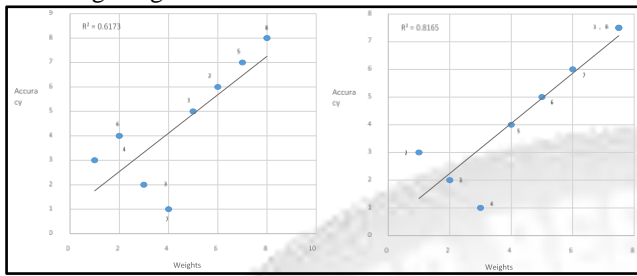


Fig. 8: Regression Graph of Features vs. Accuracy (unsupervised) for Main dataset. (see Table II for numbers)

**Features Category Analysis:** As shown in Fig. 7 there are four categories of different weights average which is very important, especially in determining which feature is more appropriate for spotting spam reviews (refer to Sec. IV-C2). Since results for different supervision are similar we have just presented the results for 5% supervision. We have analyzed features based on their categories and obtained results in all datasets show that Behavioral-based features have better weights than linguistic ones which are confirmed by [16] and [12]. Analysis of separate views shows that review based features have higher weights which leads to better performance. It is worth to mention that none of the previous works have investigated this before. Same analysis on the Main dataset shows the equal importance of both categories in finding spams. On the Other hand, in the first three datasets from Table I, RB has better weights (a bit different in comparison with RU), which means this category yields better performance than other categories for spotting spam reviews. Differently, for Main dataset UB categories has better weights and has better performance than RU category and also other categories, in all datasets behavioral-based features yield better performance with any supervision.

### 3) Unsupervised Method

One of the achievements in this study is that even without using a train set, we can still find the best set of features which yield to the best performance. As it is explained in Sec. III-A, in an unsupervised approach special formulation, is used to calculate fundamental labels and next to these labels are used to calculate the features' weight and finally review labels. As shown in Fig. 8, our observations show there is a good correlation in the Main dataset in which for NetSpam it is equal to 0.78 (p-value=0.0208) and for SPeaglePlus this value reach 0.90 (p=0.0021). As another example for the user-based dataset, there is a correlation equal to 0.93 (p=0.0006) for NetSpam, while for SPeagle this value is equal to 0.89 (p=0.0024). This perception shows NetSpam can organize highlights for the two structures. Table V shows that there is a sure grouping in include weights and it implies in spam identification issues, spammers and spam surveys have normal practices, regardless of what informal community they are composing the review for Amazon or Yelp. For every one of them, DEV is most weighted highlights, trailed by NR, ETF, and BST.

### 4) Time Complexity

If we consider the Main dataset as input to our framework, time complexity with these circumstances is equal to  $O(e2m)$  where  $e$  is the number of edges in created network or reviews number. It means we need to check if there is a metapath between a certain node (review) with other nodes which is  $O(e2)$  and this checking must be repeated for a very feature. Along these lines, our opportunity intricacy for the offline mode in which we give the Main dataset to a structure and ascertain spamicity of entire reviews, is  $O(e2m)$  where  $m$  is various highlights. In online mode, anreview is given to NetSpam to see whether it is spam or not, we have to check if there is a metapath between given survey with different surveys, which is in  $O(e)$ , and like offline mode, it must be rehased for each element and each esteem. Along these lines the many-sided quality is  $O(em)$ .

### 5) The Impact of Camouflage Strategy

One of the difficulties that spam recognition approaches confront is that spammers frequently compose non-spam surveys to conceal their actual character known as camouflage. For instance, they compose positive surveys for good eatery or negative reviews for low-quality ones; consequently, every spam indicator framework neglects to distinguish this sort of spammers or if nothing else has some inconvenience to spot them. In the past investigations, there are distinctive methodologies for dealing with this issue. For instance, in [12], the writers accept there is dependably a little likelihood that a decent review was composed by a spammer and put this supposition in its similarity framework. In this examination, we attempted to deal with this issue by utilizing weighted metapaths. Specifically, we accept that regardless of whether a survey has an almost no incentive for a specific component, it is considered in highlight weights figuring. In this way, rather than thinking about metapaths as double ideas, we take 20 esteems which signified as  $s$ . Without a doubt, if there is a camouflage its love will be decreased. As we clarified in Section III-C in such issues it is smarter to propose a fluffy system, instead of utilizing a bipolar esteem (0,1).

## V. RELATED WORKS

In the most recent decade, an extraordinary number of research thinks about the spotlight on the issue of spotting spammers and spam surveys. Notwithstanding, since the issue is non-insignificant and testing, it stays a long way from completely fathomed. We can condense our dialog about past investigations in three after classes.

### A. Semantic-Based Methods

This approach removes phonetic based highlights to find spam surveys. Feng et al. [13] utilize unigram, bigram, and their arrangement. Different examinations [4], [6], [15] utilize different highlights like pairwise highlights (includes between two reviews; e.g. content comparability), the level of CAPITAL words in anreview for finding spam surveys. Lai et al. in [33] utilize a probabilistic dialect displaying to spot spam. This examination exhibits that 2% of reviews composed of business sites are really spam.

### B. Conduct based Methods

Approaches in this group almost use reviews metadata to extract highlights; those which are the ordinary example of

commentator conduct. Feng et al. in [21] center on an appropriation of spammers rating on various items and follows them. In [34], Jindal et. al separate 36 behavioral highlights and utilize a directed technique to find spammers on Amazon and [14] demonstrates behavioral highlights demonstrate spammers' personality superior to phonetic ones. Xue et al. in [32] utilize rate deviation of a specific client and utilize a trust-mindful model to find the connection between clients for figuring the final spamicity score. Minnich et al. in [8] utilize transient and area highlights of clients to find uncommon conduct of spammers. Li et al. in [10] utilize some essential highlights (the e.g extremity of surveys) and after that run an HNC (Heterogeneous Network Classifier) to find final names on Dianpings dataset. Mukherjee et al. in [16] nearly draw in behavioral highlights like rate deviation, furthest point and so on. Xie et al. in [17] additionally utilize a transient example (time window) to find singleton surveys (reviews are composed only once) on Amazon. Luca et al. in [26] utilize behavioral highlights to indicate expanding rivalry between organizations prompts extensive development of spam surveys on items. Crawford et al. in [28] show utilizing diverse classification approach require the distinctive number of highlights to accomplish the coveted execution and propose approaches which utilize fewer highlights to achieve that execution and henceforth prescribe to enhance their execution while they utilize fewer highlights which drives them to have better-multifaceted nature. With this viewpoint our structure is doubtful. This investigation demonstrates utilizing diverse methodologies in classification yield distinctive execution as far as various measurements.

### C. Diagram based Methods

Studies in this gathering mean to make a chart between clients, reviews, and things and utilize associations in the diagram and furthermore some system based calculations to rank or name surveys (as spam or bona fide) and clients (as the spammer or legitimate). Akoglu et al. in [11] utilize a system based calculation known as LBP (Loopy Belief Propagation) indirectly versatile cycles identified with the number of edges to find final probabilities for various segments in the system. Fei et al. in [7] additionally utilize same calculation (LBP) and use burstiness of each survey to find spammers and spam reviews on Amazon. Li et al. in [10] construct a diagram of clients, surveys, clients IP and demonstrates clients with a similar IP have same names, for instance, if a client with numerous distinctive record and same IP keeps in touch with a few reviews, they should have a similar name. Wang et al. in [18] likewise make a system of clients, surveys, and things and utilize essential suspicions (for instance a commentator is more dependable in the event that he/she composes more fair reviews) and name reviews. Wahyuni in [27] proposes a half-breed strategy for spam discovery utilizing a calculation called ICF++ which is an expansion to ICF of [18] in which simply review rating is utilized to find spam identification. This work utilizes additionally estimation investigation to accomplish better exactness specifically. A more profound investigation of writing demonstrates that behavioral highlights work superior to anything phonetic ones in term of precision they yield.

There is a decent clarification for that; as a rule, spammers tend to conceal their character for security reasons. Along these lines they are not really perceived by surveys they expound on items, however, their conduct is as yet uncommon, regardless of what dialect they are composing. In result, scientists consolidated both elements composes to expand the exactness of spam identification. The way that including each element is a tedious procedure, this is the place highlight significance is helpful. In light of our insight, there is no past technique which draws the sign of highlights (known as weights in our proposed structure; NetSpam) in the classification step. By utilizing these weights, on one hand, we include highlights significance in figuring final marks and subsequently exactness of NetSpam increment, bit by bit. Then again, we can figure out which highlight can give better execution in term of their association in interfacing spam surveys (in the proposed organize).

## VI. CONCLUSION

This examination presents a novel spam location structure, in particular, NetSpam in view of a metapath idea and in addition another diagram based strategy to name reviews depending on a rank-based naming methodology. The execution of the proposed system is assessed by utilizing two certifiable named datasets of Yelp and Amazon sites. Our perceptions demonstrate that figured weights by utilizing this metapath idea can be exceptionally viable in recognizing spam reviews and prompts a superior execution. Furthermore, we found that even without a prepare set, NetSpam can compute the significance of each component and it yields better execution in the highlights' expansion procedure, and performs superior to anything past works, with just few highlights. Also, in the wake of defining, four principle classifications for highlights our perceptions demonstrate that the reviews behavioral class performs superior to anything different classes, as far as AP, AUC and additionally in the figured weights. The outcomes additionally confirm that utilizing diverse supervisions, like the semi-regulated technique, have no perceptible impact on deciding the greater part of the weighted highlights, similarly as in various datasets.

For future work, a metapath idea can be connected to different issues in this field. For instance, a comparable system can be utilized to find spammer groups. For a finding group, reviews can be associated through gathering spammer highlights, (for example, the proposed include in [29]) and surveys with the most astounding likeness in light of metapath idea are known as groups. Furthermore, using the item includes is an intriguing future work on this investigation as we utilized highlights more identified with spotting spammers and spam surveys. Additionally, while single systems have gotten extensive consideration from different controls for over 10 years, data dispersion and substance partaking in multilayer systems is as yet a youthful research [37]. Tending to the issue of spam recognition in such systems can be considered as another exploration line in this field.

## ACKNOWLEDGMENT

This work is partially supported by Iran National Science



Foundation (INSF) (Grant No. 94017889).

## REFERENCES

- [1] J. Donfro, A whopping 20 % of yelp reviews are fake. <http://www.businessinsider.com/20-percent-of-yelp-reviews-fake-2013-9>. Accessed: 2015-07-30.
- [2] M. Ott, C. Cardie, and J. T. Hancock. Estimating the prevalence of deception in online review communities. In ACM WWW, 2012.
- [3] M. Ott, Y. Choi, C. Cardie, and J. T. Hancock. Finding deceptive opinion spam by any stretch of the imagination. In ACL, 2011.
- [4] Ch. Xu and J. Zhang. Combating product review spam campaigns via multiple heterogeneous pairwise features. In SIAM International Conference on Data Mining, 2014.
- [5] N. Jindal and B. Liu. Opinion spam and analysis. In WSDM, 2008.
- [6] F. Li, M. Huang, Y. Yang, and X. Zhu. Learning to identify review spam. Proceedings of the 22nd International Joint Conference on Artificial Intelligence; IJCAI, 2011.
- [7] G. Fei, A. Mukherjee, B. Liu, M. Hsu, M. Castellanos, and R. Ghosh. Exploiting burstiness in reviews for review spammer detection. In ICWSM, 2013.
- [8] j. Minnich, N. Chavoshi, A. Mueen, S. Luan, and M. Faloutsos. Trueview: Harnessing the power of multiple review sites. In ACM WWW, 2015.
- [9] Viswanath, M. Ahmad Bashir, M. Crovella, S. Guah, K. P. Gummadi, B. Krishnamurthy, and A. Mislove. Towards detecting anomalous user behavior in online social networks. In USENIX, 2014.
- [10] H. Li, Z. Chen, B. Liu, X. Wei, and J. Shao. Spotting fake reviews via collective PU learning. In ICDM, 2014.
- [11] L. Akoglu, R. Chandy, and C. Faloutsos. Opinion fraud detection in online reviews by network effects. In ICWSM, 2013.
- [12] R. Shebuti and L. Akoglu. Collective opinion spam detection: bridging review networks and metadata. In ACM KDD, 2015.
- [13] S. Feng, R. Banerjee, and Y. Choi. Syntactic stylometry for deception detection. Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers; ACL, 2012.
- [14] N. Jindal, B. Liu, and E.-P. Lim. Finding unusual review patterns using unexpected rules. In ACM CIKM, 2012.
- [15] E.-P. Lim, V.-A. Nguyen, N. Jindal, B. Liu, and H. W. Lauw. Detecting product review spammers using rating behaviors. In ACM CIKM, 2010.
- [16] Mukherjee, A. Kumar, B. Liu, J. Wang, M. Hsu, M. Castellanos, and R. Ghosh. Spotting opinion spammers using behavioral footprints. In ACM KDD, 2013.
- [17] S. Xie, G. Wang, S. Lin, and P. S. Yu. Review spam detection via temporal pattern discovery. In ACM KDD, 2012.
- [18] G. Wang, S. Xie, B. Liu, and P. S. Yu. Review graph based online store review spammer detection. IEEE ICDM, 2011.
- [19] Y. Sun and J. Han. Mining Heterogeneous Information Networks; Principles and Methodologies, In ICCCE, 2012.
- [20] Mukerjee, V. Venkataraman, B. Liu, and N. Glance. What Yelp Fake Review Filter Might Be Doing?, In ICWSM, 2013.
- [21] S. Feng, L. Xing, A. Gogar, and Y. Choi. Distributional footprints of deceptive product reviews. In ICWSM, 2012.
- [22] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu. Pathsim: Meta path-based top-k similarity search in heterogeneous information networks. In VLDB, 2011.
- [23] Y. Sun and J. Han. Rankclus: integrating clustering with ranking for heterogeneous information network analysis. In Proceedings of the 12<sup>th</sup> International Conference on Extending Database Technology: Advances in Database Technology, 2009.
- [24] Luo, R. Guan, Z. Wang, and C. Lin. HetPathMine: A Novel Transductive Classification Algorithm on Heterogeneous Information Networks. In ECIR, 2014.
- [25] R. Hassanzadeh. Anomaly Detection in Online Social Networks: Using Datamining Techniques and Fuzzy Logic. Queensland University of Technology, Nov. 2014.
- [26] M. Luca and G. Zervas. Fake It Till You Make It: Reputation, Competition, and Yelp Review Fraud., SSRN Electronic Journal, 2016.
- [27] D. Wahyuni and A. Djunaidy. Fake Review Detection From a Product Review Using Modified Method of Iterative Computation Framework. In Proceeding MATEC Web of Conferences. 2016.
- [28] M. Crawford, T. M. Khoshgoftaar, and J. D. Prusa. Reducing Feature set Explosion to Facilitate Real-World Review Spam Detection. In Proceeding of 29th International Florida Artificial Intelligence Research Society Conference. 2016.
- [29] Mukherjee, B. Liu, and N. Glance. Spotting Fake Reviewer Groups in Consumer Reviews. In ACM WWW, 2012.
- [30] Heydari, M. A. Tavakoli, N. Salim, and Z. Heydari. Detection of review spam: A survey. Expert Systems with Applications, Elsevier, 2014.
- [31] M. Crawford, T. D. Khoshgoftar, J. N. Prusa, A. Al. Ritcher, and H. Najada. Survey of Review Spam Detection Using Machine Learning Techniques. Journal of Big Data. 2015.
- [32] H. Xue, F. Li, H. Seo, and R. Pluretti. Trust-Aware Review Spam Detection. IEEE Trustcom/ISPA . 2015.
- [33] L. Lai, K. Q. Xu, R. Lau, Y. Li, and L. Jing. Toward a Language Modeling Approach for Consumer Review Spam Detection. In Proceedings of the 7th international conference on e-Business Engineering. 2011.
- [34] N. Jindal and B. Liu. Opinion Spam and Analysis. In WSDM, 2008.
- [35] S. Mukherjee, S. Dutta, and G. Weikum. Credible Review Detection with Limited Information using Consistency Features, In book: Machine Learning and Knowledge Discovery in Databases, 2016.
- [36] K. Weise. A Lie Detector Test for Online Reviewers. <http://bloom.bg/1KAxzhK>. Accessed: 2016-12-16.

- [37] M. Salehi, R. Sharma, M. Marzolla, M. Magnani, P. Siyari, and D. Montesi. Spreading processes in multilayer networks. In IEEE

