

Affect Analysis of Radical Contents on Web Forums Using SentiWordNet

Tawunrat Chalothorn and Jeremy Ellman

Abstract—The internet has become a major tool for communication, training, fundraising, media operations, and recruitment, and these processes often use web forums. This paper presents a model that was built using SentiWordNet, WordNet and NLTK to analyze selected web forums that included radical content. SentiWordNet is a lexical resource for supporting opinion mining by assigning a positivity score and a negativity score to each WordNet. The approaches of the model measure and identify sentiment polarity and affect the intensity of that which appears in the web forum. The results show that SentiWordNet can be used for analyzing sentences that appear in web forums.

Index Terms—SentiWordNet, sentiment, analysis, web forums, radical.

I. INTRODUCTION

Web forums have become important places for social communication and discussion on the internet. Some radical groups also use them for communication and disseminating their ideologies to the public [1]. These kinds of forums can be referred to as part of the Dark Web. The Dark Web includes websites that are used by terrorists, radicals and extremist groups [2]. This paper presents the system approach of two web forums in the area of sentiment and affects analysis. Their content is related to radicalization. Many people have questioned why this research was carried out. The reason is that the United Kingdom's parliament has enacted an anti-terrorism law, the Terrorism Act 2006 [3 and 4], which extends the government's ability to outlaw terrorist organizations that promote and encourage or may be thought to encourage terrorism [5]. In 2007 they launched the 'Prevent Strategy' to prevent the radicalization of youths in Great Britain and block networks that support terrorists [6]. The internet has become the main tool used by terrorists since it can be accessed anywhere and it gives access to a wide spectrum of ideological material that may be translated into multiple languages [7]. Their main goals in using the internet are often research, communication, training, fundraising, media operations, radicalization and recruitment [8].

This paper is structured as follows: Section 2 provides some discussion on work related to sentiment analysis and SentiWordNet. SentiWordNet is a lexical resource that supports opinion mining by assigning a positivity score and a negativity score to each WordNet. Section 3 discusses the research question and this is followed by details of the data

collection in section 4. The system technique was developed to assign and measure the affect and sentiment found in the communication of web forums, as described in section 5. Finally, methods of model building and results analyses are presented in sections 6.

II. RELATED WORK

The term 'sentiment' was used by [9] and [10] in reference to the automatic analysis of evaluative text, and the tracking of predictive judgments and analysis of market sentiment in [11]. After that, the term 'opinion mining' was brought to the WWW conference by [12]. They mentioned that the ideal opinion-mining tools would press a set of search results for a given item, generating a list of product attributes and aggregating opinions about each of them [11]. Sentiment analysis has been considered in many research fields, such as [13] where sentiment analysis was used to analyze video comments and user profiles. In [14], the structure of lexical contextual sentences was used to classify sentiment classification from online customer reviews. In [15], SentiWordNet was used for classifying movie reviews in German. In addition, SentiWordNet was used in [16] for sentiment classification of reviews. As far as we are concerned, there are some papers that have used data from websites, blogs and forums but they have conducted testing using Machine Learning and there are no existing papers that have used data from radical web forums for testing with SentiWordNet

III. RESEARCH QUESTION

The internet has become the main tool of radicals, extremists and terrorists since it can be accessed anywhere and allows access to a wide spectrum of ideological material that can be translated into multiple languages [7]. Opinions and emotions are used on the internet for communication and can be related to and involve radical ideologies. The terrorists' main goals in using the internet are often research, communication, training, fundraising, media operations, radicalization and recruitment [8]. This paper presents our research on sentiment analysis and the detection of radical content. In particular, this research analyzes an existing technique in an attempt to answer the research question 'How effective is SentiWordNet for detecting opinions and emotions on the internet?'

IV. DATA

Two forums were selected for use in the research: Montada

Manuscript received September 9, 2012; revised November 15, 2012.

Tawunrat Chalothorn and Jeremy Ellman are with Computing, Engineering & Information Sciences, University of Northumbria at Newcastle, Newcastle Upon Tyne, United Kingdom (e-mail: tawunrat.chalothorn@unn.ac.uk).

and Qawem. Both of them use the Arabic language. They were selected by asking 21 people who are Arabic speakers which websites they think might have content related to radical Islamic ideologies. The results showed that Qawem and Montada are in the highest range.

V. METHODS

The overall process consisted of data collection, model building and result analysis, as shown in Fig. 1. The data collection phase has been described in the previous section. After that, 500 sentences of each forum were translated manually for use in the experiment. Model building was written using Python programming language. The model building phase was started by splitting sentences into words and reducing the high-frequency text (stopwords) in the sentences. Samples of stopwords can be found in Table 1. Words were stored in a bag of words (BOW) and part of speech (POS) was used, as shown in Table 2, for tagging words and knowing the position of each word in the sentence. Lexicon, WordNet and SentiWordNet were used for assigning positive and negative scores of each synset in each word [13].

The formulas for calculating positive and negative scores were taken from [17], as shown in (1) and (2). The final scores of sentences were calculated using a formula taken from [14], as shown in (3). The scores of sentences were applied using the rule that if the sentence had a positive score more than or equal to its negative score, then the sentence would be classified as positive. Otherwise it would be negative. Example of sentences can be found in Table 3.

$$Pos_weight = \left[\frac{pos}{senses} \right] \tag{1}$$

$$Neg_weight = \left[\frac{neg}{senses} \right] \tag{2}$$

pos is the number of lemma that have $Pos(s)(i) \geq Neg(s)(i)$ and $Pos(s)(i) \neq 0$; *neg* is the number of lemma that have $Neg(s)(i) \geq Pos(s)(i)$ and $Neg(s)(i) \neq 0$; and *senses* is the total number of lemma in synsets.

$$Sentence_score = \left[\frac{\sum_{i=1}^n Score(i)}{n} \right] \tag{3}$$

Sentence_score is positive or negative or negative scores of sentences; *Score(i)* is the positive or negative scores of the word in sentences; and *n* is the number of words in sentences.

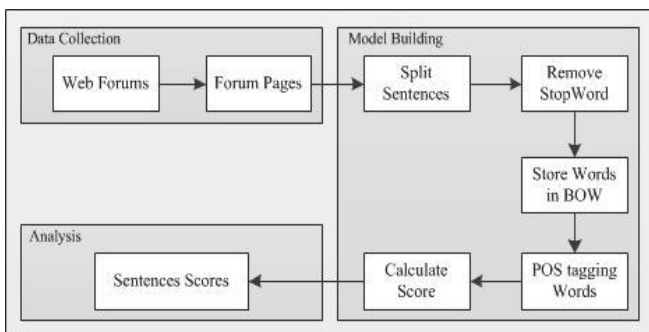


Fig. 1. Overall process of the system

TABLE I: SAMPLES OF STOPWORDS

Stopwords
['I', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'her', 'hers', 'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', ...]

TABLE II: PARTS OF SPEECH LABELS

POS Meaning	POS Tag	SentiWordNet Tag
Verb	VB, VBD, VBG, VBN, VBP, VBZ	V
Noun(s)	NN, NNS, NNP, NNPS	N
Adverb(s)	RB, RBR, RBS	R
Adjective(s)	JJ, JJR, JJS	A

TABLE III: EXAMPLE OF SENTENCES WITH SENTIMENT POLARITY¹

Arabic and English Translation	Sentiment Polarity	
	Positive	Negative
الله يلعن الوهابية والسلفية اعداء الدين Allah curse the Salafi and Wahhabi enemies of religion.	0.000	0.033
اللهم انزل سخطك على يهود آل خليفة Allah send down your wrath on the Jews of Al-Khalifa.	0.019	0.100

VI. RESULT

The model building of sentiment was applied to the web forums Montada and Qawem for analysis of the results. After removing stopwords, the rest of the sentences were used for analysis. The search function in the system was used to extract statistics of corpus for getting information about the frequency of words that were used in the forums, as shown in Fig. 2 and Fig. 3. The content in the forums was expected to be manipulated by religion and ideology. Both results showed that the top 10 most frequently used words were words related to religion, such as ‘God’ and ‘Allah’. ‘God’ was found to be the most frequently used word in both forums. In the comparison between Qawem and Montada, it was found that Qawem contained more words related to radical ideology than Montana, such as ‘curse’ and ‘enemies’. At the below, Fig. 4 and 5 show the results of the sentiment analysis of postings as percentages. The results show that the Montada forum has less negative postings than the Qawem forum. In particular, the radical affect is quite strong in the communication found in the Qawem forum. Nearly 35% of the postings in Qawem have a negative score between 0.050 and 0.100, while Montada has less than 15% of postings in the same score range. On the other hand, the positive scores of postings in the Montada forum were higher than those in

¹ These are not views expressed or implied by the author or the University of Northumbria at Newcastle.

the Qawem forum, except in the range from 0.100 to 0.150.

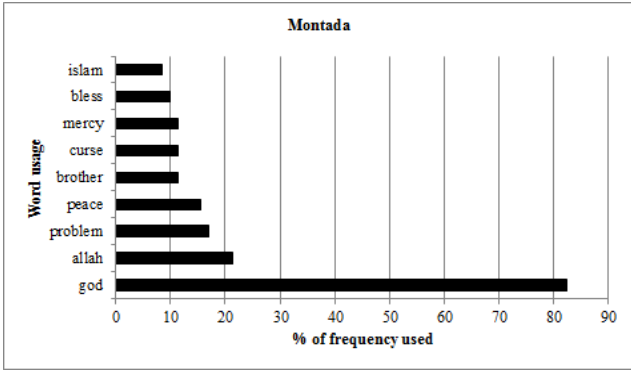


Fig. 2. Top high frequency words in Montada

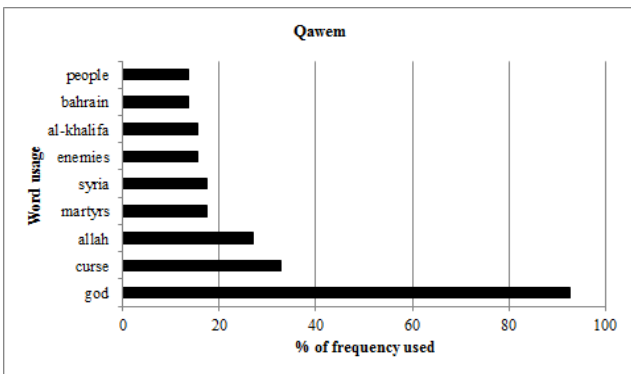


Fig. 3. Top high frequency words in Qawem

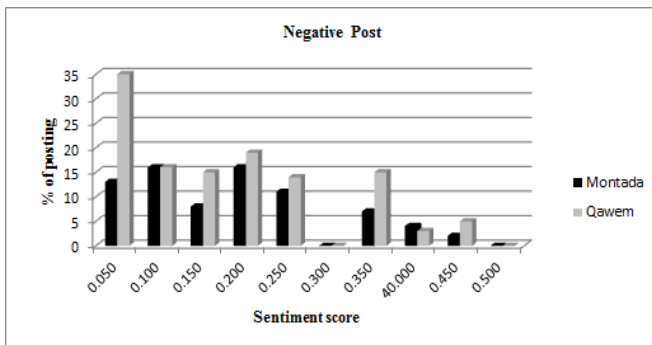


Fig. 4. Negative scores of sentiment analysis

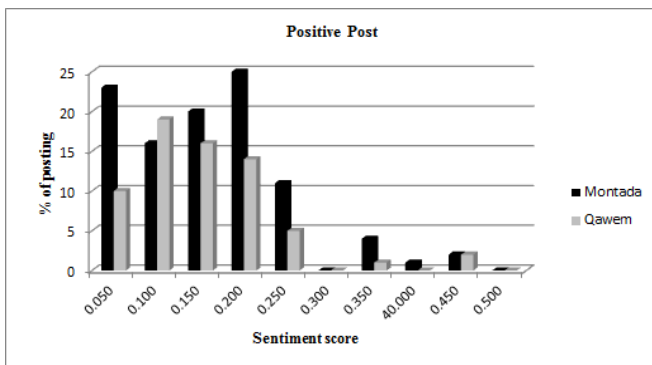


Fig. 5. Positive scores of sentiment analysis

VII. CONCLUSION

In this paper we have presented an analysis of two web forums, Montada and Qawem. They were chosen because their content relates to radicalization. The approach of model building and the results were explained. The system was developed using SentiWordNet, WordNet and NLTK for analysis of data. Overall, the results show that Qawem has more radical content than Montada. For future work, a comparative human evaluation can take place. We will ask people to rate sentences and see how their opinions on a rating scale compare to those of the model. Moreover, other techniques of sentiment analysis, such as SentiFul and SentiStrength, will be used for analyzing radical content. The aim will be to find suitable techniques for use in a model to be developed in the future.

REFERENCES

- [1] J. Glaser, J. Dixit, and D. P. Green, "Studying Hate Crime with the Internet: What Makes Racists Advocate Racial Violence?" *Journal of Social Issues*, vol. 58, pp. 177-193, 2002.
- [2] H. Chen, "Intelligence and Security Informatics For International Security," *Information Sharing and Data Mining*, Springer, 2006.
- [3] HM Government. (2006). Terrorism Act 2006. [Online]. Available: <http://www.legislation.gov.uk/ukpga/2006/11/section/1>
- [4] J. C. Paye and J. H. Membrez, *Global war on liberty*, Telos Press Pub., 2007.
- [5] E. Parker, "Implementation of the UK Terrorism Act 2006 - The Relationship between Counterterrorism Law, Free Speech, and the Muslim Community in the United Kingdom versus the United States," *Emory International Law Review* 21 *Emory Int'l L. Rev.*, pp. 711-758, 2007.
- [6] G. Morgan, "Government to block terrorist web sites," in *computing.co.uk*, ed, 2011.
- [7] P. Nessler, "Ideologies of Jihad in Europe," *Terrorism and Political Violence*, vol. 23, pp. 173-200, 2011.
- [8] B. Mantel, "Terrorism and the Internet: should web sites that promote terrorism be shut down?" ed: Washington, D. C. :CQ Press, 2009, pp. 129-155.
- [9] S. Das and M. Chen, "Yahoo! for Amazon: Extracting market sentiment from stock message boards," in *Proc. the Asia Pacific Finance Association Annual Conference APFA*, 2001.
- [10] R. Tong, "An operational system for detecting and tracking opinions on-line," in *Proc. SIGR Workshop on Operational Text Classification, New Orleans, Louisiana*, 2001.
- [11] B. Pang and L. Lee, "Opinion Mining and Sentiment Analysis," *Found. Trends Inf. Retr.*, vol. 2, pp. 1-135, 2008.
- [12] K. Dave, S. Lawrence, and D. M. Pennock, "Mining the peanut gallery: opinion extraction and semantic classification of product reviews," in *Proc. the 12th International Conference on the World Wide Web, Budapest, Hungary*, 2003.
- [13] A. Bermingham, M. Conway, L. McInerney, N. O'Hare, and A. F. Smeaton, "Combining Social Network Analysis and Sentiment Analysis to Explore the Potential for Online Radicalisation," in *Proc. the 2009 International Conference on Advances in Social Network Analysis and Mining*, 2009.
- [14] A. Khan and B. Baharudin, "Sentiment classification using sentence-level semantic orientation of opinion terms from blogs," presented at the National Postgraduate Conference (NPC), 2011.
- [15] K. Denecke, "Using SentiWordNet for multilingual sentiment analysis," presented at the IEEE 24th International Conference, the Data Engineering Workshop, 2008.
- [16] A. Hamouda and M. Rohaim, "Reviews Classification Using SentiWordNet Lexicon," *The Online Journal on Computer Science and Information Technology (OJCSIT)*, vol. 2, pp. 120-123, 2011.
- [17] A. Neviarouskaya, H. Prendinger, and M. Ishizuka, "Textual Affect Sensing for Sociable and Expressive Online Communication," in *Proc. the 2nd International Conference on Affective Computing and Intelligent Interaction, Lisbon, Portugal*, 2007.