

## NetSpam: A Network-Based Spam Detection Framework for Reviews in Online Social Media

Mr.Mrutyunjaya Yalawar<sup>1</sup> Muthyam Kavya<sup>2</sup>, P Spoorthi Reddy<sup>2</sup>, Soma Sai Amaranath<sup>2</sup>, S Venkat Sai Kumar<sup>2</sup>  
*<sup>1</sup>Asst. Professor, Computer Science and Engineering, CMR Engineering College, medchal, T.S, India*  
*<sup>2</sup>B.Tech, Computer Science and Engineering, CMR Engineering College, medchal, T.S, India*

### Abstract—

Nowadays, a big part of people rely on available content in social media in their decisions (e.g., reviews and feedback on a topic or product). The possibility that anybody can leave a review provides a golden opportunity for spammers to write spam reviews about products and services for

different interests. Identifying these spammers and the spam content is a hot topic of research, and although a considerable number of studies have been done recently toward this end, but so far the methodologies put forth still barely detect spam reviews, and none of them show the importance of each extracted feature type. In this paper, we propose a novel framework, named NetSpam, which utilizes spam features for modeling review data sets as heterogeneous information networks to map spam detection procedure into a classification problem in such networks. Using the importance of spam features helps us to obtain better results in terms of different metrics experimented on real-world review data sets from Yelp and Amazon Web sites. The results

show that *NetSpam* outperforms the existing methods and among four categories of features, including review-behavioral, user-behavioral, review-linguistic, and user-linguistic,

the first type of features performs better than the other categories.

**Index Terms**— Social media, social network, spammer, spam review, fake review, heterogeneous information networks.

### I. INTRODUCTION

ONLINE Social Media portals play an influential role in information propagation which is considered as an important source for producers in their advertising campaigns as well as for customers in selecting products and services. In the past years, people rely a lot on the written reviews in their decision-making processes, and positive/negative

reviews encouraging/discouraging them in their selection of products and services. In addition, written reviews also help service providers to enhance the quality of their products and

services. These reviews thus have become an important factor in success of a business while positive reviews can bring benefits for a company, negative reviews can potentially impact credibility and cause economic losses. The fact that anyone with any identity can leave comments as a review provides an tempting

opportunity for spammers to write fake reviews designed to mislead users' opinion. These misleading reviews are then multiplied by the sharing function of social media and propagation over the web. The reviews written to change users' perception of how good a product or a service are considered as spam [1], and are often written in exchange for money. As shown in [1], 20% of the reviews in the Yelp website are actually spam reviews.

On the other hand, a considerable amount of literature has been published on the techniques used to identify spam and spammers as well as different types of analysis on this topic [30], [31]. These techniques can be classified into different categories; some using linguistic patterns in text [2]–[4], which are mostly based on bigram, and unigram, others are based on behavioral patterns that rely on features extracted from patterns in users' behavior which are mostly metadata-based [5–6]–

[8], [9], [34], and even some techniques using graphs and graph-based algorithms and classifiers [10]–[12].

Despite this great deal of efforts, many aspects have been missed or remained unsolved. One of them is a classifier that can calculate feature weights that show each feature's level of importance in determining spam reviews. The general concept of our proposed framework is to model a given review dataset as a Heterogeneous Information Network (HIN) [19] and to map the problem of spam detection into a HIN classification problem. In particular, we model review dataset as a HIN in which reviews are connected through different node types (such as features and users). A weighting algorithm is then employed to calculate each feature's importance (or

weight). These weights are utilized to calculate the final labels for reviews using both unsupervised and supervised approaches.

To evaluate the proposed solution, we used two sample review datasets from Yelp and Amazon websites. Based on our observations, defining two views for features (review-user and behavioral-linguistic), the classified features as review-behavioral have more weights and yield better performance on spotting spam reviews in both semi-supervised and unsupervised approaches [13–15]. In addition, we demonstrate that using different supervisions such as 1%, 2.5% and 5% or using an unsupervised approach,

make no noticeable variation on the performance of our approach. We observed that feature weights can be added or removed for labeling and hence time complexity can be scaled for a specific level of accuracy. As the result of this weighting step, we can use fewer features with more weights to obtain better accuracy with less time complexity. In addition, categorizing features in four major categories (review-behavioral, user-behavioral, review-linguistic, user-linguistic), helps us to understand how much each category of features is contributed to spam detection.

In summary, our main contributions are as follows:

- (i) We propose NetSpam framework that is a novel network-based approach which models review networks as heterogeneous information
  - (ii) networks. The classification step uses different metapath types which are innovative
  - (iii) in the spam detection domain.
  - (iv) A new weighting method for spam features is proposed to determine the relative importance
  - (v) of each feature and shows how effective each of features are in identifying spams from normal reviews.
- Previous works [16–19], [20] also aimed to address the importance of features mainly in terms of obtained accuracy, but not as a build-in function in their framework (i.e., their approach is independent to ground truth for determining each feature importance). As we explain in our unsupervised approach, NetSpam is able to find feature importance even without ground truth, and only by relying on metapath definition and based on values calculated for each review.
- (vi)
  - (vii) NetSpam improves the accuracy compared to the state-of-
  - (viii) the art in terms of time complexity, which highly depends on the number of features used to identify a spam review; hence, using features with more weights will result in
  - (ix) detecting fake reviews easier with less time complexity.

## II. RELIMINARIES

As mentioned earlier, we model the problem as a heterogeneous network [21] where nodes are either real components in a dataset (such as reviews, users and products) or

spam features. To better understand the proposed framework we first present an overview of some of the concepts and definitions in heterogeneous information networks [22]–[24].

### A. Definitions

*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 metapath 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 metastructure of a given network (i.e., how many node types there are and where the possible link exist).

*Definition 3 (Metapath):* 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 metapath  $P$  is defined by a sequence of relations in the network schema  $T = (A, R)$ , denoted in the form  $A_1(R_1)A_2(R_2)\dots(R_{l-1})A_l$ , which defines a composite relation  $P = R_1 \circ R_2 \circ \dots \circ R_{l-1}$  between two nodes, where  $\circ$  is the composition operator on relations. For convenience, a metapath can be represented by a sequence of node types when there is no ambiguity, i.e.,  $PA_1A_2\dots A_l$ . The metapath 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  $C_1 \dots C_k$ , 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 metapath concept as follows. A metapath 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 user who wrote the review, the business that the review is written for, rating value of the review, date of written review and finally its 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):

- 1) *Review-Behavioral (RB) Based Features*: This feature type is based on metadata and not the review text itself. The RB category contains two features; Early time frame (ETF) and Threshold rating deviation of review (DEV) [16].
- 3) *Review-Linguistic (RL) Based Features*: Features in this category are based on the review itself and extracted directly from text of the review. In this work we use two main features in RL category; the Ratio of 1st Personal Pronouns (PP1) and the Ratio of exclamation sentences containing '!' (RES) [6].
- 4) *User-Behavioral (UB) Based Features*: These features are specific to each individual user and they are calculated per user, so we can use these features to generalize all of the reviews written by that specific user. This category has 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 [25-30].
- 5) *User-Linguistic (UL) Based Features*: These features are extracted from the users' language and shows how users are describing their feeling or opinion about what they've experienced as a customer of a certain business. We use this type of features to understand how a spammer communicates in terms of wording. There are two features engaged for our framework in this category; Average Content Similarity

**TABLE I**

**FEATURES FOR USERS AND REVIEWS IN FOUR DEFINED CATEGORIES (THE CALCULATED VALUES ARE BASED ON [12, TABLE 2])**

Spam Feature	User-based	Review-based
Behavioral-based Features	$x_{BST}(i) = \begin{cases} 0 & (L_i - F_i) \notin (0, \tau) \\ 1 - \frac{L_i - F_i}{\tau} & (L_i - F_i) \in (0, \tau) \end{cases} \quad (1)$ <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><i>Early Time Frame</i> [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> $x_{ETF}(i) = \begin{cases} 0 & (T_i - F_i) \notin (0, \delta) \\ 1 - \frac{T_i - F_i}{\delta} & (T_i - F_i) \in (0, \delta) \end{cases} \quad (2)$ <hr/> <p><i>Rate Deviation using threshold</i> [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 which is the reason they have high variance and deviation.</p> $x_{DEV}(i) = \begin{cases} 0 & \beta_1 \leq \frac{\sigma_i}{\mu_i} \\ 1 - \frac{\sigma_i - \beta_1}{\mu_i - \beta_1} & \beta_1 < \frac{\sigma_i}{\mu_i} < 1 \end{cases} \quad (3)$ <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 <math>[0, 1]</math>).</p>
Linguistic-based Features		<p><i>Number of first Person Pronouns, Ratio of Exclamation Sentences containing '!'</i> [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 <math>[0, 1]</math>).</p>

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 to 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

$$y_u = \frac{1}{L} \sum_{l \in \text{being spam according to feature } l} f(x_{lu}) \text{ where } f(x_{lu}) \text{ is the probability of } l \text{ and } L \text{ is the}$$

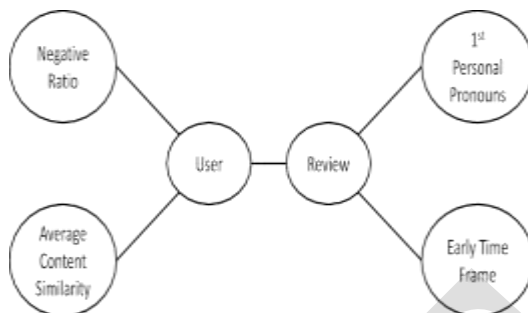


Fig. 1. An example for a network schema generated based on a given spam features list; NR, ACS, PP1 and ETF.

the metapaths used in the proposed framework. As shown, the length of user-based metapaths is 4 and the length of review-based metapaths is 2.

For metapath creation, we define an extended version of the metapath concept considering different levels of spam certainty. In particular, two reviews are connected to each other if they share same value. Hassan

### B. Network Schema Definition

The next step is defining network schema based on a given list of spam features which determines the features engaged in spam detection. This Schema are general definition of meta-  
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. Indeed, there are different levels of spam certainty. We use a step function to determine these levels. In particular, given a review  $u$ , the level of spam certainty for metapath  $p_l$  (i.e., feature  $l$ ) is calculated as  $m^{p_l}$

$$m^{p_l} = \frac{s \times f(x_{lu})}{L},$$

### C. Metapath Definition and Creation

As mentioned in Section II-

A, a metapath is defined by a sequence of relations in the network schema. Table II shows all denoted as  $m^{p_l}$   $m^{p_l}$ .

Using  $s$  with a high value will increase the number of each feature's metapaths and hence fewer reviews would be



**AlgorithmIII.1:xNetSpam()**

*Input: review-dataset,spam-feature-list,pre-labeled-reviews*

*Output:features-importance(W),spamcity-probability(Pr)*

% $u,v$ :review, $y_u$ :spamcityprobabilityofreview $u$

% $f(x_{lu})$ :initialprobabilityofreview $u$ beingspam%metapathbasedonfeature $l,L$ :featuresnumber% $n$ :numberofreviewsconnectedtoareview

% $m_u^{pl}$ :thelevelofspamcertainty

% $m_u^{pl}$ , $l$ themetapathvalue

%PriorKnowledge

**if**semi-supervisedmode

**if** $u \in \text{pre-labeled-reviews}$

$\{y_u = \text{label}(u)\}$

**else**

$\{y_u = 0\}$

**D. Classification**

Theclassificationpartof*NetSpam*includestwosteps;

(i) weight calculation which determines the importance of eachspam feature in spotting spam reviews, (ii) Labeling whichcalculates the final probability of each review being spam.NextwedescrIBethemindetail.

1) *Weight Calculation*:This step computes the weight ofeach metapath. We assume that nodes' classification is donebased on their relations to other nodes in the review network;linked nodes may have a high probability of taking the samelabels. The relations in a heterogeneous information networknot only include the direct link but also the path that can bemeasured by using the metapath concept. Therefore, we needto utilize the metapaths defined in the previous step, whichrepresent heterogeneous relations among nodes. Moreover, thisstep will be able to compute the weight of each relation path(i.e.,theimportanceofthemetapath),whichwillbeusedin thenextstep(Labeling)toestimatethelabelofeachunlabeled

: TABLEII

METAPATHSUSEDINTHE*NetSpam*FRAMEWORK

Row	Type		
	RR	Review-Threshold Rate Deviation-Review	Reviews with same Rate Deviation from average Item rate (based on recursive minimal entropy partitioning)
	UB	Review-User-Negative Ratio-User-Review	Reviews written by different Users with same Negative Ratio
	RR	Review-Early Time Frame-Review	Reviews with same released date related to Item
	UB	Review-User-Burstiness-User-Review	Reviews written by different users in same Burst
	RI.		Reviews with same number of Exclamation Sentences containing '!
	RI.	Review-first Person Pronouns-Review	Reviews with same number of first Person Pronouns
	UL.	Review-User-Average Content Similarity-User-Review	Reviews written by different Users with same Average Content Similarity using cosine similarity score
	UL.		Reviews written by different Users with same Maximum Content Similarity using cosine similarity score

TABLEIII

REVIEW DATASETSUSED IN THIS WORK

	Reviews (spam%1)	Users	Business (Resto. & hotels)
		260,277	
		48,121	

Main	608,598 (13%)		5,044
Review-based	62,990 (13%)		3,278
Item-based	66,841 (34%)		4,888
User-based	183,963 (19%)		4,568
		7685	243

lots of links with non-spam reviews, 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

This section presents the experimental evaluation part of this study including the datasets and the defined metrics as well as the obtained results. We used a dataset from Yelp, introduced in [12], which includes almost 608,598 reviews written by customers of restaurants and hotels in NYC.

##### A. Datasets

Table III includes a summary of the datasets and their characteristics. We used a dataset from Yelp, introduced in [31], which includes almost 608,598 reviews written by customers of restaurants and hotels in NYC. The dataset includes the reviewers' impressions and comments about the quality, and other aspects related to a restaurant (or hotels). The dataset also contains labeled reviews as ground truth (so-called near-ground-truth [32]), which indicates whether a review is a non-spam. The dataset was labeled using a filtering algorithm.

- Item-based dataset, composed of 10% of the randomly selected reviews of each item, also based on uniform distribution (as with Review-based dataset).

- User-based dataset, includes randomly selected reviews using uniform distribution in which one review is selected from every 10 reviews of single user and if 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 [32-34] 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 on Sec. IV-C3.

##### B. 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 value of this metric increases as the proposed method performs well in ranking, and vice-versa. Let  $A$  be the list of sorted spam reviews so that  $A(i)$  denotes a review sorted on the  $i^{th}$  index in  $A$ . If the number of spam (non-spam) reviews before review in the  $j^{th}$  index is equal to  $n$ , and the total engaged by the Yelp recommender, and although none of recommenders are perfect, but according to [36] it produces a rate of reviewers, the date of the written review, and date of actual visit, as well as the user's and the restaurant's id (name).

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

- Review-based dataset, includes 10% of the reviews integrate the area under the curve for the curve that uses their values. We obtain a value for the AUC using:

$$AUC = \sum_{i=2}^n (FPR(i) - FPR(i-1)) * (TPR(i)) \quad (7)$$

where  $n$  denotes number of reviews. For AP we first need to calculate index of top sorted reviews with spam labels. Let indexes of sorted spam reviews in list  $A$  with spam labels

$$W_{PST} = \frac{1 \times 1 \times 0.5 + 1 \times 0 \times 0.5 + 1 \times 0 \times 0.5}{0.5 + 0.5 + 0.5} = 0.33$$

$$W_{PST} = \frac{1 \times 0 \times 0.3}{0.3} = 0$$

$$Pr_{1,4} = 1 - \prod_{i=1}^m (1 - mp_{1,4}^i \times W_{PST}) = 1 - (1 - mp_{1,4}^{PAC} \times W_{PST}) = 1 - (1 - 0.5 \times 0.33) = 0.166$$

$$Pr_{2,4} = 1 - \prod_{i=1}^m (1 - mp_{2,4}^i \times W_{PST}) = 1 - (1 - mp_{2,4}^{PAC} \times W_{PST}) \times (1 - mp_{2,4}^{PST} \times W_{PST}) = 1 - (1 - 0.2 \times 0.33) \times (1 - 0.3 \times 0) = 0.164$$

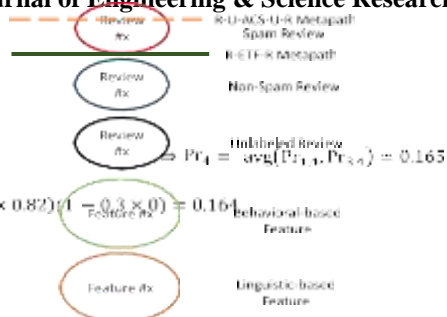


Fig.2.Anexampleofareviewnetworkanddifferentsteps ofproposedframework.

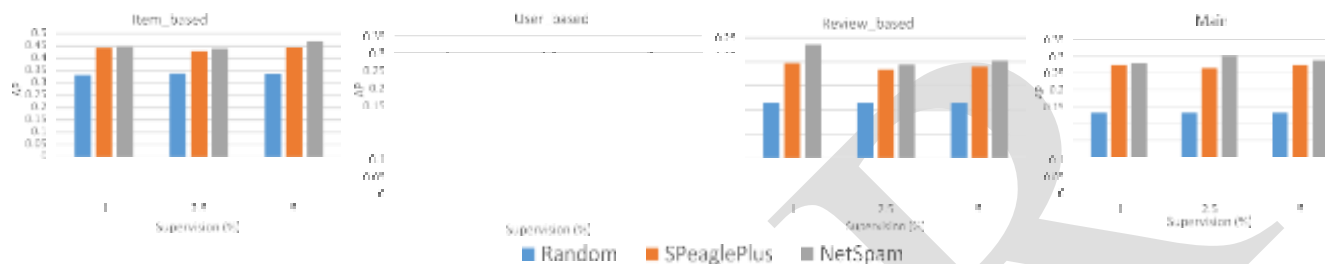


Fig.3.APforRandom,SPeaglePlusandNetSpamapproachesindifferentdatasetsandsupervisions(1%,2.5%and5%).



Fig.4.AUCforRandom,SPeaglePlusandNetSpamapproachesindifferentdatasetsandsupervisions(1%,2.5%and5%).

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 used different supervisions and we train a set for weight calculation. We also engaged these supervisions as fundamental labels for reviews which are chosen as a training set.

### C. Main Results

In this section, we evaluate NetSpam from different perspectives and compare it with two other approaches, Random approach and SPeaglePlus [36]. To compare with the first one, we have developed a network in which reviews are connected to each other randomly. Second approach uses a well-known graph-based algorithm called as "LBP" to calculate final labels. Our observations show NetSpam outperforms these existing methods. Then analysis on our observation is performed and finally we will examine our framework in unsupervised mode. Lastly, we investigate time complexity of the proposed framework and the impact of camouflage strategy on its performance.

1) *Accuracy*: Figures 3 and 4 present the performance in terms of the AP and AUC. As it is shown in all of the four



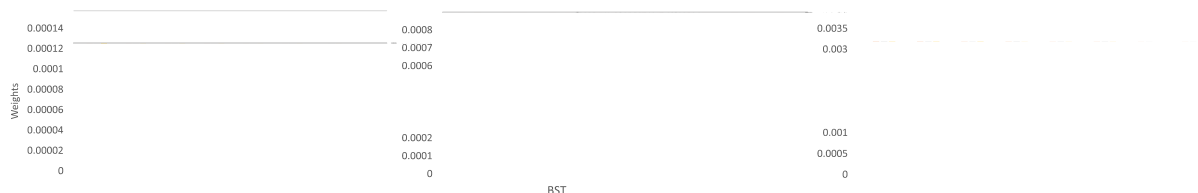


Fig.5.Features weights for NetSpam framework on different datasets using different supervisions (1%, 2.5% and 5%).

datasets NetSpam outperforms SPeagle Plus specially when number of features increase. In addition different supervisions have no considerable effect on the metric values neither on NetSpam nor SPeagle Plus. Results also show the datasets with higher percentage of spam reviews have better performance because when fraction of spam reviews in a certain dataset increases, probability for a review to be a spam review increases and as a result more spam reviews will be labeled as spam reviews and in the result of AP measure which is highly dependent on spam percentage in a dataset. On the other hand, AUC measure does not fluctuate too much, because this metric is not dependent on spam reviews percentage in dataset, but on the final sorted list which is calculated based on the final spam probability.

2) *Feature Weights Analysis*: Next we discuss about features weights and their involvement to determine spamicity. First we inspect how much AP and AUC are dependent on 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 run 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.

a) *Dataset impression on spam detection*: As we explained previously, different datasets yield different results based on their contents. For all datasets and most weighted features, there is a certain sequence for features weights. As is shown in Fig. 5 for four datasets, in almost all of them, features for the Main dataset have more weights and features for Review-based dataset stand in the second position. Third position belongs to User-based dataset and finally Item-based dataset has the minimum weights (for at least the four features with most weights).

b) *Features weights importance*: As shown in Table IV, there are 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, *N R*. From the third feature to the last feature there are different order for 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

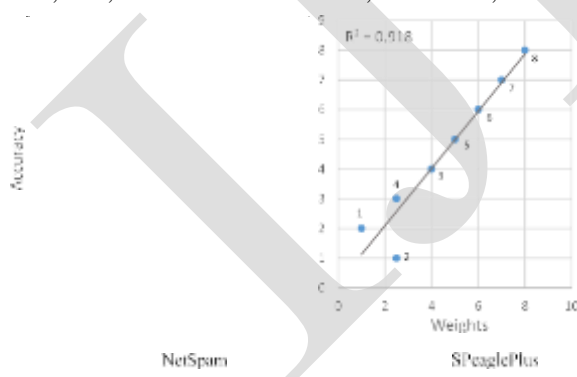


Fig. 6. Regression graph of features vs. accuracy (with 5% data as trainset) for Main dataset. (see Table II for numbers).

at rank 3. Going further, we see in the Review-based dataset all four most weighted features are behavioral-based features which show how much this type of features are important in detecting spams as acknowledged by other works as well [37], [38].

As we can see in Fig. 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 features 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 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	TF	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



Fig.7. Features weights for different features categories (RB, UB, RL and UL) with 5% supervision, on different datasets.

of AP are completely correlated. We observed values 0.958 ( $p$ -value=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. Similar result for SPEagle Plus also shows our weight 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 to a dataset used for extracting weights of features.

c) *Features category analysis*: As shown in Fig. 7 there are four categories with different weights average which is very important, specially in determining which feature is more appropriate for spotting spam reviews (refer to Sec. IV-C.2.b). resimilar 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 is confirmed by [39] and [16]. Analysis on separate views show that review-based features have higher weights which leads to better performance. It is worth to mention that none of previous work have investigated this before. Same analysis on the Main dataset shows equal importance of both category in findings spams. On the Other hand, in the first three dataset from Table I, RB has better weights (a bit difference 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.

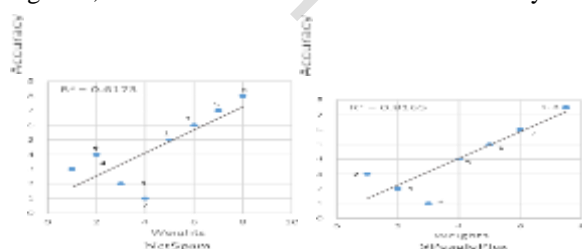


Fig.8. Regression graph of features vs. accuracy (unsupervised) for Main dataset. (see Table II for numbers).

3) *Unsupervised Method*: One of the achievement 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 unsupervised approach special formulation is used to calculate fundamental labels and next 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 SPeagle Plus this value reaches 0.90 ( $p$ =0.0021). As another example for 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 observation indicates NetSpam can prioritize features for both frameworks. Table V demonstrates that there is certain sequence in feature weights and it means in spam detection problems, spammers and spam reviews have common behaviors, no matter what social network they are writing the review for: Amazon or Yelp. For all of them, *DEV* is most weighted features, followed by *NR*, *ET* and *BST*.

4) *Time Complexity*: If we consider the Main dataset as input to our framework, time complexity with these circumstances is equal to  $O(e^2m)$  where  $e$  is number of edges in

TABLE V

WEIGHTS OF ALL FEATURES (USING UNSUPERVISED APPROACH); FEATURES ARE RANKED BASED ON THEIR OVERALL AVERAGE WEIGHTS

Dataset	DEV	NR	ET	BST	RES	PPI	MCS
M	0.0029	0.0550	0.0484	0.0445	0.0379	0.0329	0.0314
Review	0.0626	0.0510	0.0477	0.0376	0.0355	0.0346	0.0340
Item	0.0638	0.0510	0.0501	0.0395	0.0388	0.0383	0.0366
User	0.0630	0.0514	0.0494	0.0380	0.0373	0.0377	0.0366
Am	0.1102	0.0897	0.0746	0.0689	0.0675	0.0624	0.0297

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  $O(e^2)$  and this checking must be repeated for every feature. So, our time complexity for offline mode in which we give the Main dataset to framework and calculate spamicity of whole reviews, is  $O(e^2m)$  where  $m$  is number of features. In online mode, a review is given to NetSpam to see whether it is spam or not, we need to check if there is a metapath between given review with other reviews, which is in  $O(e)$ , and like offline mode it has to be repeated for every feature and every value. Therefore the complexity is  $O(em)$ .

5) *The Impact of Camouflage Strategy*: One of the challenges that spam detection approaches face is that spammers often write non-spam reviews to hide their true identity known as camouflage. For example they write positive reviews for good restaurant or negative reviews for low-quality ones; hence every spam detector system fails to identify this kind of spammers or at least has some trouble to spot them. In the previous studies, there are different approaches for handling this problem. For example, in [12], the authors assume there is always a little probability that a good review written by a spammer and put this assumption in its compatibility matrix. In this study, we tried to handle this problem by using weighted metapaths. In particular, we assume that even if a review has a very little value for a certain feature, it is considered in feature weight calculation. Therefore, instead of considering metapaths as binary concepts, we take 20 values which denoted as  $s$ . Indeed, if there is a camouflage its affection will be reduced. As we explained in Section III-C in such problems it is better to propose a fuzzy framework, rather than using bipolar values (0,1).

## V. RELATED WORKS

In the last decade, a great number of research studies focus on the problem of spotting spammers and spam reviews. However, since the problem is non-trivial and challenging, it remains far from fully solved. We can summarize our discussion about previous studies in three following categories.

### A. Linguistic-Based Methods

This approach extracts linguistic-based features to find spam reviews. Feng et al. [13] use *unigram*, *bigram* and their composition. Other studies [4], [6], [15] use other features like pairwise features (features between two reviews; e.g. content similarity), percentage of CAPITAL words in reviews for finding spam reviews. Lai et al. in [33] use a probabilistic language modeling to spot spam. This study demonstrates that 2% of reviews written on business websites are actually spam.

### B. Behavior-Based Methods

Approaches in this group almost use review metadata to extract features; those which are normal patterns of reviewer behaviors. Feng et al. in [21] focus on distribution of spammers rating on different products

and trace them. In [34], Jinda *et al.* extract 36 behavioral features and use a supervised method to find spammers on Amazon and [14] indicates behavioral features show spammers' identity better than linguistic ones. Xue *et al.* in [32] use rated deviation of a specific user and use a trust-aware model to find the relationship between users for calculating final spamicity score. Minnich *et al.* in [8] use temporal and location features of users to find unusual behavior of spammers. Li *et al.* in [10] use some basic features (e.g. polarity of reviews) and then run a HNC (Heterogeneous Network Classifier) to find final labels on Dianping dataset. Mukherjee *et al.* in [16] almost engage behavioral features like rated deviation, extremity and etc. Xie *et al.* in [17] also use a temporal pattern (time window) to find singleton reviews (reviews written just once) on Amazon. Luca and Zervas in [26] use behavioral features to show increasing competition between companies leads to very large expansion of spam reviews on products.

Crawford *et al.* in [28] indicates using different classification approach need different number of features to attain desired performance and propose approaches which use fewer features to attain that performance and hence recommend to improve their performance while they use fewer features which leads them to have better complexity. With this perspective our framework is arguable. This study shows using different approaches in classification yield different performance in terms of different metrics.

### C. Graph-Based Methods

Studies in this group aim to make a graph between users, reviews and items and use connections in the graph and also some network-based algorithms to rank or label reviews (as spam or genuine) and users (as spammer or honest). Akoglu *et al.* in [11] use a network-based algorithm known as LBP (Loopy Belief Propagation) in linearly scalable iterations related to number of edges to find final probabilities for different components in network. Fei *et al.* in [7] also use same algorithm (LBP), and utilize burstiness of each review to find spammers and spam reviews on Amazon. Li *et al.* in [10] build a graph of users, reviews, users IP and indicates users with same IP have same labels, for example if a user with multiple different account and same IP writes some

reviews, they are supposed to have same label. Wang *et al.* in [18] also create network of users, reviews and items and use basic assumptions (for example a reviewer is more trustworthy if he/she writes more honest reviews) and label reviews. Wahyuni and Djunaidy in [37] proposes a hybrid method for spam detection using an algorithm called ICF++ which is an extension to ICF of [18] in which just review rating are used to find spam detection. This work use also sentiment analysis to achieve better accuracy in particular.

Deeper analysis on literature show that behavioral features work better than linguistic ones in term of accuracy they yield. There is a good explanation for that; in general, spammers tend to hide their identity for security reasons. Therefore they are hardly recognized by reviews [38] they write about products, but their behavior is still unusual, no matter what language they are writing. In result, researchers combined both feature types to increase accuracy of spam detection. The fact that adding each feature is a time consuming process, this is where feature importance is useful. Based on our knowledge, there is no previous method which engage importance of features (known as weights in our proposed framework);

NetSpam in the classification step. By using these weights, on one hand we involve features importance in calculating final labels and hence accuracy of NetSpam increase, gradually. On the other hand we can determine which feature can provide better performance in term of their involvement in connecting spam reviews (in proposed network).

## VI. CONCLUSION

This study introduces a novel spam detection framework namely NetSpam based on a metapath concept as well as a new graph-based method to label reviews relying on a rank-based labeling approach. The performance of the proposed framework is evaluated by using two real-world labeled datasets of Yelp and Amazon websites. Our observations show that calculated weights by using this metapath concept can be very effective in identifying spam reviews and leads to a better performance. In addition, we found that even without a trainset, NetSpam can calculate the importance

of each feature and it yields better performance in the features' addition process, and performs better than previous works, with only a small number of features. Moreover, after defining four main categories for features our observations show that the reviews-behavioral category performs better than other categories, in terms of AP, AUC as well as in the calculated weights. The results also confirm that using different supervisions, similar to the semi-supervised method, have no noticeable effect on determining most of the weighted features, just as in different datasets.

For future work, metapath concept can be applied to other problems in this field. For example, similar framework can be used to find spammer communities. For finding community, reviews can be connected through group spammer features (such as the proposed feature in [39]) and reviews with highest similarity based on metapath concept are known as communities. In



addition, utilizing the product features is an interesting future work on this study as we used features more related to spotting spammers and spam reviews.

Moreover, while single network has received considerable attention from various disciplines for over a decade, information diffusion and content sharing in multilayer networks is still a young research [37]. Addressing the problem of spam detection in such networks can be considered as a new research line in this field.

## REFERENCES

- [1] J. Donfro, *A Whopping 20% of Yelp Reviews are Fake*, accessed on Jul. 30, 2015. [Online]. Available: <http://www.businessinsider.com/20-percent-of-yelp-reviews-fake-2013-9>
- [2] M. Ott, C. Cardie, and J. T. Hancock, "Estimating the prevalence of deception in online review communities," in *Proc. ACM WWW*, 2012, pp. 201–210.
- [3] M. Ott, Y. Choi, C. Cardie, and J. T. Hancock, "Finding deceptive opinion spam by any stretch of the imagination," in *Proc. ACL*, 2011, pp. 309–319.
- [4] C. Xu and J. Zhang, "Combating product review spam campaigns via multiple heterogeneous pairwise features," in *Proc. SIAM Int. Conf. Data Mining*, 2014, pp. 172–180.
- [5] N. Jindal and B. Liu, "Opinion spam and analysis," in *Proc. WSDM*, 2008, pp. 219–230.
- [6] F. H. Li, M. Huang, Y. Yang, and X. Zhu, "Learning to identify review spam," in *Proc. 22nd Int. Joint Conf. Artif. Intell. (IJCAI)*, 2011, pp. 1–6.
- [7] G. Fei, A. Mukherjee, B. Liu, M. Hsu, M. Castellanos, and R. Ghosh, "Exploiting burstiness in reviews for review spammer detection," in *Proc. ICWSM*, 2013, pp. 1–10.
- [8] A. J. Minnich, N. Chavoshi, A. Mueen, S. Luan, and M. Faloutsos, "Trueview: Harnessing the power of multiple review sites," in *Proc. ACM WWW*, 2015, pp. 787–797.
- [9] B. Viswanath et al., "Towards detecting anomalous user behavior in online social networks," in *Proc. USENIX*, 2014, pp. 1–16.
- [10] H. Li, Z. Chen, B. Liu, X. Wei, and J. Shao, "Spotting fake reviews via collective positive-unlabeled learning," in *Proc. ICDM*, Dec. 2014, pp. 899–904.
- [11] L. Akoglu, R. Chandy, and C. Faloutsos, "Opinion fraud detection in online reviews by network effects," in *Proc. ICWSM*, 2013, pp. 1–10.
- [12] S. Rayana and L. Akoglu, "Collective opinion spam detection: Bridging review networks and metadata," in *Proc. ACM KDD*, 2015, pp. 1–10.
- [13] S. Feng, R. Banerjee, and Y. Choi, "Syntactic stylometry for deception detection," in *Proc. 50th Annu. Meeting Assoc. Comput. Linguistics (ACL)*, 2012, pp. 1–5.
- [14] N. Jindal, B. Liu, and E.-P. Lim, "Finding unusual review patterns using unexpected rules," in *Proc. ACM CIKM*, 2012, pp. 1–4.
- [15] E.-P. Lim, V.-A. Nguyen, N. Jindal, B. Liu, and H. W. Lauw, "Detecting product review spammers using rating behaviors," in *Proc. ACM CIKM*, 2010, pp. 1–10.
- [16] A. Mukherjee et al., "Spotting opinion spammers using behavioral footprints," in *Proc. ACM KDD*, 2013, pp. 1–9.
- [17] S. Xie, G. Wang, S. Lin, and P. S. Yu, "Review spam detection via temporal pattern discovery," in *Proc. ACM KDD*, 2012, pp. 823–831.
- [18] G. Wang, S. Xie, B. Liu, and P. S. Yu, "Review graph based online store review spammer detection," in *Proc. IEEE ICDM*, Dec. 2011, pp. 1242–1247.
- [19] Y. Sun and J. Han, "Mining heterogeneous information networks: Principles and methodologies," in *Proc. ICC CE*, 2012, p. 159.



- [20] A.Mukherjee,V.Venkataraman,B.Liu,andN.Glance,“Whatyelpfake reviewfiltermightbedoing?”in*Proc.ICWSSM*,2013,pp.409–418.
- [21] S. Feng, L. Xing, A. Gogar, and Y. Choi, “Distributional footprints ofdeceptiveproductreviews,”in*Proc.ICWSSM*,2012,pp.98–105.
- [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*Proc.VLDB*,2011,pp.1–12.
- [23] Y. Sun, J. Han, P. Zhao, Z. Yin, H. Cheng, and T. Wu, “RankClus: Inte-grating clustering with ranking for heterogeneous information networkanalysis,” in *Proc. 12th Int. Conf. Extending Database Technol., Adv.DatabaseTechnol.*,2009,pp.1–12.
- [24] C.Luo,R.Guan,Z.Wang,andC.Lin,“HetPathMine:Anoveltransduc-tive classification algorithm on heterogeneous information networks,” in*Proc.ECIR*,2014,pp.210–221.
- [25] R. Hassanzadeh, “Anomaly detection in online social networks: Usingdatamining techniques and fuzzy logic,” School Elect. Eng. Comput.Sci.,QueenslandUniv.Technol.,Brisbane,QLD,Australia,Nov.2014.
- [26] M.LucaandG.Zervas,“Fakeittillyoumakeit:Reputation,competition,and yelpreviewfraud,”*Manage.Sci.*, vol.62, no. 16,pp.3412–3427,Jan.2016.
- [27] VM Allocation Technique and Optimized Performance Improvement for the Cloud Architecture. Authors: Dr. Md. Rafeeq, N. Navneetha, Dr. N. Subhash Chandra, M. Bhargavi,Dr.KRajeshwarRao
- [29] Vempati Krishna , G Sumalatha , A L Sreenivasulu , M Bhargavi “Robust and Effective Spam Filtering From Online Social Media Networks Using Machine Learning Strategies”,Published in Journal of Innovation in Information Technology, Vol. 5(1), Jan–Jun 2021@ ISSN:2581-723X
- [28] Teja Sree, N., Sumalatha, G. (2021). Behavioural Analysis Based Risk Assessment in Online Social Networks. In: Kumar, A., Mozar, S. (eds) ICCCE 2020. Lecture Notes in Electrical Engineering, vol 698. Springer, Singapore. [https://doi.org/10.1007/978-9815-7961-5\\_25](https://doi.org/10.1007/978-9815-7961-5_25)
- [29] Swathi Mattaparthi, Sheo Kumar, Mrutyunjaya S Yalawar, Fake Currency Detection: A Survey on Different Methodologies Using Machine Learning Techniques, 2022/4/29, International Conference on Communications and Cyber Physical Engineering 2018, 463-468, Publisher, Springer Nature Singapore. [https://doi.org/10.1007/978-981-19-8086-2\\_45](https://doi.org/10.1007/978-981-19-8086-2_45).
- [30] Sheo Kumar, Mrutyunjaya S Yalawar, K-NN SEMANTIC Inquiry on Scrambled Social Information BASE, ICDSMLA 2019: Proceedings of the 1st International Conference on Data Science, Machine Learning and Applications,982-991, 2020,Springer Singapore.
- [31] K Vijaya Babu, Mrutyunjaya S Yalawar, G Sumalatha, G Ramesh Babu, Ravi Kumar Chandu, An Overview of Various Security Issues and Application Challenges of the Attacks in Field of Blockchain Technology, 2022/5/16, ICCCE 2021: Proceedings of the 4th International Conference on Communications and Cyber Physical Engineering, 365-374, Springer Nature Singapore.
- [32] Mrutyunjaya S Yalawar, K Vijaya Babu, Bairy Mahender, Hareran Singh, A Brain-Inspired Cognitive Control Framework for Artificial Intelligence Dynamic System,2022/4/29,International Conference on Communications and Cyber Physical Engineering 2018,735-745,Springer Nature Singapore. [https://doi.org/10.1007/978-981-19-8086-2\\_70](https://doi.org/10.1007/978-981-19-8086-2_70).

- [33] Mrutyunjaya S Yalawar, L Umasankar, Mr Bottu Gurunadha Rao, RESIDUE ANALYSIS OF IMAGE QUALITY USING TENSORFLOW, JOURNAL OF CRITICAL REVIEWS, 7,ISSUE 04,1742-1746,2020
- [34] Kumar, S., Sai Lavanya, G.V.R. (2020). Face Recognition Using Open CV with Deep Learning. In: Kumar, A., Paprzycki, M., Gunjan, V. (eds) ICDSMLA 2019. Lecture Notes in Electrical Engineering, vol 601. Springer, Singapore. [https://doi.org/10.1007/978-981-15-1420-3\\_122](https://doi.org/10.1007/978-981-15-1420-3_122).
- [35] L. Pullagura, N. Kittad, G. Diwakar, V. Sathiya, A. Kumar and M. S. Yalawar, "ML based Parkinson's Disease Identification using Gait Parameters," 2022 International Conference on Automation, Computing and Renewable Systems (ICACRS), Pudukkottai, India, 2022, pp. 561-566, doi: 10.1109/ICACRS55517.2022.10029281.
- [36] Dr. C.N. RAVI , Karthikeyan Udaichi<sup>1</sup> Miguel Garcia-Torres<sup>3</sup>,Parameshchari Bidare Divakarachari<sup>4</sup> Large-scale system identification using self-adaptive penguin search algorithm, IET Control Theory & Applications,DOI: 10.1049/cth2.12479: , Received: 26 November 2022 Revised: 14, March 2023
- [37] Shrivastava, R., Jain, M., Vishwakarma, S.K., Bhagyalakshmi, L., Tiwari, R. (2023), "Cross-Cultural Translation Studies in the Context of Artificial Intelligence: Challenges and Strategies". In: Kumar, A., Mozar, S., Haase, J. (eds) Advances in Cognitive Science and Communications. ICCCE 2022. Cognitive Science and Technology. Springer, Singapore, ISBN: 978-981-19-8086-2\_9, pp 91-98. [https://doi.org/10.1007/978-981-19-8086-2\\_9](https://doi.org/10.1007/978-981-19-8086-2_9)
- [38] TarunDharDiwan, Rajesh Tiwari and VivekDubey, "Local Binary Pattern Occurrence Map Method for High Parallel Image Processing", International Conference on Advances in Computing and Communication held at NIT Hamirpur, Himachal Pradesh, India on 8 -10 April 2011, pp 538 – 540, ISBN: 978-81-920874-0-5
- [39]Naveen Kumar Sahu and Rajesh Tiwari, "Comparative Analysis of Optimization Algorithms based on Hybrid Soft Computing Algorithm", International Journal for Scientific Research & Development, Vol. 3, Issue 09, 2015, pp 33 – 39, ISSN (online): 2321-0613



**ISSN 2277-2685**

**IJESR/Sep. 2023/ Vol-13/Issue-3/1-19**

**Mr.Mrutyunjaya Yalawar *et. al.*, / International Journal of Engineering & Science Research**

IJESR



**ISSN 2277-2685**

**IJESR/Sep. 2023/ Vol-13/Issue-3/1-19**

**Mr.Mrutyunjaya Yalawar *et. al.*, / International Journal of Engineering & Science Research**

[View publication stats](#)

IJESR