

Text Summarization for Aspect-Polarity Extraction (Review Paper)

Shrutika S. Yande

Student, ME Computer Science and Engineering, Dr. BAM University Aurangabad, India

Abstract— Sentiment Analysis is the process of determining whether a piece of writing is positive, negative or neutral. It's also known as opinion mining, deriving the opinion or attitude of a speaker. The applications of sentiment analysis are broad and powerful. The ability to extract insights from social data is a practice that is being widely adopted by organizations across the world. A shift in sentiment on social media has been shown to correlate with shifts in the stock market. Sentiment analysis is not a once and done effort. By reviewing your customer's feedback on your business regularly you can be more proactive regarding the changing dynamics in the market place. In this paper we are going to study sentence compression technique which is based upon sentiment analysis. Sentence compression method used to divide the sentence into various tokens and then extract the aspect i.e. subject and the polarity i.e. positive, negative and neutral opinion of sentence. Along with this it not only extracts the polarity of whole sentence but also various aspects covered within sentence.

Keywords: - Sentiment analysis, tokens, polarity, aspect, sentence compression

I INTRODUCTION

Sentiment analysis is a way to evaluate written or spoken language to determine if the expression is favourable, unfavourable, or neutral, and to what degree. Sentiment analysis is critical because helps you see what customers like and dislike about you and your brand. A common use case for this technology is to discover how people feel about a particular topic. Sentiment analysis is extremely useful in social media monitoring as it allows us to gain an overview of the wider public opinion behind certain topics. The opinions of others have a significant influence in our daily decision-making process. These decisions range from buying a product such as a smart phone to making investments to Choosing a school—all decisions that affect various aspects of our daily life. Before the Internet, people would seek opinions on products and services from Sources such as friends, relatives, or consumer reports. However, in the Internet era, it is much easier to collect diverse opinions from different people around the world. People look to review sites (e.g., CNET,

Epinions.com), e-commerce sites (e.g., Amazon, eBay), online opinion sites (e.g., Trip Advisor, Rotten Tomatoes, Yelp) and social media (e.g., Facebook, Twitter) to get feedback on how a particular product or service may be perceived in the market. Similarly, organizations use surveys, opinion polls, and social media as a mechanism to obtain feedback on their products and services. Sentiment analysis or opinion mining is the computational study of opinions, sentiments, and emotions expressed in text. The use of sentiment analysis is becoming more widely leveraged because the information it yields can result in the monetization of products and services. For example, by obtaining consumer feedback on a marketing campaign, an organization can measure the campaign's success or learn how to adjust it for greater success. Product feedback is also helpful in building better products, which can have a direct impact on revenue, as well as comparing competitor offerings.

Sentiment analysis can occur at different levels: document level, sentence level or aspect/feature level. Document Level Classification In this process, sentiment is extracted from the entire review, and a whole opinion is classified based on the overall sentiment of the opinion holder. The goal is to classify a review as positive, negative, or neutral. Example "I bought an iPhone a few days ago. It is such a nice phone, although a little large. The touch screen is cool. The voice quality is clear too. I simply love it!" Is the review classification positive or negative? Document level classification works best when the document is written by a single person and expresses an opinion/sentiment on a single entity. Sentence Level Classification process usually involves two steps: first is Subjectivity classification of a sentence into one of two classes: objective and subjective and second is Sentiment classification of subjective sentences into two classes: positive and negative. An objective sentence presents some factual information, while a subjective sentence expresses personal feelings, views, emotions, or beliefs. Subjective sentence identification can be achieved through different methods such as Naïve Bayesian classification. However, just knowing that sentences have a positive or negative opinion is not sufficient. This is an intermediate step that helps filter out sentences with no opinions and helps determine to an extent if sentiments about entities and their aspects are positive or negative. A subjective sentence may contain multiple opinions and subjective and factual clauses. This study paper based upon sentence level sentiment analysis.



II RELATED WORK

There are many existing works on aspect extraction. "Sentence Compression for Aspect-Based Sentiment Analysis" model (Che, Wanxiang; 2015) proposes Framework for using a sentiment sentence compression model Sent Comp for aspect-based sentiment analysis. Different from the common sentence compression model, Sent Comp not only compress the redundancy in the sentiment sentence, but also needs to retain the polarity-related information to maintain the sentences' original polarities. Thus, the over-natural and spontaneous sentiment sentences will be compressed into more formal and easier-to-parse sentences after using this proposed Sent_Comp model. 90% accuracy is achieved using this model. 'Sent comp' model removes the repeated, unwanted words but it preserves polarity based information that is necessary for sentiment analysis. But this model has loop hole that is it does not provide improved Performance as the aspect-based model needs critical requirements [1].

An extractive supervised two-stage method is a method for sentence compression which is Presented by Dimitrios Galanis, Ion Androutsopoulos (2010) [18]. This method shows that it generates candidate compressions by removing branches from the source sentences dependency tree using a Maximum Entropy classifier (Berger et al., 2006) at very first stage. In next stage, it chooses and displays the best candidate compressions from given set of candidate compression using a Support Vector Machine. The candidate compressions in this type are ranked using a function. These expressions take into account the inverse document frequencies of the words, and their depths in the source dependency tree but it does not support more complex transformations, instead of only removing words and experiment with different sizes of training data.

A Rule-Based Approach to Aspect Extraction from Product Reviews proposed by Soujanya Poria and Erik Cambria in year 2010. Rule based approach used to solve the problem of aspect extraction from product reviews by proposing a method called rule-based approach that expose some common sense knowledge and tree structure for sentence dependency to detect both explicit and implicit aspects. But it may happen that user never uses common-sense words which do not indicate the aspect, so in that case this approach may create some noisy results. This method does not represent the rules for complex aspect extraction. [17] Liu, Siyuan et al. [2] have developed SVM model 'TASC' that is topic adaptive sentiment classification. The author has developed a classifier that works on dynamic tweets. The classifier works for general features and blended labelled data for diversified topics. It increases the performance by selecting more reliable tweets after selecting unlabelled data collaboratively. It fits the unlabelled data and the features without any hinge loss for all sentiment words and connection of sentiments acquired from '@' tweets and it is designated topic adaptive features. They

had even designed an TASC-t that is a timeline based model on dynamic tweets that has achieved notable accuracy and F-score. Finally a visualizing and colour gradation based graph that shows sentiment trends.

Andrea Esuli and Fabrizio Sebastiani have developed Senti Word Net that is a lexical resource for mining opinion [3]. Word Net synset have scores 'obj(s)', 'Neg(s)', 'Pos(s)' having some numerical values that outlines how negative, positive, objective is accommodated in synset. It comes with GUI. Eight Ternary classifiers are trained to get three scores for synset. Training set and learning device for each ternary classifier is distinct and hence distinct classification outcome of the like synset is produced. Scores for opinions are calculated by normalization on the basis of assigned labels. If every ternary classifier assigns same label to synset than synset will have maximal score otherwise the score will be proportional to ternary classifiers. Soo-Min Kim and Eduard Hovy [4] have developed a model that automatically identifies pros and cons reasons in the review that is the reasons behind the liking or disliking of the product. As reviewed opinion have some reasons behind recommendation or non-recommendation of the product. So developed a technique where it automatically labels pros and cons.

Compressing complicated sentiment sentence into one that is shorter to parse. We apply a discriminative conditional random field model, with certain special features, to automatically compress sentiment sentences.

Using the Chinese corpora of four product domains, Sent Comp significantly improves the performance of the aspect-based sentiment analysis. They proposed a double propagation method for extract the A-P collocations, aspects and polarity words. This scenario is based upon the observation that there are natural syntactic relations between polarities words are used to modify the aspects. Furthermore, they also discovered that the polarity words and aspects themselves had relations in certain sentiment sentences. According this approach, in the double propagation method, we first used an initial seed polarity word lexicon and syntactic relations to extract the aspects, which can fall into a new aspect lexicon. Then, the aspect lexicon and the same syntactic relations used to extract the polarity words to expand the polarity word lexicon in return. This is an iterative procedure, i.e., this method can iteratively produce the new polarity words and the aspects back and forth using the syntactic relations.

III PROPOSED SYSTEM

The proposed an effective supervised learning algorithm, called polarity semantic orientation, for classifying reviews. There are several challenges in Sentiment analysis. At beginning the opinion word that is considered to be positive in one situation may be considered negative in another situation. A second challenge is that people don't always express opinions in a same way. The proposed system aims to collect different



opinion about the product and generate the result so that it clarifies whether product is either good to buy or not. Proposed framework consists of proposed system algorithm and architecture as follows.

A. Proposed System Algorithm :

Naïve Bayes classifier Naïve Bayes is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. It is not a single algorithm for training such classifiers, but a family of algorithms based on a common principle: all naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable. For example, a fruit may be considered to be an apple if it is red, round, and about 10 cm in diameter. A naive Bayes classifier considers each of these features to contribute independently to the probability that this fruit is an apple, regardless of any possible correlations between the colour, roundness, and diameter features. For some types of probability models, naive Bayes classifiers can be trained very efficiently in a supervised learning setting. In many practical applications, parameter estimation for naive Bayes models uses the method of maximum likelihood; in other words, one can work with the naive Bayes model without accepting Bayesian probability or using any Bayesian methods. It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as 'Naive'. Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods. Bayes theorem provides a way of calculating posterior probability $P(c|x)$ from $P(c)$, $P(x)$ and $P(x|c)$. Look at the equation below:

- $P(c|x)$ is the posterior probability of class (c, target) given predictor (x, attributes).
- $P(c)$ is the prior probability of class.
- $P(x|c)$ is the likelihood which is the probability of predictor given class.
- $P(x)$ is the prior probability of predictor.

This algorithm is preferred for this study because it is easy and fast to predict class of test data set. It also performs well in multi class prediction. When assumption of independence holds, a Naive Bayes classifier performs better compare to other models like logistic regression and you need less training data. It performs well in case of categorical input variables compared to numerical variable(s). For numerical

variable, normal distribution is assumed (bell curve, which is a strong assumption).

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability
Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Bayesian classification provides practical learning algorithms and prior knowledge and observed data can be combined. Bayesian Classification provides a useful perspective for understanding and evaluating many learning algorithms. It calculates explicit probabilities for hypothesis and it is robust to noise in input data.

B. Proposed system Architecture

Text summarization is an extension of data mining which utilizes natural language processing techniques to extract people's opinion from various businesses, shopping sites. The recent trend in internet that encourages users to contribute their opinion and suggestion created a huge collection of valuable information in the web. The Opinion mining system analyse each text and see which part contain opinionated word, which is being opinionated and who has written the opinion. Sentiment analysis analyses each opinionated word or phrase and determines its sentiment polarity orientation, whether it is positive or negative or neutral. The results will obtain the summarized opinion of a writer or speaker.

Text summarization done at sentence level. Proposed framework uses Text summarization technique to extract the exact opinion of customers about the product.

Below figure shows the overall architecture of proposed system. Architecture consists various phases, such as taking an input from users, POS tagging, sentence compression, using a naive Bayes classifier algorithm, polarity detection. Each phase performs different functions. These functionalities described as:

B.1) User Input:

It takes an input from user. Input is nothing but the reviews or comments given by user for particular product.

B.2) Tokenize - Tokenizes the text into a sequence of tokens. Due to this each word classified according to their characteristic to which they are belongs.

B.3) Sentence Compression to compress the reviews of user there are various techniques will be used such as Stemming is supposed to turn inflected forms of words down to some common root. Perception words removal Words used by users in comments so it is necessary to delete certain perception words, which do not change the meaning and the sentiment orientation of the original Sentence. Because of this sentence

compression process only keeps the basic root words and those words which exactly indicate the opinion of users.

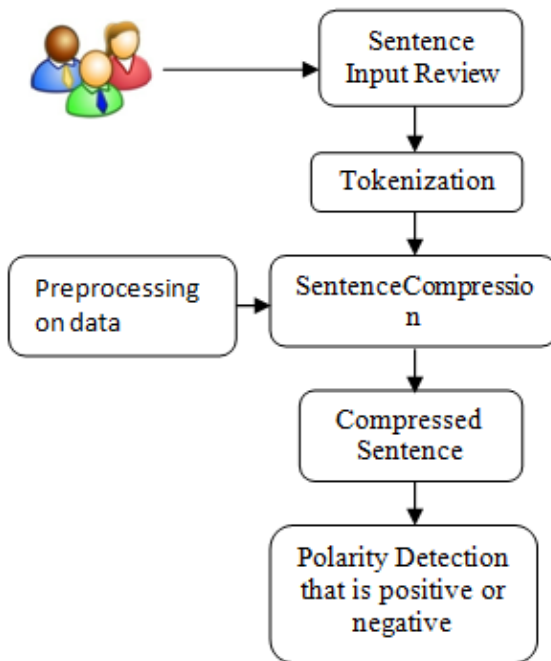


Figure 1. Proposed System Architecture

B.4) Polarity detection

It detects the polarity of sentence that is whether review is positive, negative or neutral.

V CONCLUSION

The main objective of study is to provide a summary of a large number of customer's reviews of product which sold online. It provides ease to customers to analyse single review and displays the polarity of that single statement. The proposed sentiment analysis is a subfield of opinion mining concerned with the determination of opinion and subjectivity in a text, which has many applications. This paper presents study about classifiers for sentiment analysis of user opinion and goal is to show final opinion about the product i.e. positive, negative; neutral. There are various pre-processing techniques for review compression was proposed to compress the sentence and then extract aspect and polarity from it. The proposed framework will also provide facility to determine the each aspect and its polarity independently which exist within a single statement. System aims to provide facility to users to personalize their accounts for improved performance in analysing the product opinion.

ACKNOWLEDGMENT

It is my great pleasure in expressing sincere and deep gratitude towards my guide Prof. V. S. Karwande for providing me various resources, valuable and firm suggestion, guidance and constant support throughout this work.

REFERENCES

1. Che, Wanxiang, Yanyan Zhao, Honglei Guo, Zhong Su, and Ting Liu. "Sentence Compression for Aspect-Based Sentiment Analysis." *Audio, Speech, and Language Processing, IEEE/ACM Transactions on*, 23, no. 12 (2015): 2111-2124
2. David Zajic¹, Bonnie J. Dorr¹, Jimmy Lin¹, Richard Schwartz, "Multi-Candidate Reduction: Sentence Compression as a Tool for Document Summarization Tasks." University of Maryland College Park, Maryland, USA, 2BBN Technologies 9861 Broken Land Parkway Columbia, MD 21046.
3. Kim, Soo-Min, and Eduard Hovy, "Automatic identification of Pro and Con reasons in online reviews". In *processing of the COLING/ACL on Main conference poster sessions*, pp. 483-490. Association for Computational Linguistics, 2006.
4. Hu, Mingqing, and Bing Liu, "Mining and summarizing customer reviews." In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 168-177 ACM, 2004.
5. Jin, Jian, Ping Ji, and Ying Liu. "Translating online customer opinions into engineering characteristics in QFD: A probabilistic language analysis approach". *Engineering Applications of Artificial Intelligence* 41(2015): 115-127.
6. Suanmali, Ladda, Naomie Salim, and Mohammed Salem Binwahlan. "Fuzzy logic based method for improving text summarization." *arXiv preprint arXiv:0906.4690* (2009).
7. Ghorashi, Seyed Hamid, Roliana Ibrahim, Shirin Noekhah, and Niloufar Salehi Dastjerdi. "A frequent pattern mining algorithm for feature extraction of customer reviews." In *IJCSI International Journal of Computer Science Issues*. 2012.
8. Mrs. Elakkiya.R, Mrs. Jayasudha.M2, Mr. Sivanesh Waran. "Improved Optimized Sentiment Classification On Dynamic Tweets". *IJCSMC*, Vol. 5, Issue. 6, June 2016, pg. 11 – 22, ISSN 2320-088X IMPACT FACTOR: 5.258.
9. Lu Wang, Hema Raghavan, Vittorio Castelli. "A Sentence Compression Based Framework to Query-Focused Multi-Document Summarization". Cornell University, Ithaca, NY 14853, USA, T. J. Watson Research Center, Yorktown Heights, NY 10598, USA.
10. Alejandro Molina¹, Juan-Manuel Torres-Moreno, Eric San Juan, Iria da Cunha, & Gerardo Eugenio Sierra Martínez, LIA. "Discursive Sentence Compression". *Universitat d'Avignon, IULA-Universitat Pompeu Fabra, GIL-Instituto de Ingeniería UNAM*.
11. Kapil Thadani and Kathleen McKeown. "Sentence Compression with Joint Structural Inference". Department of Computer Science Columbia University New York, NY 10025, USA.
12. Dipanjan Das Andre, F.T. Martins, "A Survey on Automatic Text Summarization" *Language Technologies Institute Carnegie Mellon University*, November 21, 2007.
13. Seyed Hamid Ghorashi, Roliana Ibrahim, Shirin Noekhah and Niloufar Salehi Dastjerdi. "A Frequent Pattern Mining Algorithm for Feature Extraction of Customer Reviews" *IJCSI International Journal of Computer Science Issues*, Vol. 9, Issue 4, No 1, July 2012 ISSN (Online): 1694-0814.



14. Kevin Knight, Daniel Marcu. "Summarization beyond sentence extraction: A probabilistic approach to sentence compression". Information Science Institute and Department of Computer Science, University of Southern California, 4676 Admiralty Way, Suite 1001, Marina Del Rey, CA 90292, USA. Received 11 May 2001.
15. Trevor Cohn and Mirella Lapata. "Sentence Compression beyond Word Deletion". Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008), pages 137–144 Manchester, August 2008.
16. Mirella Lapata. "An Abstractive Approach to Sentence Compression". ACM Transactions on Intelligent Systems and Technology, Vol. 4, No. 3, Article 41, Publication date: June 2013.
17. Soujanya Poria, Erik Cambria. "A Rule-Based Approach to Aspect Extraction from Product Reviews". Proceedings of the Second Workshop on Natural Language Processing for Social Media (SocialNLP), pages 28–37, Dublin, Ireland, August 24 2014.
18. Dimitrios Galanis and Ion Androutsopoulos. "An extractive-supervised two-stage method for sentence compression". Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the ACL, pages 885–893, Los Angeles, California, June 2010.
19. Lingling Yuan, "An Improved Naive Bayes Text Classification Algorithm In Chinese Information Processing". Proceedings of the Third International Symposium on Computer Science and Computational Technology (ISCST '10) Jiaozuo, P. R. China, 14-15, August 2010, pp. 267-269.
20. Arjun Mukharjee Bing Liu, "Aspect Extraction through Semi-Supervised Modelling". Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics, pages 339–348, Jeju, Republic of Korea, 8-14 July 2012.
21. Vidisha M. Pradhan Jay Vala Prem Balani. "A Survey on Sentiment Analysis Algorithms for Opinion Mining". International Journal of Computer Applications (0975 – 8887) Volume 133 – No.9, January 2016