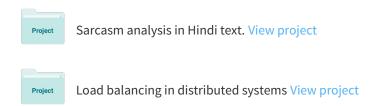
See discussions, stats, and author profiles for this publication at: https://www.researchgate.net/publication/309781364

Sarcasm Analysis on Twitter Data Using Machine Learning Approaches

Chapter	· November 2016		
CITATIONS 0	5	READS 155	
4 authors, including:			
	Santosh Kumar Bharti National Institute of Technology Rourkela 13 PUBLICATIONS 11 CITATIONS SEE PROFILE		Ramkrushna Pradhan National Institute of Technology Rourkela 2 PUBLICATIONS 0 CITATIONS SEE PROFILE
	Sanjay Kumar Jena National Institute of Technology Rourkela 182 PUBLICATIONS 1,407 CITATIONS SEE PROFILE		

Some of the authors of this publication are also working on these related projects:



All content following this page was uploaded by Santosh Kumar Bharti on 11 November 2016.

The user has requested enhancement of the downloaded file. All in-text references underlined in blue are added to the original document and are linked to publications on ResearchGate, letting you access and read them immediately.

Sarcasm Analysis on Twitter Data Using Machine Learning Approaches

SK Bharti, RK Pradhan, KS Babu and SK Jena

Abstract Sarcasm analysis, being one of the toughest challenges in natural language processing (NLP), has become a hot topic of research these days. A lot of work has already been done in the field of sentiment analysis, but there are huge challenges still being faced in identification of sarcasm. The property of sarcasm that makes it difficult to analyses and detect is the gap between its literal and intended meaning. Detecting sentiment in social media like Facebook, Twitter, online blogs and reviews have become an essential task as they influence every business organization. In this chapter, four approaches were proposed, namely parsingbased lexical generation algorithm, likes and dislikes contradiction, tweet contradicting universal facts and tweet contradicting temporary facts. The aim of the proposed methods is to extract text features such as lexical, hyperbole, behavioral and universal facts. Further, four machine learning classifiers, namely support vector machine, naive Bayes, maximum entropy and decision tree were deployed. Finally, we trained these classifiers using an extracted feature set to identify sarcasm in Twitter data. This work attains a considerable accuracy improvement over existing techniques.

1 Introduction

Sentiment analysis is a process that "aims to determine the attitude of a speaker or a writer on some topic" [1]. Social media has fueled the spread of user sentiments in the online space as ratings, reviews, comments, etc. The need for precise and reliable information about consumer preferences has led to increased interest towards analysis of social media content. For many businesses, online opinion has turned into a kind of virtual currency that can make or break a product in the mar-

Name of First Author

Santosh Kumar Bharti, NIT Rourkela e-mail: sbharti1984@gmail.com

Name of Second Author

Korra Sathya Babu, NIT Rourkela e-mail: ksathyababu@nitrkl.ac.in

Marketplace. Sentiment analysis means monitoring social media posts and discussions, then figuring out how participants are reacting to it.

Sarcasm is a particular type of sentiment which derived from the French word "Sarcasmor" that means "tear flesh" or "grind the teeth". In simple words, it means to speak bitterly. The literal meaning is different than what the speaker intends to say through sarcasm. The Random House dictionary defines sarcasm as "a harsh or bitter derision or irony" or "a sharply ironical taunt; sneering or cutting remark". Sarcasm can also be defined as a "contrast between a positive sentiment and negative situation" [8] and vice versa. For example, "I love working on holidays". In this tweet "love" gives a positive opinion but "working on holidays" is referring to a negative situation as people do not work on holidays.

For sarcasm analysis in text, it is of paramount importance to have a rudimentary knowledge of natural language processing (NLP) that aims to acquire, understand and generate human languages such as English, Chinese, Hindi, etc. Part-of-speech (POS) tagging, parsing, tokenization, etc. are the tasks performed in NLP, which are used for sarcasm detection.

Sarcasm can be detected by considering lexical, pragmatic, hyperbole or other such features of the statement. Some features can also be developed using certain patterns such as unigram, bigram, trigram, etc. There can be features based on verbal or gestural clues such as emoticons, onomatopoeic expressions in laughter, positive interjections, quotation marks, use of punctuation which can help in detecting sarcasm. But all these features are not enough to identify sarcasm in tweets until the context of the text is not known. The machine should be aware of the context of the text and relate it to general world knowledge to be able to identify sarcasm more accurately.

In this chapter, we focused on lexical (unigram, bigram and trigram), hyperbole (intensifiers and interjections), behavioral (likes and dislikes) and universal facts as the text features for sarcasm detection in tweets. For the extraction of these features from the tweets, four algorithms were proposed, namely parsing-based lexical generation algorithm (PBLGA), likes and dislikes contradiction (LDC), tweet contradicting universal facts (TCUF), and tweet contradicting time-dependent facts (TCTDF). These algorithms produce the learned feature lists to identify sarcasm in different contexts of tweets such as a contradictory sentiment and situation, likes and dislikes contradiction, universal fact negation, and time-dependent fact negation. Various machine learning classifiers namely, support vector machine (SVM), maximum entropy (ME), naive Bayes (NB), and decision tree (DT) were deployed to classify sarcastic tweets. These classifiers are trained using the extracted features lists.

This chapter is an extended version of our previous paper "Parsing-based Sarcasm Sentiment Recognition in Twitter Data" [2] on published in the proceedings of ASONAM 2015. It discussed two algorithms, PBLGA and Interjection word start (IWS) for sarcasm detection in Twitter data. In this chapter, we include three additional algorithms to cover more sarcasm types. IWS is not considered in this ex- tended version as the algorithm doesn't fit into the machine learning approach. IWS algorithm does not generate any feature list for training the classifiers mentioned earlier. The main differences lie in the experiments that have been evaluated through various machine learning techniques and have added a better comparison with existing methods, leading to new conclusions.

Rest of the chapter is organized as follows: Section 2 describes related work. Preliminary is given in Section 3. Data collection and preprocessing are discussed in Section 4. The proposed scheme and implementation details are given in Section 5. Section 6 depicts various machine learning classifiers. Experimental results are shown in Section 7. Finally, we conclude in Section 8.

8 Conclusion

Sarcasm recognition is extremely challenging work using machine learning approach. In this chapter, we make use of the supervised machine learning approach to identify sarcasm in tweets and analyzed the performances of various classifiers namely, naive Bayes, SVM, decision tree and maximum entropy. In previous re- search, authors used only SVM to identify sarcasm so we take it further and applied few more classifiers and analyzed the performances metrics of various ma- chine learning approaches. Proposed algorithms attain a much better result set as compared to the previously existing work in this domain. The PBLGA approach gives us highest accuracy and recall using the decision tree and lowest for the naive Bayes classifier. The naive Bayes classifier gives consistently worse performance for all the three out of four approaches we have used and the decision tree gives us the best results for the remaining three. Thus we see that the approaches examined by us are much better suited to resolving the sentiment from the tweets we acquire and detecting the sarcasm element if present. In future, we will target to detect sarcasm in speech data and image data. Sarcasm detection is still open for research in several domain other than Twitter data.

References

- Barbieri, F., Saggion, H., Ronzano, F.: Modelling sarcasm in twitter, a novel approach. In Proceedings of the 5th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pp. 50–58, ACL, USA (2014)
- Berger, A. L., Pietra, V. J. D., Pietra, S. A. D.: A maximum entropy approach to natural language processing. In ACL, 22, 39?-71 (1996)
- Bharti, S.K., Babu, K.S., Jena, S.K.: Parsing-based sarcasm sentiment recognition in twitter data. In International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pp. 1373–1380, IEEE/ACM, France (2015)
- 4. Bhattacharyya, P., Verma, N.: Automatic lexicon generation through wordnet. In International WordNet Conference (GWC), pp. 226–233, Brno (2004)
- Carvalho, P., Sarmento, L., Silva, M. J., De Oliveira, É.: Clues for detecting irony in usergenerated contents: oh...!! it's so easy;-). In Proceedings of the 1st international CIKM workshop on Topic-sentiment analysis for mass opinion, pp. 53–56, ACM, USA (2009)
- Chaumartin, F. R.: UPAR7: A knowledge-based system for headline sentiment tagging. In Proceedings of the 4th International Workshop on Semantic Evaluations, pp. 422–425, ACL, USA (2007)
- Davidov, D., Tsur, O., Rappoport, A.: Semi-supervised recognition of sarcastic sentences in twitter and amazon. In Proceedings of the Fourteenth Conference on Computational Natural Language Learning, pp. 107–116, ACL, USA (2010)
- 8. Esuli, A., Sebastiani, F.: Sentiwordnet: A publicly available lexical resource for opinion mining. In Proceedings of LREC, **6**, 417–422 (2006)
- Filatova, E.: Irony and sarcasm: Corpus generation and analysis using crowdsourcing. In LREC, pp. 392–398, 2012.

- Gautam, G., Yadav, D.: Sentiment analysis of twitter data using machine learning approaches and semantic analysis. In Contemporary Computing (IC3), In Seventh International Conference on Contemporary Computing (IC3), pp. 437?-442, IEEE, India (2014)
- Gonzlez-Ibnez, R., Muresan, S., Wacholder, N.: Identifying sarcasm in twitter: a closer look. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers, 2, 581–586 (2011)
- 12. Kreuz, R. J., Caucci, G. M.: Lexical influences on the perception of sarcasm. In Proceedings of the Workshop on computational approaches to Figurative Language, pp. 1–4, ACL, USA (2007)
- Kreuz, R. J., Roberts, R. M.: Two cues for verbal irony: Hyperbole and the ironic tone of voice. In International journal of Metaphor and symbol, 10, 21–31 (1995)
- Liebrecht, C. C., Kunneman, F. A., van den Bosch, A. P. J.: The perfect solution for detecting sarcasm in tweets# not. In 4th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pp. 29–37, ACL, Georgia (2013)
- Liu, P., Chen, W., Ou, G., Wang, T., Yang, D., Lei, K.: Sarcasm detection in social media based on imbalanced classification. In 15th International Conference on Web-Age Information Management, pp. 459–471, Springer, China (2014)
- Lukin, S., Walker, M.: Really? well. apparently bootstrapping improves the performance of sarcasm and nastiness classifiers for online dialogue. In Proceedings of the Workshop on Language Analysis in Social Media, pp. 30–40, ACL, Georgia (2013)
- Lunando, E., Purwarianti, A.: Indonesian social media sentiment analysis with sarcasm detection. In International Conference on IEEE on Advanced Computer Science and Information Systems (ICACSIS), pp. 195–198, IEEE, Bali (2013)
- Marcus, M. P., Marcinkiewicz, M. A., Santorini, B.: Building a large annotated corpus of English: The Penn Treebank. ACL, 19, 313–330 (1993)
- McCallum, A., Nigam, K.: A comparison of event models for naive Bayes text classification. In AAAI-98 workshop on learning for text categorization, 752, 41?-48 (1998)
- Pang, B., Lee, L., Vaithyanathan., S.: Thumbs up?: sentiment classification using machine learning techniques. In Proceedings of the Association for Computational Linguistics conference on Empirical methods in natural language processing, textbf10, 79–86 (2002).
- Pedersen, T.: A decision tree of bigrams is an accurate predictor of word sense. In Proceedings
 of the second meeting of the North American Chapter of the Association for Computational
 Linguistics on Language technologies, pp. 1–8, ACL, (2001)
- Pennebaker, J. W., Francis, M. E., Booth, R. J.: Linguistic inquiry and word count: LIWC 2001. In journal of Mahway: Lawrence Erlbaum Associates, 71, 1–11 (2001)
- Rajadesingan, A., Zafarani, R., Liu, H.: Sarcasm detection on twitter: A behavioral modeling approach. In Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, pp. 97–106, ACM, USA (2015)
- Riloff, E., Qadir, A., Surve, P., De Silva, L., Gilbert, N., Huang, R.: Sarcasm as contrast between a positive sentiment and negative situation. In proceedings of the Empirical methods in natural language processing, 13, 704–714 (2013)
- Rusu, D., Dali, L., Fortuna, B., Grobelnik, M., Mladenic, D.: Triplet extraction from sentences. In Proceedings of the 10th International Multiconference Information Society-IS, pp. 8–12, (2007)
- Strapparava, C., Valitutti, A.: Wordnet affect: an affective extension of wordnet. In LREC, 4, 1083–1086 (2004)
- Tayal, D. K., Yadav, S., Gupta, K., Rajput, B., Kumari, K.: Polarity detection of sarcastic political tweets. In Computing for Sustainable Global Development (INDIACom), pp. 625– 628, IEEE, Delhi (2014)
- Tripathy, A., Agrawal, A., Rath, S. K.: Classification of Sentimental Reviews Using Machine Learning Techniques. Procedia Computer Science, 57, 821–829 (2015)
- Tsur, O., Davidov, D., Rappoport, A.: ICWSM-a great catchy name: Semi-supervised recognition of sarcastic sentences in online product reviews. In Proceedings of the Fourth International AAAI Conference on Weblogs and Social Media, pp. 162–169, (2010)

- Tungthamthiti, P., Kiyoaki, S., Mohd, M.: Recognition of sarcasm in tweets based on concept level sentiment analysis and supervised learning approaches. In Proceedings of Pacific Asia Conference on Language, Information and Computing, pp. 404

 –413, ACL, Thailand (2014)
- 31. Utsumi, A.: Verbal irony as implicit display of ironic environment: Distinguishing ironic utterances from non-irony. In Journal of Pragmatics, 32, 1777–1806 (2000)
- 32. Verma, N., Bhattacharyya, P.: Automatic lexicon generation through wordnet. In International WordNet Conference (GWC), pp. 226–233, Brno (2004)