

WeChat Toxic Article Detection: A Data-Driven Machine Learning Approach

Yunpeng Weng*, Muhong Wu*, Xu Chen*, Qiong Wu*, Lingnan He†, Liang Chen‡,

* School of Data and Computer Science, Sun Yat-Sen University, Guangzhou, China

† Guangdong Key Laboratory for Big Data Analysis and Simulation of Public Opinion, Guangzhou, China

‡ Tencent, Shenzhen, China

Abstract—Recently, toxic information detection has attracted tremendous amounts of research interest because of the popularity of social networks and the widespread of toxic information which may have dire consequences to the public. Existing work extensively studies toxic article detection in open social networks from information diffusion perspective. However, in closed social networks as exemplified by WeChat Moments (WM), the diffusion process is uneasily visible. To tackle the toxic article detection problem in closed social networks, in this paper we empirically study the articles spread in WM which is based on the largest Chinese social platform WeChat. In particular, we systematically analyze users' behavior and text information of normal and toxic articles and identify a striking difference between them. Furthermore, we design a new model named *MAT-LSTM* which can well capture the impact of different kinds of text information. To improve the performance of automatic toxic article detection, we propose *XMATL* framework which is enhanced from *MAT-LSTM* and can utilize text information and users' behavior characteristics in a holistic manner. We conduct extensive experiments using two real-world datasets and demonstrate that our proposed model can effectively detect toxic articles in WM and achieve outstanding performance gain over the classic methods.

I. INTRODUCTION

With the emergence and proliferation of mobile internet, there are billions of users who post articles, express their opinions on social networking platforms (e.g., Twitter, WeChat, etc.) every day. Social networks have become an indispensable platform for people to communicate and consume information. Unfortunately, it has also proved to be a place where toxic information emerges and festers. Toxic information spread on social networks can be found in various forms, such as rumors, pornography, fraud and share-inducing articles. All these kinds of messages may be diffused with targeted manipulation of public opinion on specific topics and affect society in extremely worrying ways [1]. For example, in crisis situations (e.g., earthquakes, etc.), toxic information can cause wide spread panic and general chaos [2]. Therefore, it is of strategic significance to identify toxic messages and prevent them from spreading.

Towards the above goal, many existing studies mainly focus on open social networks, such as Twitter and Sina Weibo. In

This work was supported in part by the National Science Foundation of China (No. U1711265, 61502315); the Fundamental Research Funds for the Central Universities (No.17lgjc40), and the Program for Guangdong Introducing Innovative and Entrepreneurial Teams (No.2017ZT07X355). Corresponding author: Xu Chen.



Fig. 1. An example of WeChat Moment article. Users' behavior contains reading number, thumbs-up number and report information. The title, name of WeChat official account and content are three types of text information.

these networks, one's page is open for anyone's access (by default) and the path of information diffusion is visible to everyone. Therefore, with the help of a large number of users and official media, it is easy to quickly detect toxic messages and stop them from spreading. While there is another kind of social networks which provides a better protection to users' privacy, such as WeChat Moments (WM). WM is mostly based on the real relationship of daily life, which means that the information in its spread may have a strong impact on users. However, due to its enhanced privacy, we can not get access to the full path of information diffusion and users cannot see other's posted contents if they are not connected as friends, which makes it pretty difficult to disprove and delete toxic articles in WM.

To address this problem, in this paper we first analyze and identify the key differences between normal and toxic articles by the textual analysis. We also explore users' behaviors (e.g., thumbs-up, report in Fig. 1) towards normal and toxic articles and find a great difference between them, which shows that both the text information and users' behavior can contribute to detect the toxic articles in WM. Here we face two challenges. One is that for different kinds of features, we should take advantage of different models specifically. The other is how to devise a unified model that integrates these models in a holistic manner to achieve efficient detection performance.

Through addressing these challenges, we achieve our main contributions of this paper in the following aspects:

- We efficiently exploit the information provided by WM articles and analyze the difference of normal and toxic articles from a multitude of key features, such as users' behavior, the title, name of WeChat official account (WOA) and content of the article.
- We utilize different algorithms to adapt different kinds of features in order to get better prediction performance. We use XGBoost to get better prediction performance for users' behavior. While we propose *MAT-LSTM* framework to capture the predictive ability of text information. Furthermore, we develop *XMATL* model to combine users' behavior and text information together in a holistic manner.
- Extensive experiments are conducted using two realistic datasets, which demonstrate that our proposed model *MAT-LSTM* and *XMATL* can well outperform other methods, e.g., with more than 5% AUC improvement over the commonly-used method-XGBoost.

The rest of our paper is organized as follows. In Section II, we review the related work of toxic article detection in social networks. In Section III, we first give an example of WM articles and then provide a detailed analysis of users' behavior and text information of normal and toxic articles spread in WM. We introduce our proposed model *MAT-LSTM* and *XMATL* in detail in Section IV. We conduct experiments in Section V. Our final conclusion is given in Section VI.

II. RELATED WORK

Toxic information detection has been extensively studied in open social networks [1]. Since toxic information has many types (e.g., rumors, fraud, pornography, etc.), most studies are devoted to rumor detection. Takahashi et al. [2] investigate actual instances of rumors generated after a disaster on Twitter, disclose the characteristic of rumor and then find that it's useful to capture spreading topics and extract clue keywords for rumor detection. Jin et al. [3] use the epidemiological model to analyze the information diffusion pattern and demonstrate how rumor propagates. Liu et al. [4] raise their approach to do real-time rumor debunking on Twitter. They use decision tree method to do rumor detection and find it is much faster than manual verification by professionals. Yang et al. [5] not only use the previously proposed features such as content-based, account-based and propagation-based features to do toxic article detection on Weibo, but also construct some new types of features including client-based and location-based features. Wu et al. [6] focus on the propagation structure of the Weibo message, which traditional feature-based approaches always ignore. Cai et al. [7] make full use of the crowd responses of Weibo messages, which are texts of retweets and comments of the messages. Their experiments show that the crowd responses can significantly contribute to rumor detection. Most of these works focus on proposing and incorporating hand-crafted features and then using machine learning methods to determine whether an article is toxic. Ma

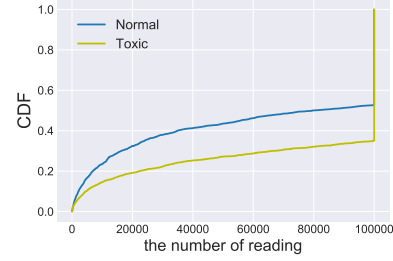


Fig. 2. The Cumulative distribution function (CDF) of reading number for normal and toxic articles. About 50% of normal articles has less than 100,000 readings while over 60% of toxic articles have more than 100,000 readings.

el al. [8] present a novel method based on recurrent neural networks (RNN) which can learn continuous representations of microblog events to detect rumors.

Due to the closed property, only limited pilot research on toxic articles in WM has been carried out. Jiang et al. [9] analyze several WM pages to study the overall diffusion spatial-temporal characteristics and find that the rumor spreading regions in different time experience little change from the very beginning. However, they do not take into account the text information and how to detect toxic articles.

In this paper, we comprehensively analyze the features that contribute to toxic article detection and then propose our methods based on the analysis to detect toxic articles in WeChat Moment.

TABLE I
DESCRIPTION OF DATASET 1.

	Normal	Toxic	Total
All articles	2337	1443	3780
Training articles	1854	1170	3024
Test articles	483	273	756

TABLE II
DESCRIPTION OF DATASET 2.

	Normal	Toxic	Total
All articles	4395	3205	7600
Training articles	3485	2595	6080
Test articles	910	610	1520

III. DATA ANALYSIS

In this section, we first give an example of a WM article and introduce the basic components of it. Then we describe our datasets used in toxic article detection. Finally, we make a preliminary analysis of the difference between normal and toxic articles based on a variety of features.

A. An Example of WM Article

Articles spread in WM are usually first published by WeChat official accounts (WOAs), which can broadcast to their followers. Then the followers can forward these articles to WM, which make it accessible to their friends. Users are able to express their attitudes to a typical article by giving a thumbs-up or reporting it as rumors, pornography, fraud, share-inducing article, etc. As illustrated in Fig. 1, we can get

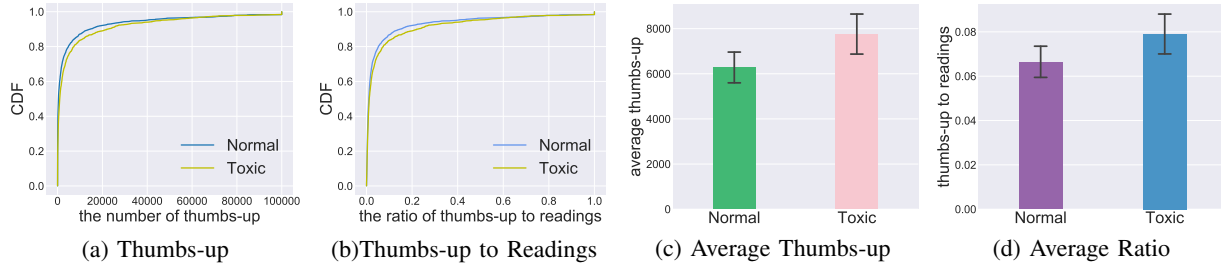


Fig. 3. Readers Thumbs-up Analytics.(a) and (b) are the Cumulative distribution function(CDF) of two types of WeChat articles' thumbs-up and the ratio of thumbs-up to readings respectively. (c) and (d) are Average Thumbs-up number and Average ratio of thumbs-up/reading of two types articles respectively. These four pictures show that toxic articles tend to receive more thumbs-up from users.

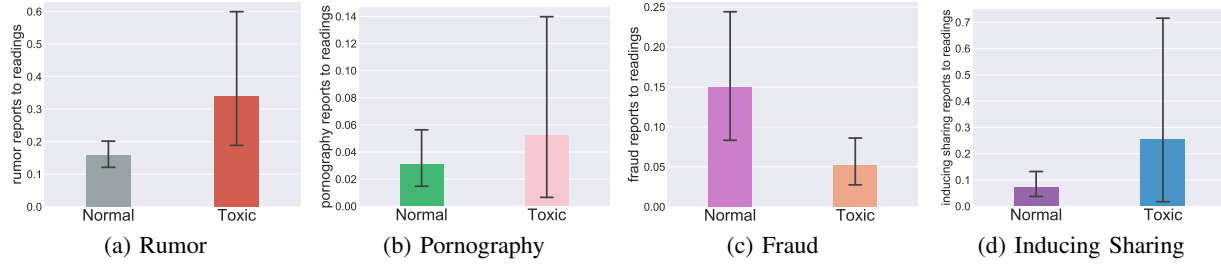


Fig. 4. Reports Analysis. (a), (b), (c) and (d) are respectively average ratios of rumors, pornography, fraud and inducing sharing reports to readings of two types of articles. These pictures show that toxic articles tend to get more reports for rumor, pornography and inducing sharing while normal articles are more likely to be misreported as fraud by users.

users' behavior such as reading number, thumbs-up number and report information as well as text information which includes the title, name of WOA and content of this article. All these characteristics of articles may have an influence on toxic article detection.

B. Dataset Description

Our study in this paper is based on two real-word datasets of articles, which are collected from WeChat, a Chinese popular social platform which has more than 900 million monthly active users. The first dataset includes users behavior characteristics and text information while the second dataset only contains text information. The details of two datasets are described in Table I and Table II. All the samples are from the WeChat article collection that are reported by many users as rumors, pornography, fraud, share-inducing article, etc. A portion of these articles have been identified as toxic articles via some manual screening processes by WeChat and will be deleted later on. Based on these data traces, we aim to devise a data-driven machine learning approach for automatic toxic article detection.

C. Analysis of Users' behavior

Users' behavior towards an article may reflect whether the article is toxic or not to some extent. We mainly focus on the reading number, thumbs-up number and report information of an article.

Reading Number. The reading number indicates the popularity of an article among users. To avoid fabricated data, the maximum reading number is limited to 100,000, which can be seen in Fig. 2 as there is a sharp rise near the upper limit.

Fig. 2 shows the cumulative distribution function of reading number of normal and toxic articles. It is obvious that about 50% of normal articles has less than 100,000 readings while over 60% of toxic articles have more than 100,000 readings. More intuitively, we can see that the curve of normal articles is always above the toxic one, which indicates that toxic articles are more likely to attract users.

Thumbs-up Number. Similar to reading number, the upper limit of thumbs-up number is 100,000. Fig. 3 illustrates users' thumbs-up behavior to normal and toxic articles. Fig. 3 (a) and (c) show the cumulative distribution function of thumbs-up number and the average number of thumbs-up for two types of articles. We can draw the conclusion that toxic articles tend to receive more thumbs-up from users. More precisely, we consider the ratio of thumbs-up to readings and find that we can get the same result as shown in Fig. 3 (b) and (d).

Report. After reading an article, users can choose to report it as rumor, pornography, fraud or a share-inducing article. The average report ratio for the four reasons of normal and toxic articles is shown in Fig. 4. It implies that toxic articles tend to get more reports for rumor, pornography and share-inducing while normal articles are more likely to be misreported as fraud by users. It can be explained that 'fraud' is the first option of reporting reasons.

D. Analysis of Text Information

Text information is one of the most essential factors which can well judge whether an article is toxic or not. As mentioned above, each WM article includes three types of text information: the title, the name of WOA and the content. In this paper, we analyze the influence of text information on normal

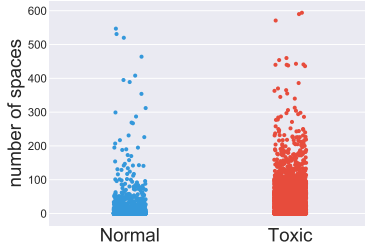


Fig. 5. Distribution of the number of spaces. The distribution of toxic articles is more dense than that of normal articles when the number of spaces exceeds 100.

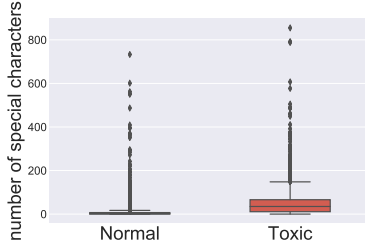


Fig. 6. Distribution of special characters. The upper edge, upper quartile and median line of toxic articles are all above those in normal articles.

and toxic articles from the perspectives of special character statistics and average word vector.

Special Character Statistics. Articles spread in WM are mostly written in Chinese. Different from English, there is no space between Chinese words and usually two or more characters are combined into a word. So we use *Jieba*¹ to do Chinese word segmentation. However, there are some toxic articles that special characters are deliberately added between words to avoid automatic recognition by the detection system. For example, the word “rumor” has two characters in Chinese, and if you add a space between two characters, the system cannot detect them. Therefore, the number of special characters in an article can reflect whether it’s toxic to some extent. Fig. 5 shows the distributions of the numbers of spaces for normal and toxic articles. It is obvious that the distribution of toxic articles is more dense than that of normal articles when the number of spaces exceeds 100.

In addition to spaces, we also choose three other special characters: ‘_’, ‘*’ and ‘.’, which are often used in English, but rarely in Chinese. Fig. 6 illustrates the distribution of these three special characters in the form of boxplot. We can see that the upper edge, upper quartile and median line of toxic articles are all above those in normal articles, which indicates that toxic articles contain more special characters than normal articles.

Average Word Vector. In order to make better use of the implied information in the text and quantify semantic similarities between different articles, we map all the words from three kinds of text information to vectors by using Word2Vec [10], an open source toolkit proposed by Google

¹<https://pypi.org/project/jieba/>

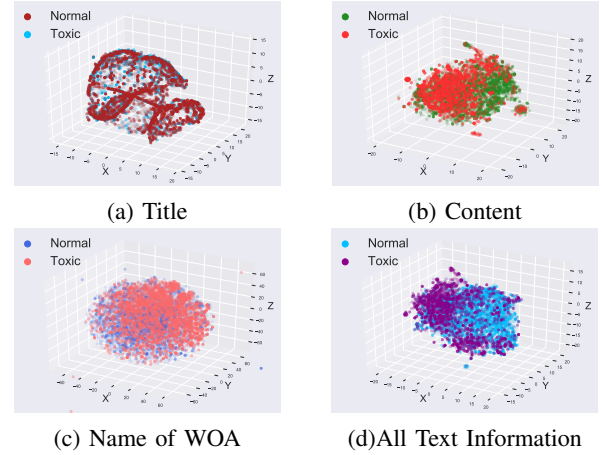


Fig. 7. Distribution of Average Word Vector. (a), (b), (c) and (d) represent the distributions of the average word vectors of title, content, name of WOA and all text information after reduction respectively. These pictures show the difference between the distribution of average word vectors of normal and toxic articles.

in 2013. The model generates a 200-dimensional word vector for each word and then we calculate the average word vector for the title, name of WOA and content of each article. The effectiveness of the learned average word vector may also be investigated qualitatively—and for this purpose we provide a visualization of the t-SNE [11]–transformed word vector representation, which is in the projected 3D space. As shown in Fig. 7, we can see the difference between the distribution of average word vector of normal and toxic articles in the field of title, content, name of WOA and all text information combined.

IV. MODELS

In this section, we first introduce our proposed *MAT-LSTM* model which integrates all the text information together for toxic article detection and then enhance it to *XMATL* framework which integrates users’ behavior and all text information combined in a holistic manner.

A. MAT-LSTM

Attention mechanism has been comprehensively applied in text classification problem as well as LSTM model. Wang etc. [12] propose AT-LSTM (Attention-based LSTM) model for sentiment classification and achieve excellent results. Inspired by their work, we propose *MAT-LSTM* model for toxic articles classification. Firstly, we use Word2Vec to get the vector representation of words in articles and utilize three different attention-based LSTM models to get the dense vector for the title, content and name of each article respectively.

In our MAT-LSTM model, the output of each single AT-LSTM is a dense vector which can be defined as

$$H = \sum_{i=1}^N \alpha_i h_i, \quad (1)$$

where $j \in \{name, title, content\}$ and N is the number of hidden states. The i -th hidden state of LSTM is h_i . For each

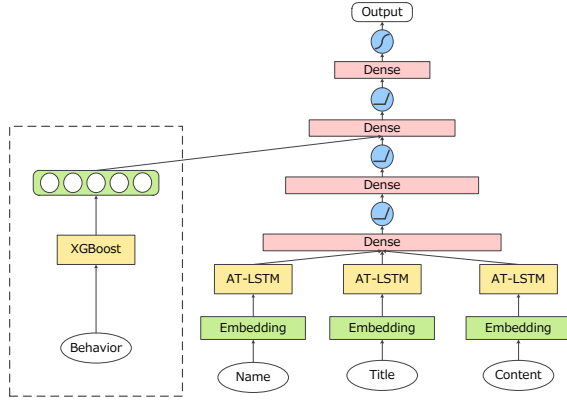


Fig. 8. Structure of MAT-LSTM and XMATL. The XMATL model use XGBoost for users' behavior feature combination as shown in the dotted box and the right part of XMATL is MAT-LSTM. In this state, we perform an attention mechanism and the attention coefficient is α_i which can be expressed as:

$$\alpha_i = \text{softmax}[\tanh(Wh_i + b)], \quad (2)$$

where W is the weight matrix of attention layer. We merge the three dense vectors by concatenating them together as Equation (3) and then fit them into a shallow neural network:

$$H = [H_{\text{name}} || H_{\text{title}} || H_{\text{content}}]. \quad (3)$$

We use joint-training to train the parameters in the model. In this way, the proposed *MAT-LSTM* model can effectively use the information of different types of text information.

B. XMATL

In order to improve the ability of toxic article detection, we enhance *MAT-LSTM* model by adding user's behavior characteristics and then propose *XMATL* framework illustrated in Fig. 8. The Dotted box part is a XGBoost model and the right part in Fig. 8 is MAT-LSTM model. We use XGBoost, a kind of tree ensemble model which has the ability of feature selection and combination, to extract advanced features from user's behavior. For text information, we use the proposed *MAT-LSTM* model. Inspired by the hybrid model proposed by He et al. [13] which use GBDT for feature combination, the leaves of XGBoost model can be treated as the combination of users behavior features and can be combined with the intermediate result of MAT-LSTM. MAT-LSTM is powerful in text-classification while XGBoost is suitable for dealing with user behavior features since the leaves of the decision trees in XGBoost are results of combination of user behavior features. Combining the advantages of XGBoost and MAT-LSTM, the expression ability of models can be enhanced. Eventually, we fit this vector into a shallow neural network and get the final classification result. In this way, we can effectively utilize users' behavior characteristics and text information to achieve efficient toxic article detection task.

V. EXPERIMENTS

In this section, we demonstrate the effectiveness of our proposed *MAT-LSTM* and *XMATL* model on toxic article detection.

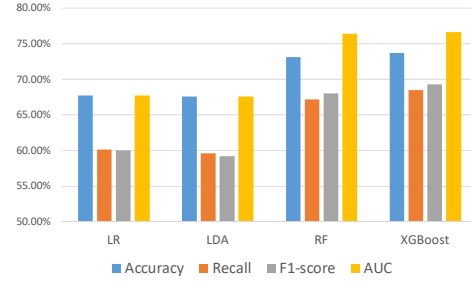


Fig. 9. Classification results based on user behavior.

A. Comparing Methods

- **Logistic Regression (LR)** estimates the parameters of a logistic model where the log-odds of the probability of an event is a linear combination of independent or predictor variables.
- **Linear Discriminant Analysis (LDA)** is a method to find a linear combination of features that characterizes or separates two or more classes of objects or events.
- **Naive Bayes classifier** applies Bayes Theorem with strong independence assumptions between the features.
- **Random Forest (RF)** is an ensemble learning method operating by constructing a multitude of decision trees.
- **XGBoost** is a gradient boosting method [14] improved from Gradient Boosting Decision Trees (GBDT).

As different models are applicable to different types of characteristics, we use different comparing methods for users' behavior and text information.

B. Experimental results

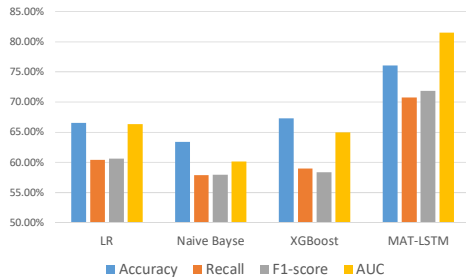
To achieve a better toxic article detection performance, we first train different classifiers based on users' behavior and text information respectively and then integrate them in an effective way.

Users' Behavior only. Based on the analysis in Section III-C, users' behavior (e.g., reading number, thumbs-up number, etc.) reflects the possibility that articles are toxic. To take full advantage of these features, we use XGBoost to train the model and finally we can achieve 73.68% in accuracy for toxic article detection in DataSet-1, which is better than all the comparing methods, as shown in Fig. 9.

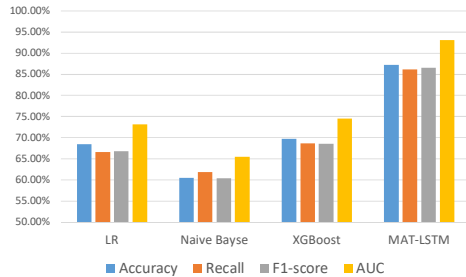
Text Information only. Considering text information only, we compare our proposed MAT-LSTM model with LR, Naive Bayes and XGBoost. We calculate td-idf^2 of the title, name of WOA and content for each article as input features of comparing methods. Furthermore, we conduct our experiments on the two datasets described in Section III-B. We can see in Fig. 10 that in the measurement of AUC, our proposed *MAT-LSTM* model can achieve 81.48% for the first dataset and 93.12% for the second dataset. Comparing with other methods, *MAT-LSTM* is much more powerful in toxic article detection.

²<https://en.wikipedia.org/wiki/Tfidf>

All Features Together. Table III lists the toxic article detection performance of *XMATL* comparing with LR, RF and XGBoost in DataSet-1. We can see that our *XMATL* method can significantly outperforms other machine learning methods in all performance metrics. For example, in terms of AUC, our proposed method which integrates users' behavior and text information achieves a significant performance improvement of 6.14% over the XGBoost algorithm. On the other hand, comparing with the best performance based on only user behavior features or textual features, *XMATL* achieves higher accuracy, which proves that *XMATL* model can effectively combines these two kinds of features and performs better in toxic article detection problem.



(a) Results of DataSet-1



(b) Results of DataSet-2

Fig. 10. Classification results based on textual information:(a) Dataset-1.(b) Dataset-2

TABLE III
THE RESULTS OF CLASSIFICATION(%)

	Accuracy	Recall	F1-score	AUC
LR	66.91	58.43	57.64	69.20
RF	71.82	64.49	64.94	71.62
XGBoost	74.74	68.20	69.18	75.57
XMATL	77.65	73.35	74.35	81.71

VI. CONCLUSIONS

In this paper, we study the toxic article detection in WM, a closed social network. We comprehensively analyze the articles from the perspectives of users' behavior and text information and find that there is a significant difference between normal articles and toxic articles which inspires us to conduct toxic article detection based on these information. To utilize the implied information in text information such as the title, name of WOA and content of an article, we propose *MAT-LSTM* model and prove that it is superior to the comparing

methods by sufficient experiments. Furthermore, we propose *XMATL* framework which is enhanced from *MAT-LSTM* by utilizing text information and users' behavior together. Extensive experiments show that our model can achieve a much better toxic article detection performance.

REFERENCES

- [1] S. Zannettou, M. Sirivianos, J. Blackburn, and N. Kourtellis, "The web of false information: Rumors, fake news, hoaxes, clickbait, and various other shenanigans," *CoRR*, vol. abs/1804.03461, 2018.
- [2] T. Takahashi and N. Igata, "Rumor detection on twitter," in *The 6th International Conference on Soft Computing and Intelligent Systems (SCIS), and The 13th International Symposium on Advanced Intelligence Systems (ISIS), Kobe, Japan, November 20-24, 2012*. IEEE, 2012, pp. 452-457.
- [3] F. Jin, E. R. Dougherty, P. Saraf, Y. Cao, and N. Ramakrishnan, "Epidemiological modeling of news and rumors on twitter," in *Proceedings of the 7th Workshop on Social Network Mining and Analysis, SNAKDD 2013, Chicago, IL, USA, August 11, 2013*, F. Zhu, Q. He, R. Yan, and J. Yen, Eds. ACM, 2013, pp. 8:1-8:9.
- [4] X. Liu, A. Nourbakhsh, Q. Li, R. Fang, and S. Shah, "Real-time rumor debunking on twitter," in *Proceedings of the 24th ACM International Conference on Information and Knowledge Management, CIKM 2015, Melbourne, VIC, Australia, October 19 - 23, 2015*, J. Bailey, A. Moffat, C. C. Aggarwal, M. de Rijke, R. Kumar, V. Murdock, T. K. Sellis, and J. X. Yu, Eds. ACM, 2015, pp. 1867-1870.
- [5] F. Yang, Y. Liu, X. Yu, and M. Yang, "Automatic detection of rumor on sina weibo," in *Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics*. ACM, 2012, p. 13.
- [6] K. Wu, S. Yang, and K. Q. Zhu, "False rumors detection on sina weibo by propagation structures," in *31st IEEE International Conference on Data Engineering, ICDE 2015, Seoul, South Korea, April 13-17, 2015*, J. Gehrke, W. Lehner, K. Shim, S. K. Cha, and G. M. Lohman, Eds. IEEE Computer Society, 2015, pp. 651-662.
- [7] G. Cai, H. Wu, and R. Lv, "Rumors detection in chinese via crowd responses," in *2014 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2014, Beijing, China, August 17-20, 2014*, X. Wu, M. Ester, and G. Xu, Eds. IEEE Computer Society, 2014, pp. 912-917.
- [8] J. Ma, W. Gao, P. Mitra, S. Kwon, B. J. Jansen, K. Wong, and M. Cha, "Detecting rumors from microblogs with recurrent neural networks," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, S. Kambhampati, Ed. IJCAI/AAAI Press, 2016, pp. 3818-3824.
- [9] W. Jiang, B. Chen, L. He, Y. Bai, and X. Qiu, "Features of rumor spreading on wechat moments," in *Web Technologies and Applications - APWeb 2016 Workshops, WDMA, GAP, and SDMA, Suzhou, China, September 23-25, 2016, Proceedings, ser. Lecture Notes in Computer Science*, A. Morishima, L. Chang, T. Z. J. Fu, K. Liu, X. Yang, J. Zhu, R. Zhang, W. Zhang, and Z. Zhang, Eds., vol. 9865, 2016, pp. 217-227.
- [10] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *CoRR*, vol. abs/1301.3781, 2013.
- [11] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579-2605, 2008.
- [12] Y. Wang, M. Huang, X. Zhu, and L. Zhao, "Attention-based LSTM for aspect-level sentiment classification," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, J. Su, X. Carreras, and K. Duh, Eds. The Association for Computational Linguistics, 2016, pp. 606-615.
- [13] X. He, J. Pan, O. Jin, T. Xu, B. Liu, T. Xu, Y. Shi, A. Atallah, R. Herbrich, S. Bowers, and J. Q. Candela, "Practical lessons from predicting clicks on ads at facebook," in *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising, ADKDD 2014, August 24, 2014, New York City, New York, USA*, E. Saka, D. Shen, K. Lee, and Y. Li, Eds. ACM, 2014, pp. 5:1-5:9.
- [14] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, B. Krishnapuram, M. Shah, A. J. Smola, C. C. Aggarwal, D. Shen, and R. Rastogi, Eds. ACM, 2016, pp. 785-794.