

The Authentication Framework for Spam (Unwanted) Messages On Micro Blogging WebSite

Aher Prajkata P., Ahire Priti T., Awade Smita S., Shinde Ekta S.

Abstract— One fundamental issue in today On-line Social Networks (OSNs) is to give users the ability to control the messages posted on their own private space to avoid that unwanted content is displayed. Up to now OSNs provide little support to this requirement. To fill the gap, in this paper, we propose a system allowing OSN users to have a direct control on the messages posted on their walls. This is achieved through a flexible rule-based system, that allows users to customize the filtering criteria to be applied to their walls, and Machine Learning based soft classifier automatically labeling messages in support of content-based filtering.

Index Terms— On-line Social Networks, Information Filtering, Short Text Classification, Policy-based Personalization.

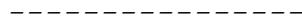
**1 INTRODUCTION**

ONLINE Social Networks (OSNs) are today one of the most popular interactive medium to share, communicate, and distribute an important amount of human living information. On a daily basis and continuous messages involve the swap of several types of content, including free content, image, audio, and video information. Along with Facebook information¹ average user creates 90 pieces of substance every month, while more than 30 billion quantity of substance (web links, news stories, notes, blog posts, photo albums, etc.) are distributed every month. The vast and dynamic character of this information produces the premise for the employment of web content mining strategies aimed to automatically discover useful information dormant contained by the information. They are instrumental to give a dynamic support in complex and sophisticated tasks involved in OSN administration, for example such as access power or information filtering. Information filtering has been significantly searched for what concerns textual documents and, more recently, web content.[1-3]

However, the aim of the majority of these proposals is mainly to provide users a classification mechanism to avoid they are overwhelmed by unsuccessful information. In OSNs, information filtering can also be exploited for a dissimilar, more responsive, purpose. This is due to the

fact that in OSNs there is the possibility of posting or commenting other posts on exacting public/private regions, called in common walls. Information filtering can therefore be used to provide users the capability to automatically control the messages written on their individual walls, by filtering out surplus communication. We believe that this is a key OSN service that has not been offered so far. Certainly, in the present day OSNs provide very tiny maintain to prevent unwanted messages on user walls. For instance, Facebook permits users to status who is allowed to insert messages in their walls (i.e., friends, defined groups of friends or friends of friends). Though, no content-based preferences are maintained and therefore it is not possible to prevent undesired messages, for instance political or offensive ones, no matter of the user who posts them. Providing this service is not only a topic of using previously defined web content mining methods for a different purposes, rather it entails to propose adhoc categorization strategies. This is because wall messages are represented by tiny text for which traditional classification methods have serious limitations since short texts do not provide sufficient word occurrences.

The aim of the present work is therefore to propose and experimentally evaluate an automated system, called Filtered Wall (FW), able to filter unwanted messages from OSN user walls. We exploit Machine Learning (ML) text categorization techniques [4] to automatically assign with each short text message a set of categories based on its substance. The most important efforts in building a robust small text classifier (STC) are concentrated in the extraction and selection of a set of characterizing and discriminant aspects. The resolutions examined in this paper are an extension of those adopted in a previous work by us [5] from which we inherit the learning model and the elicitation procedure for generating preclassified information.



• Aher Prajkata P., Ahire Priti T., Awade Smita S., Shinde Ekta S. is research scholar at Jayhind College of Engineering, Kuran, Pune.

2 RELATED WORK

The main contribution of this paper is the design of a system providing customizable content-based message filtering for OSNs, based on ML techniques. As we have pointed out in the introduction, to the best of our knowledge we are the first proposing such kind of application for OSNs [6] However, our work has relationships both with the state of the art in content-based filtering, as well as with the field of policy-based personalization for OSNs and, more in general, web contents. Therefore, in what follows, we survey the literature in both these fields.

A. Content-based filtering Information filtering systems are designed to classify a stream of dynamically generated information dispatched asynchronously by an information producer and present to the user those information that are likely to satisfy his/her requirements .In content-based filtering each user is assumed to operate independently. As a result, a content-based filtering system selects information items based on the correlation between the content of the items and the user preferences as opposed to a collaborative filtering system that chooses items based on the correlation between people with similar preferences While electronic mail was the original domain of early work on information filtering, subsequent papers have addressed diversified domains including newswire articles, Internet "news" articles, and broader network resources , Documents processed in content-based filtering are mostly textual in nature and this makes content-based filtering close to text classification. [7-9] The activity of filtering can be modeled, in fact, as a case of single label, binary classification, partitioning incoming documents into relevant and non-relevant categories. More complex filtering systems include multi-label text categorization automatically labeling messages into partial thematic categories. Content-based filtering is mainly based on the use of the ML paradigm according to which a classifier is automatically induced by learning from a set of pre-classified examples. A remarkable variety of related work has recently appeared, which differ for the adopted feature extraction methods, model learning, and collection of samples. The feature extraction procedure maps text into a compact representation of its content and is uniformly applied to training and generalization phases. Several experiments prove that Bag of Words (BoW) approaches yield good performance and prevail in general over more sophisticated text representation that may have superior semantics but lower statistical quality.[10-13]

3 PROJECT SCOPE

The aim of the present work is therefore to propose and experimentally evaluate an automated system, called Filtered Wall (FW), able to filter unwanted messages from OSN user walls. We exploit Machine Learning (ML) text categorization techniques to automatically assign with each short text message a set of categories based on its content

- Automatic filtering of unwanted messages
- Avoid messages from undesired creators
- To avoid overwhelming users of micro blogging services

4 SYSTEM FEATURES

Functionality Requirements

- **Log In**
 - **Description:** If the user is already registered then Login is the feature of the system which allows the user to enter into system and will provide him access to the system.
- **Registration:**
 - **Description:** If the user is new to the system then he need to register first for getting access to the system. Register will provide user a option for registration.
- **Select add friend menu**
 - **Description:** After click on add friend menu then he/she require clicking on add friend button.
- **View requested friend details**
 - **Description:** If user wants to add new friend then he/she needs to confirm the friend request.
- **View friend list**
 - **Description:** This menu provide friend list.
- **Approved blog**
 - **Description:** User can approved unwanted message with help of this blog.
- **Add new blog**
 - **Description:** This menu allows the user to add new blogs.

5 DICOMFW

DICOM Filtering Wall

Marco Vanetti's wall

The screenshot shows a Facebook-style wall interface. At the top, it says "Marco Vanetti's wall". Below that is a text input field with a "Send" button. A red dashed box highlights a message from "Moreno Carullo": "Ciao Marco, visto? Sembra tutto funzionare - Hi Marco, did you see? Everything seems to work". Below this message, there is a small note: "Your message can't be posted because it was filtered! Non-Neutral = 1 Neutral = 0 Violence = 0.316783 Vulgar = 1 Offensive = 0.221593 Hate = 0 Sex = 0.601016". Further down, another message from "Antonio Tirrendi" is shown: "Si scrive qui? Ciao, Antonio - Do i have to write here? See you, Antonio". Both messages have timestamped "Posted" details and "Reply", "Delete", and "Filtering metadata" links.

Fig. 1. DicomFW: a message filtered by the wall's owner FRs (messages in the screenshot have been translated to make them understandable)

DicomFW is a prototype Facebook application [8] that emulates a personal wall where the user can apply simple combination of the proposed FRs. Throughout the development of the prototype we have focused our attention only on the FRs, leaving BL implementation as a future improvement. However, the implemented functionality is critical, since it permits the STC and CBMF components to interact. Since this application is conceived as a wall and not as a group, the contextual information (from which CF are extracted) linked to the name of the group are not directly accessible. Contextual information that is currently used in the prototype is relative to the group name where the user that writes the message is most active. As a future extension, we want to integrate contextual information related to the name of all the groups in which the user participates, appropriately weighted by the participation level. It is important to stress that this type of contextual information is related to the environment preferred by the user who wants to post the message, thus the experience that you can try using DicomFW is consistent with what described and evaluated in Section VI-C.

To summarize, our application permits to:

- 1) view the list of users' FWs;
- 2) view messages and post a new one on a FW;
- 3) define FRs using the OSA tool.

When a user tries to post a message on a wall, he/she receives an alerting message (see Figure 3) if it is blocked by FW.

6 EVALUATION

In this section, we illustrate the performance evaluation study we have carried out the classification and filtering modules. We start by describing the dataset. A. Problem and Dataset Description. The analysis of related work has highlighted the lack of a publicly available benchmark for comparing different approaches to content based classification of OSN short texts. To cope with this lack, we have built and made available a dataset D of messages taken from Facebook. The dataset, called WmSnSec 2, is available online at www.dicom.uninsubria.it/~marco.vanetti/wmsnsec/

91266 messages from publicly accessible Italian groups have been selected and extracted by means of an automated procedure that removes undesired spam messages and, for each message, stores the message body and the name of the group from which it originates. The messages come from the group's web page section, where any registered user can post a new message or reply to messages already posted by other users. e. The role of the group's name within the text representation features was explained in Section IV-A. The set of classes considered in our experiments is $\{f_{Neutral}, f_{Violence}, f_{Vulgar}, f_{Offensive}, f_{Hate}, f_{Sex}\}$, where $f_{Neutral}$ are the second level classes. The percentage of elements in D that belongs to the Neutral class is 31%.

In order to deal with intrinsic ambiguity in assigning messages to classes, we conceive that a given message belongs to more than one class. Each message has been labeled by a group of five experts and the class membership values $y_{j1}, y_{j2}, \dots, y_{j5}$ for a given message m_j were computed by a majority voting procedure. After the ground truth collection phase, the messages have been selected to balance as much as possible second-level class occurrences. The group of experts has been chosen in an attempt to ensure high heterogeneity concerning sex, age, employment, education and religion. In order to create a consensus concerning the meaning of the Neutral class and general criteria in assigning multi-class membership we invited experts to participate to a dedicated tuning session. Issues regarding the consistency between the opinions of experts and the impact of the dataset size in ML classification tasks will be discussed and evaluated in Section VI-B. We are aware of the fact that the extreme diversity of OSNs content and the continuing evolution of communication styles create the need of using several datasets as a reference benchmark. We hope that our dataset will pave the way for a quantitative and more precise analysis of OSN short text classification methods.

7 CONCLUSIONS

In this paper, we have presented a system to filter undesired messages from OSN walls. The system exploits a ML soft classifier to enforce customizable content-

dependent FRs. Moreover, the flexibility of the system in terms of filtering options is enhanced through the management of BLs.

8 REFERENCES

- [1] A. Adomavicius and G. Tuzhilin, "Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions," *IEEE Trans. Knowledge and Data Eng.*, vol. 17, no. 6, pp. 734-749, June 2005.
- [2] M. Chau and H. Chen, "A Machine Learning Approach to Web Page Filtering Using Content and Structure Analysis," *Decision Support Systems*, vol. 44, no. 2, pp. 482-494, 2008.
- [3] R.J. Mooney and L. Roy, "Content-Based Book Recommending Using Learning for Text Categorization," *Proc. Fifth ACM Conf. Digital Libraries*, pp. 195-204, 2000.
- [4] F. Sebastiani, "Machine Learning in Automated Text Categorization," *ACM Computing Surveys*, vol. 34, no. 1, pp. 1-47, 2002.
- [4] M. Vanetti, E. Binaghi, B. Carminati, M. Carullo, and E. Ferrari, "Content-Based Filtering in On-Line Social Networks," *Proc. ECML/PKDD Workshop Privacy and Security Issues in Data Mining and Machine Learning (PSDML '10)*, 2010.
- [5] N.J. Belkin and W.B. Croft, "Information Filtering and Information Retrieval: Two Sides of the Same Coin?" *Comm. ACM*, vol. 35, no. 12, pp. 29-38, 1992.
- [6] P.J. Denning, "Electronic Junk," *Comm. ACM*, vol. 25, no. 3, pp. 163-165, 1982.
- [7] P.W. Foltz and S.T. Dumais, "Personalized Information Delivery: An Analysis of Information Filtering Methods," *Comm. ACM*, vol. 35, no. 12, pp. 51-60, 1992.
- [8] P.S. Jacobs and L.F. Rau, "Scisor: Extracting Information from On- Line News," *Comm. ACM*, vol. 33, no. 11, pp. 88-97, 1990.
- [9] S. Pollock, "A Rule-Based Message Filtering System," *ACM Trans. Office Information Systems*, vol. 6, no. 3, pp. 232-254, 1988.
- [10] P.E. Baclace, "Competitive Agents for Information Filtering," *Comm. ACM*, vol. 35, no. 12, p. 50, 1992.
- [11] P.J. Hayes, P.M. Andersen, I.B. Nirenburg, and L.M. Schmandt, "Tcs: A Shell for Content-Based Text Categorization," *Proc. Sixth IEEE Conf. Artificial Intelligence Applications (CAIA '90)*, pp. 320-326, 1990.
- [12] G. Amati and F. Crestani, "Probabilistic Learning for Selective Dissemination of Information," *Information Processing and Management*, vol. 35, no. 5, pp. 633-654, 1999.
- [13] M.J. Pazzani and D. Billsus, "Learning and Revising User Profiles: The Identification of Interesting Web Sites," *Machine Learning*, vol. 27, no. 3, pp. 313-331, 1997.