

An Intelligent Online Grooming Detection System Using AI Technologies

Philip Anderson, Zheming Zuo, Longzhi Yang
Department of Computer and Information Sciences,
Faculty of Engineering and Environment
Northumbria University, United Kingdom
Email: longzhi.yang@northumbria.ac.uk

Yanpeng Qu
Information Science and Technology College
Dalian Maritime University
Dalian, P. R. China
Email: yanpengqu@dlnu.edu.cn

Abstract—The rapid expansion of the Internet has experienced a significant increase in cases of child abuse, as more and more young children have greater access to the Internet. In particular, adults and minors are able to exchange sexually explicit messages and media via a variety of online platforms that are widely available, which leads to an increasing concern of child grooming. Traditionally, the identification of child grooming relies on the analysis and localisation of conversation texts, but this is usually time-consuming and associated with other implications such as psychological pressure on the investigators. Therefore, automatic methods to detect grooming conversations have attracted the attention of many researchers. This paper proposes such a system to identify child grooming in online chat conversations, where the training data of the system were harvested from publicly available information. The data processing is based on a group of AI technologies, including fuzzy-rough feature selection and fuzzy twin support vector machine. Evaluation shows the promise of the proposed approach in identifying online grooming conversations to be implemented in the future after further development to support real-world cases.

I. INTRODUCTION

Digital forensics is widely used by security sectors for the identification, acquisition and analysis of evidence originating from digital devices to prevent, investigate or endorse the conviction of criminal events [1]. Indeed, digital evidence has now become a prominent aspect of many criminal investigations. Typically, the examination of digital data uses a variety of forensically sound methods and techniques to identify, capture and analyse data. In an attempt to recover relevant evidence, these could include employing forensics software to search for the files and data using specific keywords, keyword lists or file types. The analysis of such evidence is the process of interpreting the extracted data to determine its significance in the investigation.

Due to the time-consuming nature of the conventional digital forensics investigative process, it can be susceptible to errors occurring. In addition, the human aspect and manual examination process of investigations can potentially lead to examiner bias, either unintentional or intentional. The potential for investigator exhaustion and burnout from the intense repetitive approach adds further to the potential for errors to occur during investigations. Such processes, in turn, impact the workloads of investigators arguably subjecting them to increased stress as workloads increase whilst working to given time-frames [2]. What is more, investigators working on child

exploitation cases are at an increased risk for experiencing psychological distress [3].

The investigation of online child grooming, as a common task of digital forensics, is no exception. Online child grooming is a process of approaching, persuading, and engaging a child in sexual activity. The perpetrator approaches the child to initially build an emotional relationship before it becomes sexual [4]. Guidance and research show that groomers exploit any vulnerability to increase the child's dependence on them as part of the grooming process [5]. It is therefore essential to understand sexual grooming and identify when it is happening from both a psychological or social perspective of the child to prevent sexual abuse occurring [6]. This usually relies on careful analysis of the location and content of conversation texts, which often requires a significant amount of time [7].

This paper presents an intelligent online grooming detection system using a range of AI technologies, so that human errors can be prevented, time can be saved, and other implications can also be addressed. In particular, the proposed system can automatically identify online child grooming based on the bag of words (BoW) approach [8], which identifies a list of words for text classification in the context of grooming. This is followed by the application of fuzzy-rough feature selection in choosing the most important features as the input of classifiers. A chosen set of classifiers, such as the fuzzy twin support vector machine, are then applied based on the extracted features for grooming text classification.

The training data were harvested from publicly available information from two sources. One source was the archive page of the perverted-justice website, where more than 600 archived texts of actual child grooming conversations involving perpetrators and adult volunteers acting as children were available [9]. The other source was conversations taken from the PAN13 data set which was originally proposed for predicting age and gender [10]. The combination of the training data sets from different sources provides a more balanced and diverse means to empower the proposed system with enhanced performance, as witnessed by the evaluation results.

The remainder of this paper is organised as follows. Section II reviews the related work, and Section III describes the approach followed in this research. The experimental results are analysed in Section IV while Section V presents conclusion and suggestions for future work.

II. BACKGROUND

Conventional grooming detection approaches are briefly reviewed first in this section, which is followed by some representative AI-based grooming detection methods.

A. Conventional Approaches

Keyword searches are commonly used by investigators in various digital forensic tools for detecting grooming conversations. The tools are designed to find all instances of a given text string or strings, using single keywords or a list of words which contains multiple keywords around a single area of interest, such as grooming. However, not all data can be instantly and easily searched, often it must be decoded and presented in a text-searchable form or extracted from areas of unused data or deleted data [11]. The results of these keyword searches are often presented by the forensic tool in tables. To improve keyword searching, forensic software tools started to index the data, identifying the location of each keyword within the data source. This significantly improves the speed of the overall search process and enables the construction of complex search criteria and receive immediate results [12].

However, there are issues to using the traditional approach of keyword searching, as it often returns a large amount of irrelevant information which in turn wastes the time of the investigators as they review the findings. Furthermore, keyword searches in digital forensic tools often do not offer grouping or ordering functionality for the search results to help investigators retrieve the relevant information in an efficient manner. In addition, the determination of key works for searching are often subjective, which introduces noise in the process.

B. AI-based Approaches

The work [13] investigated the effectiveness of text classifiers to identify child grooming in online chat conversations. This was achieved by applying a new technique alongside three traditional text categorisation techniques such as linear or logistic regression, decision trees, and Naive Bayes. For improving the classification performance, psychometric and categorical information techniques (e.g. linguistic inquiry and word count) were employed. For evaluation, the chat logs were collected from various publicly accessible websites. This was followed by data pre-processing in which data cleaning (by removing user names) and format conversion (to generate string vectors) were conducted. The two types of features were selected for two sets of experiments. In one set of experiment the term-based features were used. The other set of experiment used psychometric and categorical information from LIWC. In the drawn conclusions, it has proven that the performance of Naive Bayes and (linear or logistic) regression classifiers in predicting the type of the chats which could also be enriched by involving psychometric and categorical information.

A low computational cost classification method based on the number of existing grooming conversation characteristics was presented in [14]. Two types of conversations were used: child grooming and non-grooming conversations but with grooming characteristics included. Texts were pre-processed and transformed into a vector space model. The features are words or combinations of words that formed the word list. The method was evaluated using 150 conversation texts

in which 105 texts were grooming and 45 ones were non-grooming. For extracting text features, 17 attributes were extracted for each data instance. According to the conclusion, classifiers such as support vector machine (SVM) and k-nearest neighbouring (KNN) were suggested for yielding reasonable text classification precision.

The work of [15] reported a grooming detection approach to address uncertainty based on a highly imbalanced real-world data set adopted from a publicly available PAN'13 profiling data set. The work reconstructed a new data set by randomly selecting 1,000 from the total 262,254 XML documents. Bow and TFIDF features were adopted in identifying the word-wise significance in each document. In eliminating redundancy and noise in the extracted features, fuzzy-rough feature selection (FRFS) approach was employed for dimensionality reduction. Based on the intensively conducted experiments, the conclusion given by the work was that performance with the implementation of feature selection was generally better than those without using feature selection.

III. GROOMING DETECTION FRAMEWORK

The proposed grooming detection framework is illustrated in Figure 1. Briefly, this framework consists of six consecutive phases, which are, data collection from various sources, pre-processing to guarantee all the documents are with unified format, feature extraction to build concise representation for each document, feature selection to remove redundant or noisy features in the generated document representation, normalisation to make the data more sensitive to classifiers, and finally classification for predicting the class label of a given online conversation.

A. Data Harvesting

In this work, the data was collected from two sources: child grooming (Internet source) [9] and non-grooming (dataset source) [10] conversations. The former contains more than 600 archived texts (*i.e.* .TXT files) of child grooming conversations involving perpetrators and adult volunteers acting as children, from which 200 texts ordered by the chat usernames of perpetrators were selected. These conversations vary in size from 1KB to 495KB giving suitable content variety. The later source consisted of a total 262,254 XML files in which, 236,814 and 25,440 XML files were contained in the training and testing sets of the English corpus, from which 1,000 XML files were randomly selected. Thus, the data set used in this work contains 1,200 documents in total.

B. Data Pre-processing

The harvested data are in different formats, as 1,000 documents are in the form of XML and 200 were TXT files. In this project, the 200 TXT files were converted to XML files that share the same format as the 1,000 XML files by wrapping the TXT files as follows:

```
<author uuid="Bpm020" lang="en">
<conversation id = "11">
... .. //the TXT file
</conversation>
</author>
```

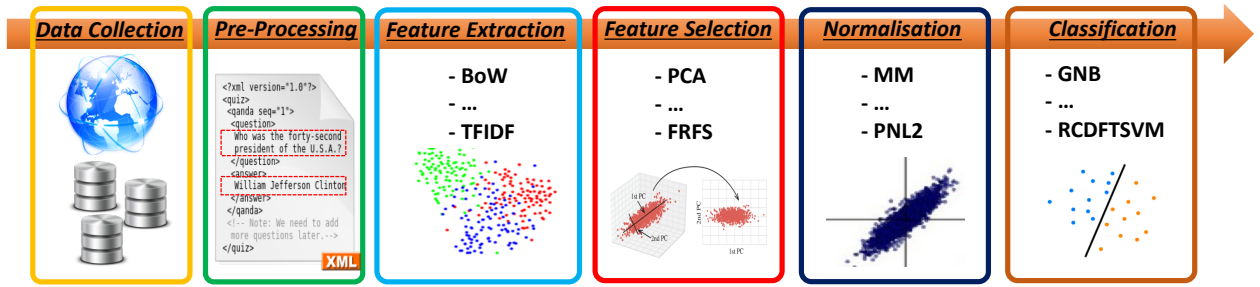


Fig. 1. The framework of online grooming text classification

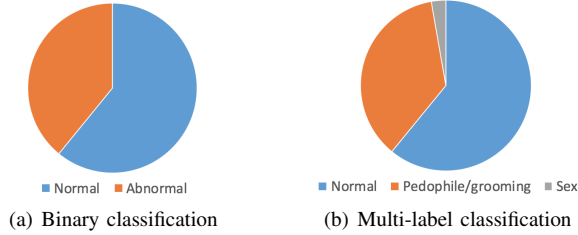


Fig. 2. Data instances percentage for both learning tasks

where `uuid` refers to the original TXT file name and all the conversations are included in the node of `conversation`.

It is essential to convert any online conversation logs into the format of the training data set. The formats of chat logs vary from one social network to another. Skype and WhatsApp, save conversations to a generic database file (.db) and there are several software programs available that can view the raw database file whereas Facebook Messenger conversations are downloaded in an HTML file format. The proposed system takes XML as its input, due to its popularity. There are several existing tools which can be directly used for file format conversation. For instance, there are online resources and software tools that convert different flat files to XML, such as Stylus Studio [16].

C. Feature Extraction and Selection

Based on the training dataset, a set of text features are extracted. For building discriminative document representation and considering the diversity of chat contents, bag of words (BoW) [8] and term frequency–inverse document frequency (TF-IDF) [17] feature extraction approaches were adopted. Redundant and noisy information are commonly included in the generated document features, thus feature selection is required. Briefly, feature selection selects a small proportion of attributes by mapping high dimensional data into the lower one based on the training data set. A number of feature selection techniques are available, such as principal component analysis (PCA), linear discriminant analysis (LDA), local linear embedding (LLE), locality preserving projection (LPP), stochastic proximity embedding (SPE), stochastic neighbour embedding (SNE), and fuzzy-rough feature selection (FRFS) approaches ([18], [19], [20], [21], [22]). Due to its high performance and efficiency, fuzzy-rough feature selection is used to develop the proposed online grooming detection system. The optimal number of features used in the proposed system is 150 for binary classification and 30 for multi-label classification.

For any online grooming detection request regarding a given pre-processed chat log, the system extracts important features based on the developed model. In this case, each chat log is concisely represented by its feature vector.

D. Feature Normalisation

The value ranges of the selected features may vary significantly, as they have very different physical meanings. This results in deteriorating the performance of classification in the subsequent stage of the framework. To cope with this challenging factor, data is normalised first. Denoting the selected features for the dataset as X , then the normalised features \hat{X} can be calculated by using different feature normalisation techniques as summarised in Table I. Of course, the users of the proposed system do not need to use all of these normalisation approaches for one application, but the proposed system provides the flexibility to users when they have customisation preferences.

TABLE I. COMMON NORMALISATION APPROACHES

No.	Normalisation Technique	Formulation
1	Min-Max (MM)	$\hat{X} = X - \min(X) / \max(X) - \min(X)$
2	ℓ_1 norm.	$\hat{X} = X / \ X\ _1$
3	ℓ_2 norm.	$\hat{X} = X / \ X\ _2$
4	Power Norm. (PN)	$\hat{X} = \text{sign}(X) X ^\alpha, \alpha \in [0, 1]$
5	ℓ_1 PN	$\hat{X} = \text{sign}(X / \ X\ _1) X / \ X\ _1 ^\alpha$
6	ℓ_2 PN	$\hat{X} = \text{sign}(X / \ X\ _2) X / \ X\ _2 ^\alpha$
7	PN ℓ_1	$\hat{X} = \ \text{sign}(X) X ^\alpha\ _1$
8	PN ℓ_2	$\hat{X} = \ \text{sign}(X) X ^\alpha\ _2$

E. Classification

A number of classifiers have been implemented in the proposed system as listed in Table II. Note that the system is developed based on the initial work of [15]. In the initial work, the first four classifiers in the table have been detailed, and thus these are omitted here. The last two classifiers are variants of the conventional SVM. In particular, both linear and nonlinear (*i.e.* using RBF kernel) coordinate descent fuzzy twin support vector machine (CDFTSVM) has been employed in the proposed system to enhance the classification performance, by removing the noise using a fuzzy membership function and reducing the computational complexity using a coordinate descent strategy [23], [24].

Different to the conventional SVM (where one decision plane is used, *i.e.* $y = ax + b$) for binary class labels ‘+’ and ‘-’, twin SVM (TSVM) aims to generate two optimal solutions, *i.e.* (a_+^*, b_+^*) and (a_-^*, b_-^*) , of two non-parallel decision

TABLE II. COMMON CLASSIFIERS FOR ONLINE GROOMING DETECTION

No.	Classifier	Abbreviation
1	Gaussian Naïve Bayes	GNB
2	Random Forest	RF
3	AdaBoost	AB
4	Logistic Regression	LR
5	Linear Coordinate Descent Fuzzy Twin SVM [24]	LCDFTSVM
6	RBF Coordinate Descent Fuzzy Twin SVM [24]	RCDFTSVM

planes (*i.e.* $a_+^T x + b_+ = 0$ and $a_-^T x + b_- = 0$) for predicting the label y of the instance x :

$$y = \arg \min_{\pm} \frac{|a_{\pm}^{*T} x + b_{\pm}^*|}{\|a_{\pm}^*\|}. \quad (1)$$

On this basis, by assigning a fuzzy membership to each training sample in each of the two classes, fuzzy TSVM (FTSVM) is established for robust learning [24]. For instance, in the nonlinear FTSVM, two decision planes become $\kappa(x, \hat{X}^T) a_+ + b_+ = 0$ and $\kappa(x, \hat{X}^T) a_- + b_- = 0$, where $\kappa(\cdot)$ is a kernel function, then the class label y of data instance x can be predicted by:

$$y = \arg \min_{\pm} \frac{|\kappa(x, \hat{X}^T) w_{\pm}^{*T} + b_{\pm}^*|}{\sqrt{a_{\pm}^{*T} \kappa(\hat{X}, \hat{X}^T) a_{\pm}^*}}. \quad (2)$$

Then, for speeding up the FTSVM, CDFTSVM [24] is constructed by using coordinate descent strategy with active set shrinking [25].

The proposed system includes all of the 6 classifiers to offer flexibility to users. All of these classifiers have been trained using the harvested data set as introduced in Section III-B. The detailed configuration of these classifiers is detailed in the next section.

IV. EVALUATION AND DISCUSSION

The proposed system was evaluated in this section. The detailed system setup and configuration of the applied AI techniques are presented.

A. System and Evaluation Setup

The proposed system can be used with major existing operating systems such as Windows, Ubuntu, Mac OS X, as it only requires the installation of Python 2.7.14, except the CDFTSVM [24] where Matlab 2018b is used. A piece of live future work is to implement the proposed CDFTSVM in Python to make an ‘all-in-one’ software solution for user convenience. For hardware, an ordinary Notebook is sufficient to run the proposed system. All of the evaluation experiments were conducted using a MacBook Pro with Intel(R) Core(R) i7 processor (3 GHz) and 16 GB RAM.

Given the limited size of the training data set and also the time-consuming nature in capturing more labelled data, 10-Fold cross-validation is adopted here for all the experiments [15]. For feature extraction, the approach of BoW and TFIDF were selected. For feature selection, FRFS was adopted (60% number of attributes were selected). For feature normalisation,

the techniques listed in Table I were employed. For classification, all of the classifiers adopted in the proposed system were summarised in Table II.

B. Results and Discussion

To systematically evaluate the proposed system, extensive experiments were conducted, each with a different number of extracted features, text feature normalisation techniques, and the classification approaches. The detailed experimental results for binary classification are presented in Figure 3; and the experiments for multi-label classification are summarised in Table III without the use of FRFS, and Table IV with FRFS.

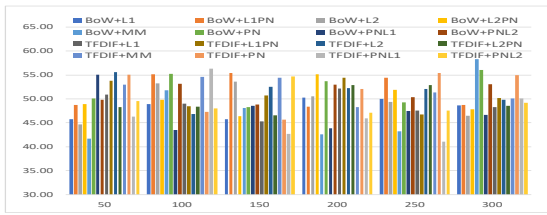
As delineated in Figure 3, when FRFS was not utilised, GNB achieved its best performance of 58.33% using the BoW features and MM normalisation technique, RF yielded 57.50% by the combination of BoW and PN ℓ 1, 60.75% mean accuracy was generated by AB in using TFIDF plus PN ℓ 1, while LR reached 61.92% by consecutively applying BoW and MM, 60.93% and 60.79% were separately generated by LCDFTSVM and RCDFTSVM, where the same combination of TFIDF and PN ℓ 1 was employed.

With the involvement of FRFS for the binary classification task, the performance of GNB was boosted up to 59.08% using BoW and ℓ 1, RF was also improved to 57.66% with the help of BoW and ℓ 1PN, AB arrived at the best accuracy of 60.34% by BoW and MM, LR yielded 61.08% with TFIDF and PN ℓ 1 employed, while 60.90% and 60.74% were generated by LCDFTSVM (BoW+ ℓ 1) and RCDFTSVM (TFIDF+PN ℓ 2).

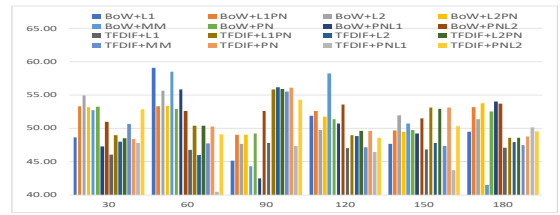
For the best performances obtained in classifying the multi-label instances, without using the FRFS, GNB obtained its peak performance of 47.84% (using TFIDF features that normalised by PN), RF achieved 56.85% (by classifying BoW features which further normalised using ℓ 2PN), AB yielded 60.34% when BoW features and PN ℓ 1 technique were utilised, LR produced 61.58% prediction precision using the combination of BoW and PN ℓ 2, while LCDFTSVM (TFIDF+PN ℓ 1) and RCDFTSVM (BoW+PN ℓ 2) were respectively with 60.96% and 60.78%.

After using the FRFS approach, the performance of GNB increased to 56.84% using the combination of BoW and MM, while RF reached its best performance of 55.99% on the BoW features that normalised using MM, AB yielded the mean accuracy of 57.34% when BoW features and PN ℓ 1 technique were used, LR yielded 61.42% (TFIDF+ ℓ 1), LCDFTSVM and RCDFTSVM achieved 60.96% (BoW and PN ℓ 1) and 60.68% (BoW+MM).

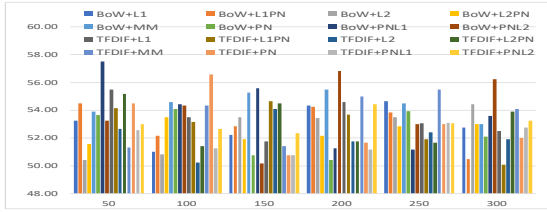
Based on the above discussion, it can be concluded that it is possible for the AI-based digital forensics tools to be used for chat log screening, but it needs further improvement before it can be applied in the real-world to support investigators. For instance, the evolving of grooming chat logs along with the development of Internet language requires the application of adaptive AI approaches [26], which remains as a piece of active future work. From a technical point of view, most classifiers could achieve their best performance using PN normalisation technique or its variants (*i.e.* PN ℓ 1, PN ℓ 2, ℓ 2PN) when FRFS is not used, while conventional normalisation techniques (*i.e.* MM and ℓ 1) tend to be better in improving the accuracy when FRFS is adopted.



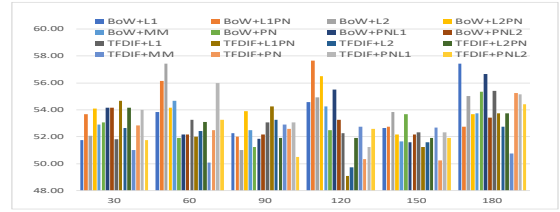
(a) GNB without FRFS



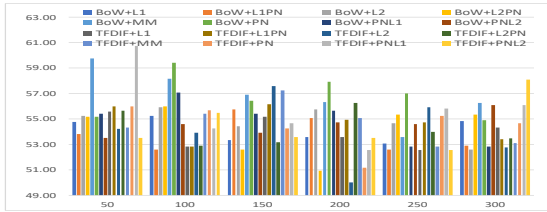
(b) GNB with FRFS



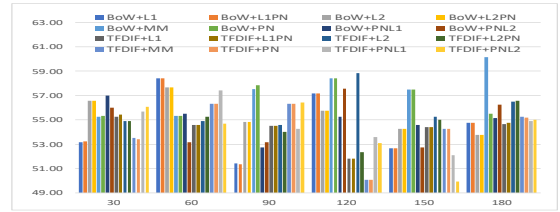
(c) RF without FRFS



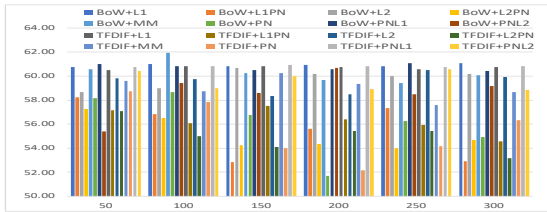
(d) RF with FRFS



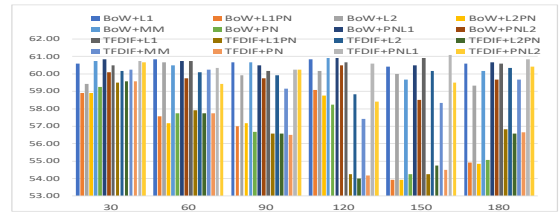
(e) AB without FRFS



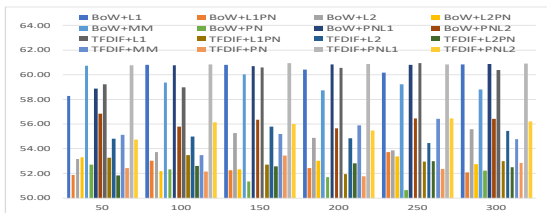
(f) AB with FRFS



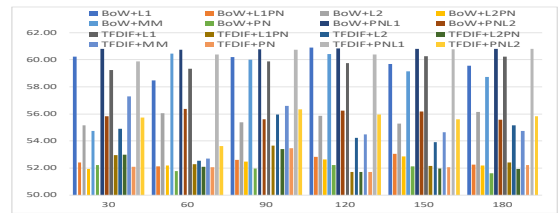
(g) LR without FRFS



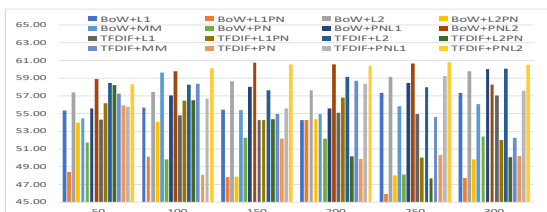
(h) LR with FRFS



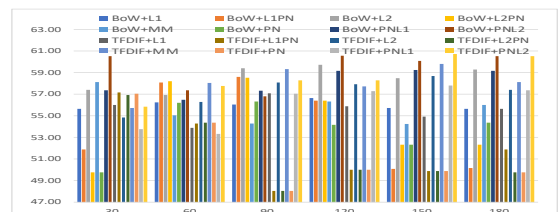
(i) LCDTSVM without FRFS



(j) LCDTSVM with FRFS



(k) RCDTSVM without FRFS



(l) RCDTSVM with FRFS

Fig. 3. Performance in mean accuracy with and without fuzzy-rough feature selection applied under different feature normalization strategies and feature set sizes for binary classification.

TABLE III. PERFORMANCE SUMMARY WITHOUT USING FRFS WITH BEST PERFORMANCE MARKED IN **BOLD**

Classification Type	Classifier	#feature extracted					
		50	100	150	200	250	300
Binary	GNB	55.58	56.34	55.42	55.16	55.41	58.33
	RF	57.50	56.58	55.58	56.84	55.50	56.25
	AB	60.75	59.41	57.58	57.92	57.00	58.08
	LR	61.00	61.92	60.92	60.92	61.08	61.08
	LCDFTSVM	60.77	60.82	60.93	60.93	60.86	60.90
	RCDFTSVM	58.90	60.09	60.75	60.56	60.79	60.50
Multi-label	GNB	47.84	46.58	43.34	36.76	40.92	45.59
	RF	55.01	54.17	56.33	55.74	58.85	55.76
	AB	60.34	56.77	56.33	57.40	57.17	56.90
	LR	61.09	61.58	60.84	61.09	61.09	61.17
	LCDFTSVM	60.82	60.94	60.89	60.88	60.96	60.87
	RCDFTSVM	59.42	59.81	60.57	60.69	60.71	60.78

TABLE IV. PERFORMANCE SUMMARY WITH FRFS ADOPTED WITH BEST PERFORMANCE MARKED IN **BOLD**

Classification Type	Classifier	#feature extracted					
		30	60	90	120	150	180
Binary	GNB	54.92	59.08	56.16	58.25	53.09	54.00
	RF	54.67	57.41	54.25	57.66	53.84	57.42
	AB	57.00	58.41	57.84	58.83	57.50	60.17
	LR	60.83	60.83	60.67	60.92	61.08	60.83
	LCDFTSVM	60.81	60.75	60.77	60.90	60.81	60.81
	RCDFTSVM	60.51	58.91	59.42	60.58	60.74	60.53
Multi-label	GNB	45.23	56.84	41.75	42.25	48.10	46.32
	RF	55.99	54.68	55.01	55.83	54.68	55.08
	AB	57.34	55.76	55.81	55.92	56.67	54.58
	LR	60.75	61.42	61.17	61.42	60.92	61.25
	LCDFTSVM	60.96	60.87	60.94	60.91	60.94	60.96
	RCDFTSVM	60.68	59.87	60.34	60.59	60.56	60.68

V. CONCLUSION

The automatic system to detect child grooming is expected to play an increasingly important role in analysing the text from online conversations. The proposed system showed that, given text chats between a sexual predator and a pseudo victim, it is possible to automatically distinguish between child grooming conversations and non-grooming conversations, as implicated by the evaluation results. Although promising, the system should be further improved by integrating and adapting more AI techniques and better exploring their parameter settings. Also, it is expected that a larger data set with better coverage can lead to performance enhancement. What is more, the proposed system should also be extended to deal with live evolving streaming chat conversations along with the evolution of Internet language.

REFERENCES

- [1] X. Du, N. Le-Khac, and M. Scanlon, "Evaluation of digital forensic process models with respect to digital forensics as a service," in *Proceedings of the 16th European Conference on Cyber Warfare and Security*, 2017.
- [2] G. Horsman, C. Laing, and P. Vickers, "A method for reducing the risk of errors in digital forensic investigations," in *Communications and Multimedia Security*, vol. 7394. Berlin, Heidelberg: Springer, 2012.
- [3] K. Seigfried-Spellar, "Assessing the psychological well-being and coping mechanisms of law enforcement investigators vs. digital forensic examiners of child pornography investigations," *Journal of Police and Criminal Psychology*, vol. 33, pp. 1–12, 2017.
- [4] Childline, "Online grooming," 2018-12-19. [Online]. Available: <https://www.childline.org.uk/info-advice/bullying-abuse-safety/online-mobile-safety/online-grooming/>.
- [5] NSPCC, "Grooming - what it is, signs and how to protect children," 2018-12-19. [Online]. Available: <https://www.nspcc.org.uk/preventing-abuse/child-abuse-and-neglect/grooming/>.
- [6] D. Pollack and A. MacIver, "Understanding sexual grooming in child abuse cases," *ABA Child Law Practice*, vol. 34, pp. 165–168, 11 2015.
- [7] A. Kontostathis, L. Edwards, and A. Leatherman, "Text mining and cybercrime," in *Text Mining: Application and Theory*. John Wiley & Sons, Ltd, 2010.
- [8] D. Jurafsky and J. H. Martin, *Speech and language processing*. Pearson London:, 2014, vol. 3.
- [9] Perverted Justice, "The largest and best anti-predator organization online," 2018-12-12. [Online]. Available: <http://www.perverted-justice.com/>.
- [10] F. Rangel, P. Rosso, M. Koppel, E. Stamatatos, and G. Inches, "Overview of the Author Profiling Task at PAN 2013," in *CLEF 2013 Evaluation Labs and Workshop – Working Notes Papers, 23-26 September, Valencia, Spain*, 2013.
- [11] A. Bogen, "Selecting keyword search terms in computer forensics examinations using domain analysis and modeling," Ph.D. dissertation, 2006, aAI3241379.
- [12] N. Beebe and J. Clark, "Digital forensic text string searching: Improving information retrieval effectiveness by thematically clustering search results," *Digital Investigation*, vol. 4, pp. 49–54, 2007.
- [13] M. D. Rahman Miah, J. Yearwood, and S. Kulkarni, "Detection of child exploiting chats from a mixed chat data set as a text classification task," in *Australian Language Technology Association Workshop*, 2011, pp. 157–165.
- [14] F. E. Gunawan, L. Ashianti, S. Candra, and B. Soewito, "Detecting online child grooming conversation," in *11th International Conference on Knowledge, Information and Creativity Support Systems*, 2016, pp. 1–6.
- [15] Z. Zuo, J. Li, L. Yang, P. Anderson, and N. Naik, "Grooming detection using fuzzy-rough feature selection and text classification," in *In IEEE World Congress on Computational Intelligence*. IEEE, 2018, pp. 1–8.
- [16] StylusStudio, "Convert text files to xml (any flat file to xml)," 2019-01-23. [Online]. Available: <http://www.stylusstudio.com/text-file-to-xml.html>.
- [17] H. C. Wu, R. W. P. Luk, K. F. Wong, and K. L. Kwok, "Interpreting tf-idf term weights as making relevance decisions," *ACM Transactions on Information Systems (TOIS)*, vol. 26, no. 3, p. 13, 2008.
- [18] L. Van Der Maaten, E. Postma, and J. Van den Herik, "Dimensionality reduction: a comparative," *J Mach Learn Res*, vol. 10, pp. 66–71, 2009.
- [19] Y. Qu, Y. Rong, A. Deng, and L. Yang, "Associated multi-label fuzzy-rough feature selection," in *Joint World Congress of International Fuzzy Systems Association and International Conference on Soft Computing and Intelligent Systems (IFSAS-SCIS)*, 2017.
- [20] R. Cameron, Z. Zuo, G. Sexton, and L. Yang, "A fall detection/recognition system and an empirical study of gradient-based feature extraction approaches," in *UK Workshop on Computational Intelligence (UKCI)*, 2017.
- [21] Q. Sun, Y. Qu, A. Deng, and L. Yang, "Fuzzy-rough feature selection based on λ -partition differentiation entropy," in *2017 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD)*, July 2017, pp. 1222–1227.
- [22] Z. Zuo, L. Yang, Y. Peng, F. Chao, and Y. Qu, "Gaze-informed egocentric action recognition for memory aid systems," *IEEE Access*, vol. 6, pp. 12 894–12 904, 2018.
- [23] R. Khemchandani, S. Chandra *et al.*, "Twin support vector machines for pattern classification," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 29, no. 5, pp. 905–910, 2007.
- [24] B. Gao, J. Wang, Y. Wang, and C. Yang, "Coordinate descent fuzzy twin support vector machine for classification," in *IEEE 14th International Conference on Machine Learning and Applications*, 2015, pp. 7–12.
- [25] Y.-H. Shao and N.-Y. Deng, "A coordinate descent margin based-twin support vector machine for classification," *Neural networks*, vol. 25, pp. 114–121, 2012.
- [26] L. Yang, F. Chao, and Q. Shen, "Generalized adaptive fuzzy rule interpolation," *IEEE Transactions on Fuzzy Systems*, vol. 25, no. 4, pp. 839–853, 2017.