

SPSB: Semi-Permeable Spam Banner

Parth S. Shah¹, Vishal R. Andodariya²

¹PG Student, ²Asst. Professor

Department of Computer Engineering

Shri Pandit Nathulalji Vyas Technical Campus, Wadhwan - India

Abstract— In the past few years Online Social Networking has gained importance because of its cheap & easily accessible features. Various online communication media platforms such as Social Networking, Electronic Mail Services (i.e. eMails) and various online chat messaging applications allows online users to reach & communicate with number of people at a barely negligible & affordable cost. Apart from this, it also enables the information to roam freely in network. Thus resulting into undesired communication happening on a larger scale. With increase in the number of eMail users there is a significant increase in irrelevant messages called Spam, Junk Mail or Ham. Hence, the objective is to develop a Semi Permeable System which bans the spams from entering the network with great efficiency hence, preventing Unwanted Communication and reducing cost of the network too. Various remedies taken earlier to eradicate this problem use filtering techniques based on the content, header & classification. Hence, SPSB is a proposed solution which will minimize the routing of irrelevant content in the network. It will act as a Semi Permeable Banner at the client side which will ban the Spams to further roam in the network and allow legitimate mails to enter the system. Hence SPSB will be a good authentication standard which will differentiate between mails which are legitimate & the ones that are not.

Index Terms—Spam Mail; SPSB (Semi Permeable Spam Banner); Legitimacy; Unwanted Communication; Ham; Spam Filtering; Types of Spam.

I. INTRODUCTION

This generation can be named as the Online Generation. With increase in the facilities provided over the net, there is a significant increase in the number of Internet users. The Online media has proved to be a one stop solution to all problems lately. With its immensely popular and easily accessible facilities the Internet Applications are full heartedly accepted by the users worldwide. In this current generation it is next to impossible to find people who are novice to the Internet applications.

Various Internet based communication platforms are Social Networking Websites & Applications, Electronic Mails (i.e. eMails), and Content Sharing sites, Shopping Portals, Video Portals, Search Engines, Blogging Sites and many more. These facilities are provided and available globally 24*7. The active Internet Users worldwide in 2014 are 2.5 Billion^[15] approx.

This concludes that social media has now become an engrained part of the lives of many people across different demographic regions. This increased ubiquity may result in some changes to the specific demographic bases of individual platforms, but even if people's habits are changing, it appears that people are moving from one social platform to another, rather than deserting social media in its entirety.

The most important feature of internet is anytime anywhere availability. This enables users to access the desired information 24*7 with a prime requisite of data internet connection. With the use of an inexpensive Internet connection any user has the potential to reach millions of users by posting messages to an email list or by uploading content to a sharing site. This feature comes with a drawback that it has a Democratic Content publication, as it gives the freedom to post anything with a visibility to all global users.

As a result, this same attribute is used with negative harmful intentions of interference of legitimate communication, unsolicited marketing messages, enter unauthenticated areas and fetch important information. All these irrelevant, unwanted and unsolicited communications are known as **Spam Communications**.

The Internet Spam Communications can be classified in the following categories:

- **Social Network Spam**
- **Email Spam**
- **Image Spam**
- **Click Spam**
- **Content Spam**
- **Cloaking & Redirection Spam**
- **Link Spam**

From the 2.5 billion^[15] approx. Internet Users worldwide there are approximately 1.06 billion^[16] eMail Users who communicate via eMails on daily basis. The sad part being that 69.6%^[17] of the eMails used for communication this year were Spams.

The cost incurred to the society due to spams is around \$200 Billion^[18] every year. According to Spam Experts Justin Rao from the Microsoft Research and David Riley from the Google, the cost of spam incurred was in a 100 to 1 ratio.

Hence the objective is to devise a more optimized **Selectively Permeable Spam Banner** which bans the illegitimate communications and decrease the cost caused on the network bandwidth & protect users from unsolicited communication. It is also desirable that SPSB never stops the legitimate senders and never catch legitimate mail as a spam mail.

II. PREVIOUS APPROACHES TO COUNTER SPAM

Unwanted, illegitimate & unsolicited communication has been a major problem in the form of Spam. There have been many attempts & measures taken to eliminate this unwanted communication. Various Spam Filters previously made are:

TABLE 1: SPAM FILTER MANUFACTURERS

Manufacturer	Product name
Sonic Wall	Email Security Appliances
Symantec	Norton Anti-Spam
Google	Mail
spamcop.net	Spam Cop
Apple	Mac Mail
mozilla.com	Thunderbird built-in filter
Microsoft	Exchange Server spam filter
McAfee	Spam Killer 6

Usually all the measures taken to build Spam Filters follow either of these techniques:

- **Content Rating Approach**
- **Header Based Approach**
- **Content Based Approach**
- **Protocol Based Approach**

A. Content Rating Approach

The Content Rating Filtering ^[1] approach is used by many content-sharing sites (e.g., YouTube ^[21]). In this approach, the users have the facility to rate the content they watch on parameters such as the level of interest, relevance of the content, and legitimacy of content item they have viewed. Furthermore, the content is tagged with the users' ratings, likes & dislikes. This feature of text mining is used in many content sharing sites to check the legitimacy & appropriateness of the content.

This approach would guide the user to differentiate between the legitimate content and help him discard unwanted data. Apart from this, this rating based approach also enables the administrators to monitor illegitimate accounts & discard or remove them. It is mainly used for One to Many i.e. Broadcasting Communication systems. The drawback here is that it enables to filter & differentiate between the legitimate & illegitimate content but doesn't remove it from the network. Hence adding more to the bandwidth & network cost.

B. Header Based Approach

This approach works on a mechanism to fetch & examine the Headers of eMail messages for identifying its legitimacy. It has a concept named **Blacklist**. This blacklist will work as a database which stores IP addresses of all known Spammers (detected previously) and doesn't accept any eMail messages arriving from those IP addresses. Furthermore, it also has a **Whitelist**. Whitelist is again a database which according to the admin are Non Spammers & legitimate eMail senders. This is done to decrease the number of checks every time an eMail is arrived.

To attain maximum accuracy the user can manually create his own Blacklist & Whitelist. The drawback here is that it is quite a burdening task to maintain & update the list regularly. There is an alternate automatic list creator approach which uses results from previous results. It is named as *autowhitelists* in Spam Assassin ^[9]. It is hard to maintain both Blacklists and Whitelists because IP addresses can be forged easily by Spammers & Hackers.

C. Content Based Approach

The mechanism used by this approach is to analyse the subject content of the eMail message and search for certain **Keywords** (Generated using a Bayesian filter or by the statistics provided) or **Patterns** which look like a typical spam mail (e.g., URLs with numeric IP addresses in the email body, Various eMails showing fake prize money award, etc.).

The advantage about this approach is its ability to filter Spam mails to a greater extent. The drawback it faces is that it needs to update the set of keywords regularly because the spammers won't use the same techniques again and again ^[2].

D. Protocol Based Approach

The mechanism used by this approach is to add an authentication check in the underlying eMail Protocol. It has a Challenge-Response scheme wherein it require a manual effort from the sender to send the first email to a particular recipient. For example, the sender has to go to a certain web page and activate the email manually, which might involve answering a simple question (such as solving a simple mathematical equation).

Once this has been done, the recipient would check his authentication and then the sender will be added to the recipient's whitelist. Once the sender has been added to the whitelist the next communication via emails can be done without the activation procedure. The benefit of this approach is that it would decrease the amount of spams on larger scale because it would require the Spammer to pass through the activation task with every new user. Hence, the option of sending spam eMails to millions would be eliminated.

III. LITERATURE REVIEW

Many approaches have been conducted by various researchers to counter spam because of its parasitic nature of

increasing day by day. There were different approaches wherein some proved fruitful while others didn't. Few Such techniques are shown below:

A. SOAP: Social Network aided personalized and effective spam filter to clean your eMail inbox.

At current may spam filters uses social networks itself to monitor spam detection. To develop the perfect spam filter this paper lightens the way. They proposed a new filter called SOAP: That is n/w aided spam filter.

As seen in previous papers many of filters (Bayesian) emphasis on static keywords or lists (Black or White). Unlike many of filters SOAP not depends on a single methods to filter spams rather it uses more than one technique to filter spams. The system integrates trust management, social relations and basic one that is Bayesian filter.

This system also checked with real dataset of Facebook profiles, which includes both regular and spam profiles. The system proves better to scan the spams.

B. Detecting Spammer on Twitter

This paper discuss about to deal with spammers on Twitter. To cope with spams on Twitter manually classified the legitimate users and spammers. For that real dataset of Twitter about 54 million users is collected, along with 1.9 billion links, and almost 1.8 billion tweets.

To detect spammers they identifies number of attributes or behaviours related to content and behaviour. This is very much useful to detect spammers. To detect spammers or non-spammers uses this attributes to MLP (Machine Learning Process) for classification.

This strategy succeeds to detect irrelevant data or spam data (content) with great percentage approx. 70% of dummies and 96% of regular one.

C. Mail Rank: rank based Spam Detection

This technique uses ranking system to rate the emails which are arrived. As a result from that rank sender can be identified as spam or non-spam. There are two possibilities for Mail-Rank system:

- **Basic Mail-Rank**, which calculates an overall (global) rank for every mail address.
- **Personalized Mail-Rank**, in which for every mail address score is different.

The system, Mail-Rank is very much reliable and highly resistant against spam attack. In sparse network, the network of a small set of peers, Mail-Rank can also performs well.

D. Personalized eMail Network

To find trusted networks of friends in cyberspace personal email network provide automated graph theoretic method. Network keeps history of users. Mail user can use their mail network to differentiate irrelevant or can say unsolicited mail, named spam. Now this mail network is generally constructed from historical information available in the header of email.

Paper focus to construct a trusted like of network in which network must know about all the users resides in the network. This personalized network thus helps to identifies legitimate data and spam data. With 100% accuracy, algorithm of this tool can classify approx. 53% of all emails as spam or non-spam.

E. Detecting Spammer using SNARE

SNARE is the type of reputation engine that uses more than one method to classify spammers and non-spammers. The regular spam filtering technique like listing is not easy to maintain and error prone also if attacker attacks on lists.

SNARE examines features rather than contents that's why it is very much lightweight. They incorporate this feature in classification algorithm and tests whether it can classify as spammer or legitimate one. SNARE is build using this feature kept in mind. This engine can be used as first pass in the blacklists.

F. Markov Clustering Approach

The study is based on a real dataset of Facebook profiles, which includes both regular and spam profiles. Paper uses weighted graph technique to model social network as profiles are represented as nodes and their interactions represented as edges. To calculate weight of an edge which connects user profiles as a pair is calculated as a function of the real social interactions in terms of shared URLs, page likes also active friends in the network.

MCL is applied on the weighted graph to generate different clusters containing different categories of profiles. Majority voting is applied to handle the cases in which a cluster contains both spam and normal profiles.

Experimental results of this paper show that majority voting not only reduces the number of clusters to a minimum, but also increases the performance.

IV. SPSB: SEMI PERMEABLE SPAM BANNER

The problem definition has made it clear that Spam Filtering is yet an area which needs proper optimization. Having studied all the related work done and the approaches taken to filter spam, I realized that there is still more scope for further improvement & optimization in Spam filtering.

Hence a Proposed Solution to this problem definition would be **SPSB: Semi Permeable Spam Banner**.

The parameter that differentiates SPSB from other approaches is that it would be implemented on the Mail Server i.e. Server side instead of prior approaches which were implemented on the client side. Let us understand the architecture of SPSB in the below figure:

A. The SPSB architecture is basically divided into three parts:

1. **Cluster creation**
2. **Main User Node**
3. **Authenticator**

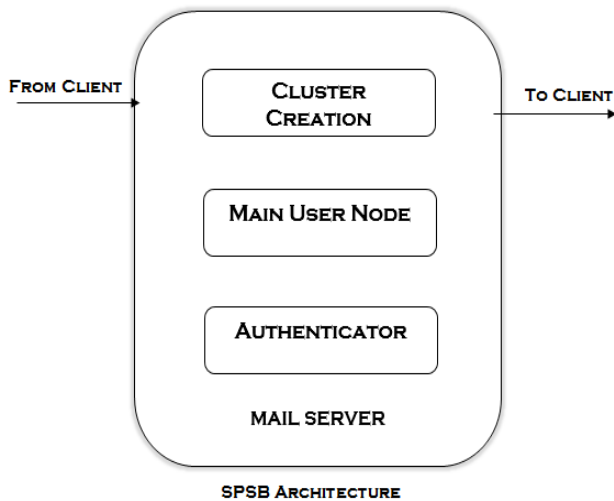


Fig 1: SPSB Architecture

1. Cluster creation:

The Cluster creation is done on the basis of dividing the members into two types of friends. They are named as:

- **Direct Friend:** Nodes directly connected to the user are Direct Friends [DF].
- **Neighbor Friend:** Nodes connected to the Direct Friend of the User are called Neighbor Friends [NF].

The combination of both will create a **Cluster**. The NF will act as a guard of the group from outliers. The group creation should be done carefully for proper management. Direct Friends are friends who are present in the users' **Whitelist**.

2. Main User Node (MUN):

The **Main User Mode** is the most important part of this architecture because it is the main node in the cluster. The MUN is not the head of the cluster, neither will it manage the flow of mails of the other nodes in the cluster.

The MUN is selected on the basis of the node having maximum number of nodes in the entire cluster. MUN is selected carefully as it gives the reference of the all nodes resides in the group.

There are four occasions where MUN is selected:

- **Select MUN within the group.**
Node with the highest number of nodes as DF is select as MUN.
- **Select MUN for adjacent group.**
To select a Node which is outside the group, take help of NF in adjacent cluster. From many any one must be connected with any DF or NF. Thus with that it will suggest the MUN in adjacent cluster.
- **Select MUN for different MS group.**
Here after the previous process completes the MS of both sender and receiver communicate with each other to form a new group.
- **Select MUN for new node.**

Here for this problem first it starts communication with some authentication process. The authentication process may be any to check legitimacy of user. May ask questions, or to do calculation or to identify numbers etc.

Hence, MUN is an important parameter of the SPSB architecture & thus needs to be obtained carefully.

3. Authenticator:

It is this part of the SPSB Architecture that checks the legitimacy of the selection of the user amongst the group. This part of the architecture will classify the user type i.e. **Regular User, Spammer and New user**.

B. Proposed Algorithm

The SPSB will implement the SAA (Sender Authentication Algorithm). It uses the following steps to check the sender's authenticity:

Step 1:

DF[], NF[]; //List of Direct Friends and Neighbor Friends
Get the sender = sender_id

Step 2:

```
if (sender is in common friends (DF)) // Check whether
from own group
    then allow sender(msg) to go through
goto step 6
```

Step 3:

```
if (sender is in common neighbor_friends (NF)) // whether
from adjacent group
    then allow sender(msg) to go through
goto step 6
```

Step 4:

```
if (sender is not from current community circle)
    Check the validation (whether regular user or
    spammer) of the sender // Check whether it from validate or
    not
goto step 5
```

Step 5:

```
Check validation from MUN.
if validated then
    goto step 6
else discard the message, report about
```

sender

Step 6:

Send message to Receiver.

V. CONCLUSION

Prevention of Spam is an issue tried by many but yet has not profound its optimal solution. This problem increases with increase in users. Hence, with SPSB, a Semi Permeable Spam Banner a prevention filter is devised which will counter Spam messages by differentiating it with illegitimate & legitimate users.

With the help of Main User Node (MUN) it counters spams that are not from individuals' social circle. SPSB, will counter

spams with a higher reliability. This approach helps the user to be free from spam attacks and free from attacks of data which are totally irrelevant. It is also desirable that SPSB never stops the legitimate senders and never catch legitimate mail as a spam mail.

ACKNOWLEDGMENT

I am heartily thankful to my guide **Asst. Prof. Vishal R. Andodariya**, the person who made me follow the right steps during this research. I deeply express my gratitude to him for his friendly guidance, valuable suggestions and expertise at every phase of this research.

I would also like to extend my heartfelt gratitude to the Head of Our Department **Asst. Prof. Rushirajsinh Zala** for providing me this opportunity to work over this research and for his tremendously helpful nature throughout the research.

REFERENCES

A. PAPERS:

- [1]. Gray and M. Haahr. Personalised, Collaborative Spam Filtering. In *Proc. of the Conference on Email and Anti-Spam (CEAS)*, Mountain View, CA, USA, July 2004.
- [2]. Isode: Benchmark and comparison of Spamassassin and m-switch anti-spam, Technical report, Isode, April 2004.
- [3]. An MCL-Based Approach for Spam Profile Detection in Online Social Networks Faraz Ahmed, Muhammad Abulaish, IEEE 2012, Computing and Communications
- [4]. Z. Li and H. Shen. SOAP: A social network aided personalized and effective spam filter to clean your e-mail box. In *Proc. of IEEE INFOCOM*, 2011
- [5]. Preventing Unwanted Social Inferences with Classification Tree Analysis. Sara Motahari, Sotirios Ziavras, Quentin Jones, IEEE, 2009
- [6]. Shuang Hao, Nadeem Ahmed Syed, Nick Feamster, Er G. Gray, and Sven Krasser. Detecting spammers with SNARE: Spatio-temporal network-level automatic reputation engine. In *USENIX Security*, 2009.
- [7]. Paul Alexandru Chirita, Jörg Diederich, and Wolfgang Nejdl. Mail Rank: using ranking for spam detection. In *Proc. of CIKM*.
- [8]. Enrico Blanzieri and Anton Bryl, "A survey of learning based techniques of Email spam filtering", A Technical Report, University of Trento, 2008

- [9]. P. O. Boykin and V. Roy Chowdhury. Personal email networks: An effective anti-spam tool. *IEEE COMPUTER*, 2004.
 - [10]. S. Garriss, M. Kaminsky, M. J. Freedman, B. Karp, D. Mazières, and H. Yu. Re: Reliable email. In *Proc. of NSDI*, 2006.
 - [11]. A. Mislove, A. Post, P. Druschel, and KP Gummadi. Ostra: Leveraging trust to thwart unwanted communication. In *Proc. of NSDI*, 2008.
 - [12]. M. Sirivianos, K. Kim, and X. Yang. Introducing social trust to collaborative spam mitigation. In *Proc. of IEEE INFOCOM*, 2011.
 - [13]. M. Perone. An overview of spam blocking techniques. Technical report, Barracuda Networks, 2004.
 - [14]. J. Golbeck and J. Hendler Reputation Network Analysis for Email Filtering. In *Proc. of the Conference on Email and Anti-Spam (CEAS)*, Mountain View, CA, USA, July 2004.
- ##### B. WEBSITES:
- [15]. <http://wearesocial.net/blog/2014/01/social-digital-mobile-worldwide-2014/>
 - [16]. <http://emailclientmarketshare.com>
 - [17]. <http://www.zdnet.com/article/worldwide-spam-rate-falls-2-5-percent-but-new-tactics-emerge/>
 - [18]. <http://cleanmessage.com/index.php/cost-of-spam/>
 - [19]. http://www.cisco.com/web/offer/gist_ty2_asset/Cisco_2014_ASR.pdf
 - [20]. <http://royal.pingdom.com/2013/01/16/internet-2012-in-numbers/>
 - [21]. <http://krebsonsecurity.com/2013/01/spam-volumes-past-present-global-local/>
 - [22]. <https://www.youtube.com/>
 - [23]. <http://spamassassin.apache.org/>