# ijRASET

International Journal For Research in
Applied Science and Engineering Technology

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

www.ijraset.com

Call: ⓒ08813907089  |  E-mail ID: ijraset@gmail.com

# Spammer Community Used for Spam Detection in Online Social Media Reviews

Arun K[1], Manoj M[2]

[1]PG Scholar, Dept. of Computer Science and Engineering, Jawaharlal College of Engineering and Technology, Lakkidi, Palakkad, Kerala

[2]Asst. Professor, Dept. of Computer Science and Engineering, Jawaharlal College of Engineering and Technology, Lakkidi, Palakkad, kerala

*Abstract: Nowadays, Social media and internet usage increases in day by day. All the facilities are available in a just a click on the internet. The daily usages of items are also available in the social media and E-Commerce websites. The peoples are given an importance to the E-commerce products. The user initially checks the social media reviews of that product and rating of product before buying a product. There is problem if you check a social media review of a product somebody are write a spam reviews of that product. The spam review gives a bad impression of the product to a user. Identifying the spammers and the spam content is an important research area and although a number of studies and different methodology is used to detect the spam reviews. The proposed work is spammer community, which is used to the decrease the negative reviews writing in a social media and E-commerce websites. The spammer community frameworks use the importance of spam features for review datasets as heterogeneous information networks to detect the spam in social media and E-commerce. Its also helps to getting better results in by using different parameters experimented on the review datasets taken from the amazon and Yelp websites.*
*Keywords: Social Media, E-commerce, Spam Review, Heterogeneous Information Network, Social Media Review*

## I. INTRODUCTION

In present years, customers are more important to make decisions for buying products from on E-commerce sites or retail stores. The spam review writing is an important and interesting challenge for success or failure of a product. The reviews are common factor in e-commerce system and it includes positive or negative opinions. The truthful and untruthful reviews[1] are a factor to both customers and companies in digital environment. The important and difficult task is to identify these reviews. The spammers[5] are paid for writing fake reviews for a particular product. It is a difficult task for a customer to differentiate fake reviews from the genuine reviews. It has been serious problem in multi-national companies and they are competing with other companies and also defamation in same products in the same sector. The customers write negative reviews for a product so the e-commerce profit system is decreased. Nowadays one of the most popular e-commerce and most popular site is amazon.com included the fake reviews on its website and their websites accused of provided fake reviews[1] from the datasets. Fake reviews are detected from the amazon dataset and it helped in ecommerce sites to provide the percentage of fake reviews and rate. Ratings systems and reviews are directly proportional to the customer purchased options. To write the positive reviews, it given good ratings are provide for financial performance. The negative reviews are effected the reputation of the company and its move to the economic loss. Ratings can fake reviews are down a business status. So that, its very difficult task to identify the fake reviews. By using the traditional method[12][13][14], the data analysis method is used to detect the fake reviews. Before few years, data analysis techniques are used to extract the qualitative and quantitative data. These methods are very informative to facilitate the data interpretation and which helps to get good outcomes into behind the data process. The huge amount of data can be taken and it will be considerable for the background data from the dataset. It can be able to perform the task involved by the reasoning method. This goal researcher has to be turned the fields of machine learning approach and artificial intelligence methods. By using supervised method a review can be classified as truthful review. By using the unsupervised learning method that review can be classified as fake review. By using these methods used to analyse the customer reviewer's profile, data's used for writing reviews and reviewer activity on the internet used by the cookies and generated the user profiles. The methods can be classified into two types supervised method and unsupervised method which gives spam detection probability rate. In spite of this incredible arrangement of endeavors, numerous viewpoints have been missed or stayed unsolved. One of them is a classifier that can ascertain include weights that demonstrate each element's level of significance in deciding spam surveys. The general idea of the proposed structure is to show a given survey dataset as a Heterogeneous Information Network (HIN) [11] and to outline issue of spam identification into a HIN order issue. Specifically, we

International Journal for Research in Applied Science & Engineering Technology (IJRASET)
*ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 6.887*
*Volume 6 Issue IV, April 2018- Available at www.ijraset.com*

show survey dataset as a HIN in which surveys are associated through various hub writes (for example, highlights and clients). A weighting calculation is at that point utilized to compute each element's significance (or weight). These weights are used to ascertain the last marks for audits utilizing both unsupervised and directed methodologies. In summary, the contributions of Spammer community framework a network based approach which is a review network model as heterogeneous information networks. By using different metapath types which used to detect the spam in online social medias. A review length method introduced and One Time Registration(OTP) method also introduced to write a review about a product in online social medias and E commerce websites .The accuracy also improves compare with the previous work.

## II.     RELATED WORKS

In previous works, number of research studies and experimented on the problem is to determine spammers and spam detection in online social medias and E-commerce websites. The summarize discussion about previous studies in the following section. By using the clustering[15] technique finding the  Different interaction patterns can be observed for different groups of users, characterizing and identifying user profiles in online social networks. exploit the user behavior to display more appropriate advertisements. By Clever Ant Colony Metaphor[16] find out the  to cluster social network structure through maximum clique and sub grouping criteria.

By using the Graph based and Content based and the methodology[17] used in this method is Compared Naive Bayesion , Neural Network ,SVM & Decision Tree  The other work is User Based[19] for working the experiment based on Compared Decorate, Simple Logistic, FT, Logi Boost RandomsubSpace,Bagging,j48,LibSVM

In behaviour approach used to extract the features in metadata review concept. Feng et al in [23] focus different products rating given by the spammers. In[25] Jindal et. al [27] features are extracted by using the supervised learning approach and using the amazon dataset. The rate deviation of a specific user and a trust aware model used to the spamacity calculation Li et al. in[20] use some common features and runs on heterogeneous network classifier. In[8] almost all the combination of behavioural feature are used for finding single ton reviews using a temporal pattern[24].By using different classification approaches for different number of features used to getting the high performance.  In graph based method s, aims to graphical representation of users, reviews and use some network based algorithms for ranking. The network based algorithm known as LBP(Loopy Belief Propagation)[5] used to find out the final probabilities for components are used in network. Build a graph of users, reviews, users IP are connected to each other[20]. The Proposed a hybrid method and use the ICF algorithm to find out spam detection and the latest method by using the ICF extension algorithm for find out spam detection based on review rating concept. In Linguistic based feature is used to find the spam reviews. The unigram, bigram and the composition of unigram and bigram features[22] used to find the spam reviews. Here also used to find the spam by using pairwise feature and their studies[4][6][5].In business websites written are reviews are 2%  are spam [21][18] studied in probabilistic analysis.

## III.     BASIC DEFINITIONS

The model is proposed as a Heterogeneous Information Network(HIN). The heterogeneous network nodes are in dataset as real components (such as reviews, users and products or spam features. The following concepts and basic definitions are used to understand the heterogeneous information networks[2][3][4]

*A.  Definitions*

1) *Heterogeneous Information Network:* Suppose there are $r(> 1)$ types of nodes and $s(> 1)$ types of relation links between the nodes, then it gives a heterogeneous information network.  In 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 are the same type, the starting node and ending node of links are the same.

2) *Network Schema:* The given HIN, G = (V,E), a network schema T =(A,R) is a metapath with the object type mapping $\tau$ : V -> A and link mapping $\theta$: E ->R, which is defined over object type A, with links as relations from R.

3) *Metapath:* There are no edges between two nodes of the same type, but there are paths. In HIN G =(V,E), a metapath P is defined by a sequence of relations in the network schema T = (A,R), it can be represented in the form of $A_1(R_1)A_2(R_2)\ldots$    $(R_{(l-1)})A_l$, which defines a composite relation $P = R_1 o R_2 o \ldots o R_{(l-1)}$ between the two nodes, where o is the composition operator on relations. The metapath can be represented by a sequence of node types and therefore there is no ambiguity.

4) *Classification problem in heterogeneous information networks:* Given a heterogeneous information network G = (V,E), suppose $V_0$ is a subset of V that contains the types of nodes can be classified in this network, the target type denoted as k, the number of the class, and for each class, say $C_1 \ldots C_k$, there are included some pre-labeled nodes in $V_0$ associated with a single user.

### B. Feature Types

In this paper, use extended definition of the metapath concept. A metapath is defined as there is a path between the two nodes and the representation of the connected two nodes and the shared features.

In this case, the data is the written review, and the metadata concept is data about the reviews, which includes the user who wrote the review, the review is written by business purpose for rating value of the review, date of written review and finally representation is label as spam or genuine review.

In particular, in this work features for users and reviews fall into the categories as follows (shown in Table I):

The features can be classified into four categories based on review and user in Spammer community.

1) *User Behavioural Category:* It specified to each individual user and they are calculated per user, and these features to generalize all of the reviews written by that specific user. It can be classified into two main features; the Burstiness of reviews written by a single user [7], and the average of a users' negative ratio given to different businesses [6].

2) *User-Linguistic (UL):* Its based features. It describes extracted feature from the users language and shows that users are describing their feeling or opinion about what they've experienced as a customer of a certain business. This type of feature is understood how a spammer communicates in terms of wording. There are two features engaged for this framework in this category; Average Content Similarity (ACS) and Maximum Content Similarity (MCS).

3) *Review-Behavioral (RB):* It based features. This feature is based on metadata. It contains two features; Early time frame (ETF) and Threshold rating deviation of review (DEV) [8].

4) *Review-Linguistic (RL):* In this classification depend on the audit itself. There are two primary highlights in RL classification; the Ratio of first Personal Pronouns (PP1) and the Ratio of outcry sentences containing '!' (RES) [9].

### IV. SPAMMER COMMUNITY

The detailed description of the proposed solution is explained with the basis of the algorithm and the Fig.1 shows the architecture diagram of the spammer community.

### A. Prior Knowledge

The initial step is processing earlier learning, i.e. the underlying likelihood of survey u being spam which meant as $y_u$. The proposed system works in two renditions; semi-supervisied[30] learning and unsupervised learning. In the semi-supervisied technique, $y_u = 1$ if audit u is named as spam in the pre-marked audits, generally $y_u = 0$. In the event that the name of this audit is obscure due the measure of supervision. Consider $y_u = 0$, In the unsupervised technique, our earlier learning is acknowledged by utilizing $y_u = (^1/_L) \sum_{l=1}^{L} f(xlu)$

where f(xlu) is the likelihood of survey u being spam as indicated by highlight l and L is the number of all the utilized highlights[10]

### B. Network Schema Definition

The following stage is characterizing system pattern in light of guaranteed rundown of spam highlights which decides the highlights occupied with spam identification. This Schema are general meanings of metapaths what's more, appear when all is said in done how extraordinary system parts are associated.

### C. Metapath Creation

For metapath creation, characterize an expanded form of the metapath idea thinking about various levels of spam assurance. Specifically, two surveys are associated with each other on the off chance that they share same esteem. Hassanzadeh et al. [26] propose a fluffy based structure and show for spam location, it is better to utilize fluffy rationale for deciding a survey's mark as a spam or non-spam. Without a doubt, there are diverse levels of spam sureness. It utilize a stage capacity to decide these levels. In specific, given an audit u, the levels of spam conviction for metapath pl is computed as $m^{pl}_{u,v} = \lfloor s * f(xlu) \rfloor / s$ where s indicates the quantity of levels. Subsequent to registering $m^{pl}_u$ for all surveys and metapaths, two audits u and v with the same metapath esteems (i.e., $m^{pl}_u = m^{pl}_v$ ) for metapath pl are associated with each other through that metapath and make one connection of survey

organize. The metapath esteem between them signified as $m^{pl}_{u,v} = m^{pl}_{u}$. Utilizing s with a higher esteem will build the quantity of each element's metapaths and subsequently less surveys would be associated with each other through these highlights. Then again, utilizing lower an incentive for s drives us to have bipolar esteems (which implies audits take esteem 0 or 1). Since we require enough spam and non-spam audits for each progression, with less number of audits associated with each other for each progression, the spam likelihood of surveys take uniform conveyance, however with lower estimation of s we have enough surveys to compute last spamicity for each survey. Accordingly, exactness for bring down levels of s diminishes on account of the bipolar issue, and it decades for higher estimations of s, since they take uniform dissemination.

TABLE 1: Features for users and reviews

| Spam Features | User-Based | Review Based |
|---|---|---|
| Behavioural- Based Feature | Burstiness [6]: Spammers, can be calculated the spam reviews in short period of time ,it can be first calculated by the impact readers and other users, and second one to be calculated as temporal users, they have to write as much as reviews they can in short time.<br><br>$X_{BST}(i) = \begin{cases} 0, (Li - Fi) \notin (0, \tau) \\ 1 - \frac{(Li - Fi)}{\tau} (Li - Fi) \in (0, \tau) \end{cases}$<br><br>Where Li-Fi describes the day between last and first review for $\tau = 28$<br>Users calculated value greater than 0.5 take value 1 and others take 0.<br><br>Negative Ratio [6]: Spammers are write reviews with low score is given for the products. Users with average rate equal to 2 or 1 take 1 and others take 0. | Early Time Frame[8]: when a product is launched in that time the spammers write spam reviews about the products.<br><br>$X_{ETF}(i) = \begin{cases} 0, (Ti - Fi) \notin (0.\delta) \\ 1 - \frac{Ti - Fi}{\delta} (Ti - Fi) \epsilon (0, \delta) \end{cases}$<br><br>Where Li-Fi describes days between the specified written review and first written review for a specific business. Take the value $\delta = 7$.<br>Users calculated value greater than 0.5 take value 1 and others take 0.<br><br>Rate Deviation using Threshold[8]: Spammers are given high scores with certain products<br><br>$X_{DEV}(i) = \begin{cases} 0, otherwise \\ 1 - \frac{|Rij - avg \in E*j R(e)|}{5} > \\ \qquad\qquad \beta1 \end{cases}$<br><br>Where $\beta1$ is some threshold determined by recursive minimal entropy partitioning. The values are taken in the range of (0,1) |
| Linguistic Based Feature | Average Content Similarity [7], Maximum Content Similarity [8]: Spammers are writing their reviews with same style and they do not to waste their time write an original review. Users have close calculated values take same values in [0, 1). | Number of first Person Pronouns, Ratio of Exclamation Sentences containing '!' [9]: using exclamation sentences are taken and calculate the spam reviews. Reviews are close to each other based on their calculated value, take same values in [0,1]. |

The following steps are used to implementing this architecture:

1) Stage 1: Not-suggested reviews are extricated from yelp.com utilizing crawlers. Content pre-preparing is done to evaluate all the undesirable characters and discover the surveys as it were.
2) Stage 2: Genuine/Truthful reviews are taken from Yelp scholarly test dataset. Since these reviews are cleaned, pre-processing isn't required.

3) Stage 3: The utilized unigram nearness, unigram frequency, bigram presence, bigram frequency and review length as highlights for used this model.

4) Stage 4: Training information acquired in the past advances is utilized to prepare the Naive Bayes Classifier, Support Vector Machines and Logistic Regression classifiers. This preparation information has a proportion of 50:50 i.e. it contains half of fake review and half of genuine review.

5) Stage 5: Once the Naive Bayes Classifier (NBC)[29], Support Vector Machines (SVM)[31] and Logistic Regression classifiers are prepared independently for unigram, bigram and review length, it is presently used to create the discovery exactness. This test information has 80% of prepared information and 20% of test information.

6) Stage 6: Here the prepared Naive Bayes Classifier (NBC)[29], Support Vector Machines (SVM)[31] and Logistic Regression classifiers give both test accuracy exactness and test frequency precision. n-gram: An n-gram [28] is a contiguous sequence of n items from a given sequence of text or speech. The items can be phonemes, syllables, letters, words or base pairs according to the application. These n-gram's typically are collected from a text or speech corpus.

7) Unigram Frequency: It is a component that arrangements with number of times each word unigram has happened in a specific review. Unigram Presence: It is an element that primarily sees whether a specific word unigram is available in a review.

8) Bigram Frequency: It is an element that arrangements with number of times each word bigram has happened in a specific review.

9) Bigram Presence: It is a component that predominantly sees whether a specific word bigram is available in a review.

10) Review length (RL): Review length is the normal number of words introduce in a review [6]. Typically the length of fake review will be on the lesser side due to the accompanying reasons :

a) Reviewer won't have much learning about the item/business.

b) Reviewer tries to accomplish the goal with as few words as could be expected under the circumstance.

The review length can be calculated as follows:
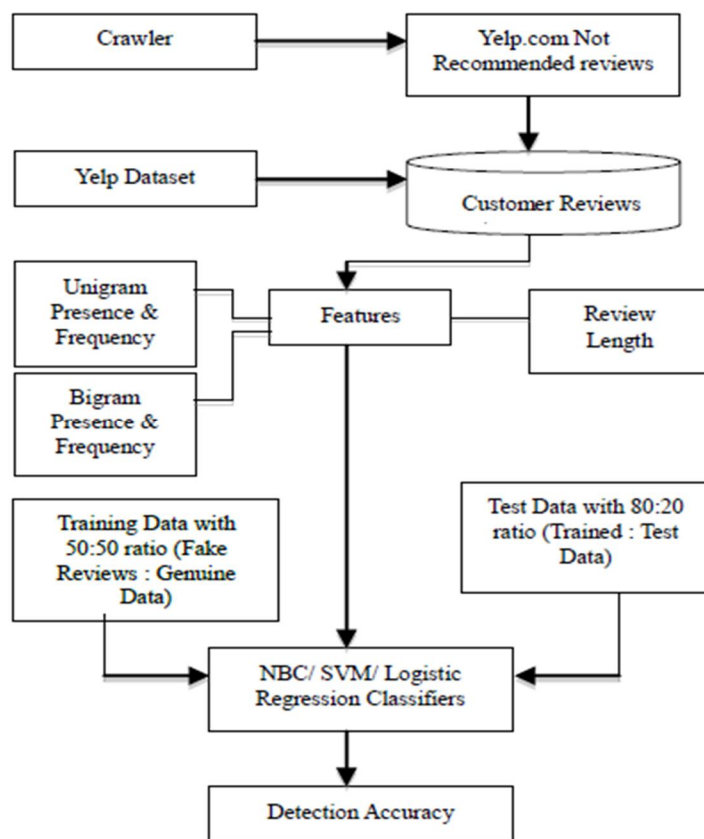
$$f_{rl} = length(rl) \qquad (1)$$



Fig 1: Architecture Diagram

The following algorithm explain the spammer community proposed work.

*Algorithm 1:* Spammer Community()

% u; v: review, $y_u$: spamicity probability of review u

% $f(x_{lu})$: initial probability of review u being spam

% $p_l$: metapath based on feature l, L: features number

% n: number of reviews connected to a review

% $m^{pl}_u$ : the level of spam certainty

% $m^{pl}_{u,v}$: the metapath value

% Earlier Knowledge

if semi-supervised mode

$$\left\{ \begin{array}{c} u \in pre\ labeled\ reviews \\ \{yu = label(u)\} \\ else \\ \{yu = 0\} \end{array} \right\}$$

Else % unsupervised mode

$$\left\{ yu = \frac{1}{L} \sum_{l=1}^{L} f(xlu) \right\}$$

% Network Schema Definition

Schema=defining schema based on spam feature list

 % Metapath Definition and Creation

For $p_l \in schema$

$$do\left\{ \begin{array}{c} for\ u,v\ \in review - dataset \\ do\left\{ \begin{array}{c} mpl,u = \frac{\lfloor s*f(xlu) \rfloor}{s} \\\\ mpl,v = \frac{\lfloor s*f(xlu) \rfloor}{s} \\ if\ mpl,u = mpl,v \\ \{mp\ pl,u,v = mpl,u\} \\ else \\ \{mp\ pl,u,v = 0\} \end{array} \right\} \end{array} \right\}$$

% Review Length Calculation

$f_{rl=}length(rl)$

% Classification- Weight Calculation

for $pl \in schemes$

do $\left\{ Wpl = \frac{\sum_{r1}^{n} \sum_{S=1}^{n} mp\ pl,r,s*yr*ys}{\sum_{r=1}^{n} \quad \sum_{S=1}^{n} mp,pl,r,s} \right\}$

% Classification –Labelling

For u,v$\in review - dataset$

$$do\left\{ \begin{array}{c} Pru,v = 1 - \prod_{pl=1}^{L} 1 - mpl\ u,v * Wpl \\\\ Pru = avg(Pru,1,Pru,2,....Pru,n) \end{array} \right\}$$

return(W, Pr)

#### D. Classification

The classification of the spammer community includes two steps: (i) Weight Calculation used to calculate the importance of spam reviews. (ii) Labelling is final probability of spam review.

*1) Weight Calculation*: Computes the weight of each metapath used in the spammer community. Its based on their relations from the review networks to other nodes. The relations are heterogeneous information networks include either direct or indirect relations by using the metapath.

To calculate the metapath weight $p_i$ for i=1,...., L where L is the number of metapathsby using this equation.

$$Wpl = \frac{\sum_{r1}^{n}\sum_{s=1}^{n} mp\, pl,r,s*yr*ys}{\sum_{r=1}^{n}\quad \sum_{s=1}^{n} mp,pl,r,s} \qquad (2)$$

*2) Labeling:* To calculate the unlabelled probability review. Consider Pr $_{u,v}$, u is represented the unlabelled review being spam and its relatios belongs to spam review v. To calculate the labelling by following the equation

$$Pru, v = 1 - \prod_{pl=1}^{L} 1 - mpl\, u, v * Wpl \qquad (3)$$

$$Pru = avg(Pru, 1, Pru, 2, ....Pru, n) \qquad (4)$$

## V. EXPERIMENTAL RESULTS

This section describes the experimental results including the datasets and the obtained results.

### A. Datasets

Table II includes dataset values and their characteristics. The dataset from amazon[10] which includes almost 60000 reviews are written by customers about the products. Three others datasets from the main dataset as follows:

Review based includes, 20% of the review from the dataset. Item based dataset, written reviews of each product items. User based dataset, one review is selected from 10 reviews of single user and the reviews less than 10.

Table II: Review Dataset

| Dataset | Spam Reviews | Users Review | Products Review |
|---|---|---|---|
| Main | 600000 | 250984 | 4000 |
| Review based | 75000 | 8000 | 2800 |
| Item based | 80000 | 9500 | 3400 |
| User based | 150000 | 19000 | 3500 |
| Amazon | 8000 | 7500 | 1000 |

### B. Evaluation Metrics

The two evaluation metrics based on the Average Precision(AP) and Area Under the Curve(AUC). It takes the value in y-axis is False Positive Ratio(FPR) against the x-axis parameter is True Positive Ratio(TPR). The equations for calculating these evaluations metrics as follows:

$$AUC = \sum_{i=1}^{n} \left|\left(FPR(i) - FPR(i-1)\right) * (TPR(i))\right| \qquad (5)$$

For calculating AP, sorted the top reviews with the spam labels.The equation as follows:

$$AP = \sum_{i=1}^{n} \frac{i}{I(i)} \qquad (6)$$

### C. Main Results

To evaluate Spammer Community from different perspective and compare it two other approaches, like NetSpam and SPeagle[10]. The accuracy and feature weight analysis graphical representation is shown in Figures 2.
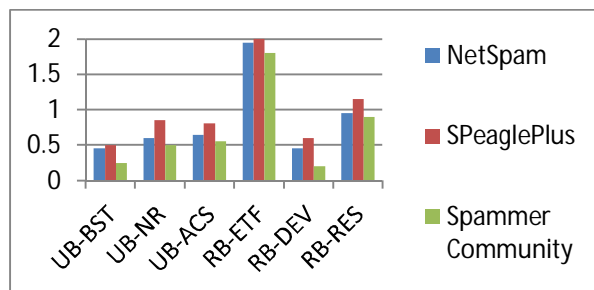
Fig. 2: Graphical representation of evaluation  metrics in different approaches in supervised.

In figure 2 , the accuracy graphical representation of each spam features used in the spammer community. Results shows that the highest spam occurs in the Early Time Framework(ETF), which compares the proposed frame work results the lowest value compare to previous works. In ETF its exceeds the value 1. In the spammer community if the value is exceeds value 1, its represents the highest spam occurs.

In Figure 3, its graphical representation of the overall performance of the spammer community in supervised mode and semi-supervised mode using the metapath concept.

In semi supervised mode, that is the dataset is taken from the real world amazon website. In figure 3 represents the Early Time Frame(ETF) and the Busrstiness(BST) shows high spammers occur in the semi-supervised mode. When the semi-supervised mode compares to the supervised mode, its spam percentage very low and the feature values are cannot exceeds 1. Other four features are cannot exceed the value 0.5. If the value above 1 means that it occurs high spam.
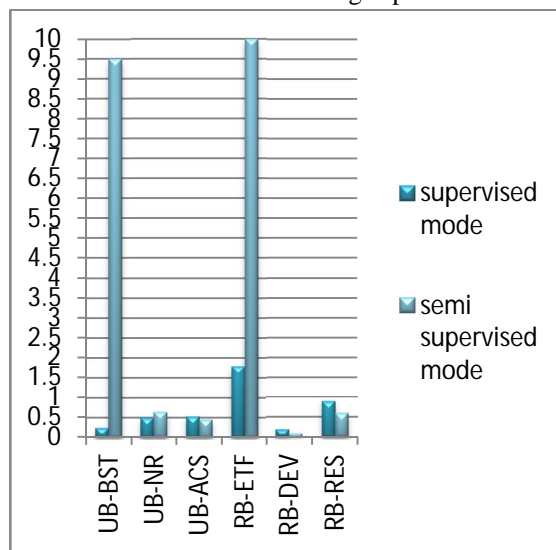


Fig 3: Spammer community in supervised and semi supervised mode

## VI.    CONCLUSION

The introduced novel framework for spam detection is Spammer Community based on the graph based as well as metapath concept. To evaluate the performance of the proposed framework is using the real world dataset. In classification approach the weights calculated by using the metapath concept can be very effective for identifying the spam detection on soaial medias. The spammer community also use a trained dataset and evaluate the each spam features which gives the better performance and its better to the previous works. To avoid the spam, OTP(One Time Password) introduced in this framework through the registered mail. The results also confirm by using the supervised and semi-supervised method by using the different datasets.

For future work, the metapath concept and graph based model can be applied to other problems. The problems can be implemented in heterogeneous information networks using the metapath concept for detecting the spammers. For detecting the spammers in multilayer networks is new research in this field.

## VII. ACKNOWLEDGMENT

## REFERENCES

[1] N. Jindal and B. Liu Opinion Spam and Analysis, In WSDM, 2008

[2] Y. Sun and J. Han, X. Yan, P. S. Yu and T. Wu Pathsim: Meta path-based top-K similarity search in heterogeneous .informayion networks. In VLDB 2011

[3] Y. Sun and J. Han, Rankclus: intergrating clustering with ranking for heterogeneous information network analysis. In Proceedings of the 12thInternational Conferenceon Extending Database Technology: Advances in Database Technology,2009

[4] C. Luo, R.Guan, Z. Wang and C. LinHet PathMine: A Novel Tranductive Classification Algorithm Heterogeneous Information networks. In ECIR 2014.

[5] L. Akoglu, R. Chandy, and C. Faloutsos, Opinion Fraud Detection in Online Reviews by Network Effects. In ICWSM, 2013.

[6] A. Mukerjee, V. Venkataraman, B. Liu, and N. Glance. What Yelp Fake Review Filter Might Be Doing?, In ICWSM, 2013.

[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] A. Mukherjee, A. Kumar, B. Liu, J. Wang, M. Hsu, M. Castellanos, and R. Ghosh. Spotting opinion spammers using behavioral footprints. In ACM KDD, 2013

[9] 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

[10] R. Shebuti and L. Akoglu. Collective opinion spam detection: bridging review networks and metadata. In ACM,KDD,2015

[11] Y. Sun and J. Han, Mining Hetergeneuous Information Networks; Principles and Methodologies, In ICCCE 2012

[12] M. Ott, C. Cardie, and J. T. Hancock. Estimating the prevalence of deception in online review communities. In ACM WWW, 2012

[13] M. Ott, Y. Choi, C. Cardie, and J. T. Hancock. Finding deceptive opinion spam by any stretch of the imagination.In ACL, 2011

[14] Ch. Xu and J. Zhang. Combating product review spam campaigns via multiple heterogeneous pairwise features. In SIAM International Conference on Data Mining,2014

[15] Marcelo Maia, Jussara Almeida, Virgílio Almeida, Identifying User Behavior in Online Social Networks, SocialNets'08, April 1, 2008 , Glasgow, Scotland, UK Copyright 2008 ACM ISBN 978-1-60558-124-8/08/04.

[16] Mohammad Al-Fayoumi, Soumya Banerjee, Jr., and P. K. Mahanti, "Analysis of Social Network Using Clever Ant Colony Metaphor", Proceedings Of World Academy Of Science, Engineering And Technology Volume 41 May 2009 Issn: 2070-3740.

[17] Alex Hai Wang, Security and Cryptography (SECRYPT), Don't Follow Me: Spam Detection in Twitter, Proceedings of the 2010 International Conference, Pages 1-10, 26-28 July 2010, IEEE.

[18] 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

[19] Kyumin Lee, James Caverlee, Steve Webb, Uncovering Social Spammers: Social Honeypots + Machine Learning, Proceeding of the 33rd International ACM SIGIR conference on Research and development in information retrieval, 2010, Pages 435–442, ACM, New York (2010).

[20] H. Li, Z. Chen, B. Liu, X. Wei, and J. Shao. Spotting fake reviews via collective PU learning. In ICDM, 2014.

[21] C. 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.

[22] 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

[23] S. Feng, L. Xing, A. Gogar, and Y. Choi. Distributional footprints of deceptive product reviews. In ICWSM, 2012.

[24] S. Xie, G. Wang, S. Lin, and P. S. Yu. Review spam detection via temporal pattern discovery. In ACM KDD,2012.

[25] 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.

[26] R. Hassanzadeh. Anomaly Detection in Online Social Networks: Using Datamining Techniques and Fuzzy Logic. Queensland University of Technology, Nov.201

[27] K. Weise. A Lie Detector Test for Online Reviewers. http://bloom.bg/1KAxzhK. Accessed: 2016-12-16.

[28] Wikipedia- n-gram http://en.wikipedia.org/wiki/N-gram.

[29] Gauri Joshi, Mr. Samadhan Sonawane, Filtering and Classification Of User Based On Social Media Data Using Memetic and Naïve Bayes Methods, (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 6 (5) , 2015, 4795-4798

[30] Mita K. Dalal and Mukesh A. Zaveri, Semisupervised Learning Based Opinion Summarization and Classification for Online Product Reviews, Hindawi Publishing Corporation Applied Computational Intelligence and Soft Computing Volume 2013.

[31] Shilpa Radhakrishnan, A Survey on Text Filtering in Online Social Networks, (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 6 (4) , 2015, 3874-3876.

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  ⓢ (24*7 Support on Whatsapp)