

Opinion Targets and Opinion Words from Online Reviews Based on the Word Alignment Model

Syama .N¹, Mr. G. Pandiyan²

M.Phil Scholar, Department of Computer Science, RVS College of Arts and Science (Autonomous), Coimbatore¹

M.Phil, B.Ed, Assistant Professor, School of Computer Studies²

Abstract: Opinion mining has gained increasing attention and shown great practical value in recent years. Extracting opinion words and targets is a main task in opinion mining. For the purpose of customer and business perspective, the task of scanning these reviews manually is computational burden. Hence, to process reviews automatically and summarizing them in suitable form is more efficient. The distinguished problem of producing opinion summary addresses is how to determine the mood, and opinion expressed in the review with respect to a numerical feature value. This paper proposes a novel approach with a hybrid algorithm which combines Expectation Maximation (EM) algorithm. It focused on the main task of opinion mining called as opinion summarization. The extraction of product feature, technical feature value and opinion are critical for opinion summarization as they affect the performance significantly. The proposed approach consists of a software system in which mining of product feature, technical feature value and opinion is performed. The main motto of this software system is to recognize the technical feature value depending on review, which the reviews are summarized. This software is helpful for humans to understand the technical values expressed in the reviews. It represents relations between opinion words and targets, which is employed to measure the confidence of each candidate from opinion words and targets datasets. The words or targets with high confidence are kept in their respective datasets and the rest are removed as false results which are used to refine extraction rules. K-nearest neighbor classifier - used for classify the extracted data's in an opinion mining. Experimental results Shows the effectiveness of proposed method and finally, candidates with higher confidence are extracted as opinion targets or opinion words.

Index Terms: Opinion Mining, Opinion Targets Extraction, Opinion Words Extraction, KNN.

1. INTRODUCTION

In recent year, online product reviews have been considered as a valuable source of information to assist people in making buying decisions. Most of prior studies on the effect of online product reviews have utilized the factors which manufactures cannot control by themselves, such as the number of reviews, the average review rating, as independent variables in their regression models. However, those factors cannot provide direct implications for manufacturers. For example, managers cannot easily increase the number of reviews to rise the product price or demand. In contrast, they have to trace the causes of why the amount of reviews grows. Thus, in order to offer more straightforward suggestions, we adopt the concept of hedonic analysis which decomposes the demand of a commodity into several product features to identify which of them impact its demand mostly.

Mining Opinion Words and Targets from Online Reviews in a hybrid Framework [1] In the following refinement process, an Opinion Relation Graph (ORG) is modeled to represent relations between opinion words and targets, which is employed to measure the confidence of each candidate from opinion words and targets datasets. The words or targets with high confidence are kept in their respective datasets and the rest are removed as false results which are used to refine extraction rules with an

Automatic Rule Refinement (ARR) method. Update ORG model and repeat the joint process of propagation and refinement until ORG model reaches stable. Experimental results on both English and Chinese datasets demonstrate the effectiveness of proposed method comparing with the-state-of-the-art methods. Hedonic Analysis for Consumer Electronics Using Online Product Reviews [2] we take smartphone market as our research target and there are reasons why we select it. Frist, it is a booming market in the United States and there may be a large number of reviews or comments about smartphones strewed online.

Second, it matches the definition of high involvement product and there are some researches also classify smartphone as a kind of high involvement. A Unified Framework for Fine-Grained Opinion Mining from Online Reviews [3] unified framework for fine-grained opinion mining, combining propagation with refinement in a dynamic and iterative process. In the propagation process, syntactic patterns are chosen as opinion relations to extract new opinion words and targets. Besides, syntactic patterns are further generalized to make them more flexible and scalable. In the refinement process, a three-layer opinion relations graph (ORG) model is constructed based on three types of candidates: opinion word candidates, opinion target candidates and syntactic pattern candidates. A

sorting algorithm based on ORG model is proposed to rank all the candidates in their own type, and low-rank candidates are removed from candidate datasets. Extracting Opinion Explanations from Chinese Online Reviews [4] the explanations of opinions, which are potentially valuable for many applications, are totally ignored. To address this specific research challenge, we propose an approach to extract the explanation of reason and/or consequence behind an opinion via learning word pairs and using causal indicators from Chinese online reviews. We also improve our word pair based method by constructing clusters of word paris. Experiments on a Chinese business review corpus show that our method is feasible and effective. Estimating the Sentiment of Social Media Content for Security Informatics Applications [5] hey outperform several standard techniques for the task of inferring the sentiment of online movie and consumer product reviews. Additionally, we illustrate the potential of the methods for security informatics by estimating regional public opinion regarding Egypt's unfolding revolution through analysis of Arabic, Indonesian, and Danish (language) blog posts.

The proposed method to extracting opinion words from the online review is a main task in opinion mining. A method a hybrid algorithm which combines Expectation Maximization (EM) algorithm used for mining and k-nearest neighbor classifier used for classification.

2. RELATED WORKS

Existing approaches on extracting opinion words and opinion targets basically follow two frameworks. One is the pipeline. Under this framework, candidates of opinion words and opinion targets are firstly generated, and then false candidates are filtered by using refinement methods. Our proposed method used for mining a online reviews using a hybrid algorithm which combines Expectation Maximization (EM) algorithm for identify the opinion of the users. K-nearest neighbor classifier used for classification an approach of Opinion Mining for online marketing Using Sentiment Thesaurus and Concept Search Engine [1] The crux of this research work is to do a summarization of all the customer reviews of a product. This summarization task calls out the specific feature details like opinions of the product unlike the conventional text summarization including positive and negative. No original sentences of reviews are summarized by selecting or rewriting to identify the important concepts as in the classic text summarization. The interest is limited to the mining of opinion and product features captured as part of the summarization task. Cross-Domain Sentiment Analysis of Product Reviews by Combining Lexicon-based and Learn-based Techniques [2] Combines Lexicon-based and Learn-based techniques (CLL) to analyze the cross-domain sentiment of Chinese product reviews. We first build three domain lexicons based on the basic lexicon and corpus from three domains containing books, hotels and electronics. Furthermore, we use four categories

of features (including 16 features in total) to build six classifiers. We conduct a series of experiments to evaluate our proposed CLL by using different lexicons and different classifiers. Our experimental results show that domain lexicons outperform the basic lexicon no matter in which domain. Our method CLL performs better than state-of-the-art methods in domains of books and hotels, and is slightly inferior in the domain of electronics.

OpinMiner: Extracting Feature-Opinion Pairs with Dependency Grammar from Chinese Product Reviews [3] Chinese product review. We propose a method based on Chinese dependency grammar to extract feature-opinion word pairs. Specifically, we use Chinese dependency grammar to set several rules, and then we make use of these rules to extract candidate feature-opinion word pairs. Finally, we filter out mismatched feature-opinion words pairs by feature ranking and Named Entity Recognition (NER) system. Mining consumer's opinion target based on translation Model and word representation [4] a system using translation model as well as word representation method to obtain user's interests on dataset in Chinese.

To release the word segmentation error, a finely generated system with new Chinese word detection module is proposed. The experiments on two corpus subjected on digital product verify the effective of our method. Recommending Products to Customers using Opinion Mining of Online Product Reviews and Features [5] We have used natural language processing to automatically read reviews and used Naive Bayes classification to determine the polarity of reviews. We have also extracted the reviews of product features and the polarity of those features. We graphically present to the customer, the better of two products based on various criteria including the star ratings, date of review, the helpfulness score of the review and the polarity of reviews.

OpinMiner: Extracting Feature-Opinion Pairs with Dependency Grammar from Chinese Product Reviews [6] method based on Chinese dependency grammar to extract feature-opinion word pairs. Specifically, we use Chinese dependency grammar to set several rules, then we make use of these rules to extract candidate feature-opinion word pairs. Finally, we filter out mismatched feature-opinion words pairs by feature ranking and Named Entity Recognition (NER) system. Experiment shows that our method in Precision is rather high.

A Logic Programming Approach to Aspect Extraction in Opinion Mining [7] Double propagation (DP) method is implemented using 8 ASP rules that naturally model all key ideas in the DP method. Our experiment on a widely used data set also shows that the ASP implementation is much faster than a Java-based implementation. Syntactical approach has its limitation too. To further improve the performance of syntactical approach, we identify a set of general words from Word Net that have little chance to be an aspect and prune them when extracting aspects.

3. PROPOSED WORK

Online reviews usually have informal writing styles, including grammatical errors, typographical errors, and punctuation errors. This makes prone to generating errors. We present the main framework of our method. As mentioned before, we regard extracting opinion target words/sentence as a co-ranking process. We assume that all nouns/noun phrases in sentences are opinion target candidates, and all adjectives/verbs are regarded as potential opinion words, which are widely adopted by previous methods. Each candidate will be assigned a confidence, and candidates with higher confidence than a threshold are extracted as the opinion targets or opinion words. To assign a confidence to each candidate, is our basic motivation.

- Opinion system finds and extracts important topics in the text that will then be used to summarize. This system present a technique based on a hybrid algorithm which combines Expectation Maximation (EM) algorithm.
- This system helps to find the opinion from online reviews to specifies rating of the particular product, movie etc. which is give a confident for buy a products.
- We are implementing some preprocessing methods to remove the noise in the sentence and easily filter-out the words.
- Extracting features from the sentence and after that applying K-nearest neighbor classifier used for classification.

- Effectiveness of the proposed method, we select real online reviews from different domains and languages as the evaluation datasets. We compare our method to several state-of-the-art methods on opinion words extraction.

3.1 Data Collection

Opinion text in blog, reviews, comments etc. contains subjective information about topic. Reviews classified as positive or negative review. Opinion summary is generated based on features opinion sentences by considering frequent features about a topic. It is the process of collecting review text from review websites. Information retrieval techniques such as web crawler can be applied to collect the review text data from many sources and store them in database. This step involves retrieval of reviews, micro-blogs, and comments of user.

Blogs have become popular because of the niche comments shared by readers in a lucid and lively format. Textual is the norm for many of the blogs though some are art blogs, photo blogs, video blogs or vlogs music blogs or mp3 blogs and audio podcasts. The text content in the blog deals with various topics for eg comments about airways deals with hospitality, food, service etc.

3.2 Preprocessing

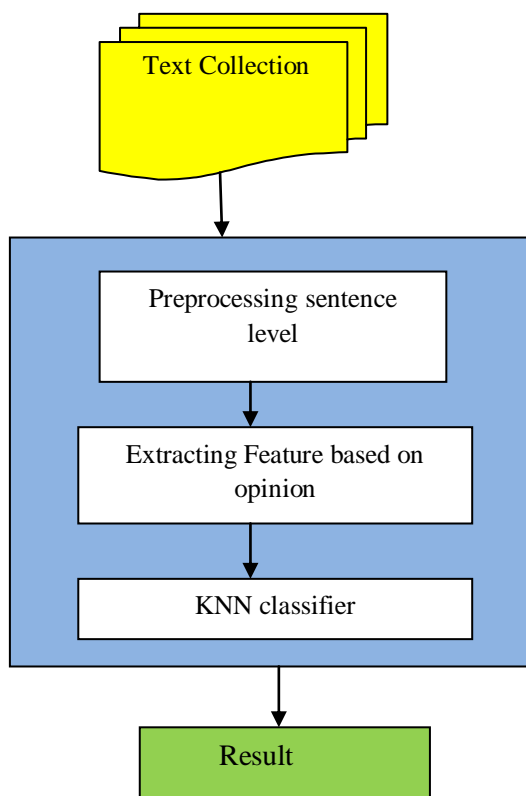
Preprocessing Algorithm receives user opinions in raw form. We implement some form of preprocessing in order to filter-out noise. Sentence splitting is a critical step in this module (opinion delimitation) since double propagation takes into account neighborhood sentences in order to propagate sentiment. Additionally in order to increase the efficiency of the extraction process we have adopted an on-line stemmer engine.

3.3 Feature Classification

It defines the polarity of document, but a positive phrase does not indicate that the user likes everything and similarly a negative phrase does not indicate that the opinion holder dislikes everything. It is a fine-grained level of classification in which polarity of the sentence can be given by three categories as positive, negative and neutral. It is defined as product attributes or components. In this approach positive or negative opinion is identified from the already extracted features. It is a fine grained analysis model among all other models. It is having a drawback that it could really cut very badly if there used any grammatically incorrect text.

3.4 Sentence level Opinion Mining

In sentence level Opinion Mining, the polarity of each sentence is calculated. The same document level classification methods can be applied to the sentence level classification problem also but Objective and subjective sentences must be found out. The subjective sentences contain opinion words which help in determining the sentiment about the entity. After which the polarity classification is done into positive and negative classes.



3.5 Feature Opinion

The knowledge resource is useful for improving the performance of the opinion mining. Opinion words lexicon is adopted in the stage of identifying opinions regarding the product features. A domain independent Lexicon and manually constructed Emoticon dictionary is used to assign polarity score (positive, negative or neutral) to opinionated words and sentences. For deciding correct polarity class of such words, revised mutual information concepts are used. These words could strengthen, weaken the surrounding opinion words' extent or even transit its sentiment orientation.

In our proposed system we implement a a hybrid algorithm which combines Expectation Maximation (EM) algorithm with k-nearest neighbor classifier.

4. METHODOLOGY

Expectation Maximation (EM) algorithm

The proposed EM algorithm, the expectation step (E-step) computes expected statistics over completions rather than explicitly forming probability distribution over completions. The system's E-step consists of storing the extracted product candidate feature, related feature opinion and technical feature value. Similarly, for the maximization step (M-step) consists of model re-estimation which can be thought of as „maximization“ of the expected log-likelihood of the data. In the system, M-step consists of following step:

- The stored technical feature values of particular product feature candidate are clustered in one group.
- Statistical calculations are carried out on those technical feature values as they need to be grouped into three different classes as best, average and poor so that the summary for the particular product feature candidate is generated
- For grouping the technical feature values, the statistical method called standard deviation is used. Standard deviation is basically shows how much variation exists from the average (mean) or expected value so that the values get distributed into classes. The standard deviation formula is as follows: For grouping the technical feature values, the statistical method called standard deviation is used. Standard deviation is basically shows how much variation exists from the average (mean) or expected value so that the values get distributed into classes. The standard deviation formula is as follows:

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

Where $\{x_1, x_2, \dots, x_n\}$ are the technical feature values extracted from the reviews and \bar{x} is the mean value of these technical feature values, while the denominator N stands for the number of reviews and S is the standard deviation.

- In the proposed system, the standard deviation is calculated using largest technical feature value, smallest technical feature value and mean of technical feature value. The smallest standard deviation value calculated among these is considered and the product feature candidate is assigned to that particular class which can be best, average or poor which is considered as processed opinion by the system.

- The summary is generated from the above statistical calculations is in tree view form in which the sentiment processed by the system depending on the technical feature value extracted is rated into three classes as good, average and poor.

k-Nearest Neighbour (k-NN) classifier in order to combine it with the lexicon based. The k-Nearest Neighbour (k-NN), a popular example-based classifier, is also known as lazy learning because it postpones the decision to make generalizations beyond the training data until it has located every single new incidence. In order to classify a review, the k-NN classifier roughly ranks the review among the training reviews, before classifying it according to the k most similar neighbours. When presented with a test review d , the classifier will locate the k nearest neighbours among training reviews. The score of each nearest neighbour review that is the most similar to the test review is used as the weight of the classes of the neighbour reviews. The weighted sum in k-NN classification can be represented.

$$score(d, t_i) = \sum_{d_j \in kNN(d)} sim(d, d_j) \delta(d_j, c_i)$$

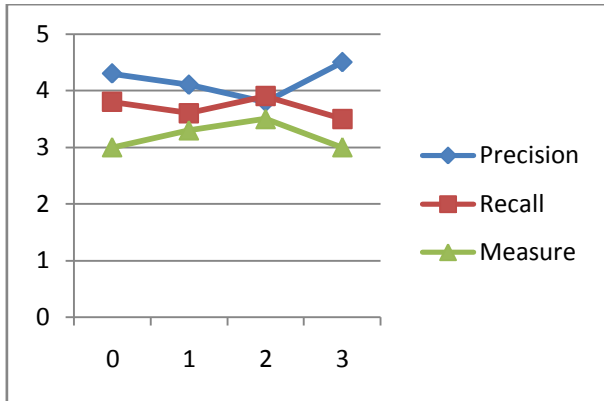
5. EXPERIMENT RESULT AND DISCUSSION

In order to evaluate the performance, we conducted an experiment on feature identification. Then, we performed an experiment for extracting feature-opinion pairs. First, we measure the identification of feature from feature-opinion pairs. We use Recall, Precision and F-measure as evaluation criteria. Then, we measure the identification of feature opinion pairs and we use feature-opinion pairs recall and precision as our evaluation criteria. We prune feature-opinion pairs by different thresholds which is percent of the number of feature-opinion pairs.

Features	Average Score	Bleu	Accuracy
Camera	0.8024		80.24%
Processor	0.7708		77.08%
Screen	0.9112		91.12%
RAM	0.7532		75.32%

We identify feature by extracting feature-opinion pairs, which can improve the precision of feature detection.

Feature-opinion pairs impose restrictions on feature so that feature and opinion word co-occur in the same sentence by certain relation.



6. CONCLUSION

In this paper we propose a method a hybrid algorithm which combines Expectation Maximation (EM) algorithm and k-nearest neighbor classifier. Which is used for extracting the users opinions through the online reviews of customer and generating the summary for those reviews by using modified Expectation Maximization(EM) algorithm. This method summarizes review depending on features and technical feature value extracted from the reviews. Then, we use these rules to extract candidate feature-opinion pairs directly. Finally, we filter out mismatched feature-opinion pairs by feature ranking and k-nearest neighbor classifier used for classification. Experimental results produced by the system shows the accuracy of the proposed algorithm.

REFERENCES

[1] M. Hu and B. Liu, "Mining and summarizing customer reviews," in Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, Seattle, WA, USA, 2004, pp. 168–177.

[2] F. Li, S. J. Pan, O. Jin, Q. Yang, and X. Zhu, "Cross-domain coextraction of sentiment and topic lexicons," in Proc. 50th Annu. Meeting Assoc. Comput. Linguistics, Jeju, Korea, 2012, pp. 410–419.

[3] L. Zhang, B. Liu, S. H. Lim, and E. O'Brien-Strain, "Extracting and ranking product features in opinion documents," in Proc. 23th Int. Conf. Comput. Linguistics, Beijing, China, 2010, pp. 1462–1470.

[4] K. Liu, L. Xu, and J. Zhao, "Opinion target extraction using wordbased translation model," in Proc. Joint Conf. Empirical Methods Natural Lang. Process. Comput. Natural Lang. Learn., Jeju, Korea, Jul. 2012, pp. 1346–1356.

[5] M. Hu and B. Liu, "Mining opinion features in customer reviews," in Proc. 19th Nat. Conf. Artif. Intell., San Jose, CA, USA, 2004, pp. 755–760.

[6] A.-M. Popescu and O. Etzioni, "Extracting product features and opinions from reviews," in Proc. Conf. Human Lang. Technol. Empirical Methods Natural Lang. Process., Vancouver, BC, Canada, 2005, pp. 339–346.

[7] G. Qiu, L. Bing, J. Bu, and C. Chen, "Opinion word expansion and target extraction through double propagation," *Comput. Linguistics*, vol. 37, no. 1, pp. 9–27, 2011.

[8] B. Wang and H. Wang, "Bootstrapping both product features and opinion words from chinese customer reviews with crossinducing," in Proc. 3rd Int. Joint Conf. Natural Lang. Process., Hyderabad, India, 2008, pp. 289–295.

[9] B. Liu, *Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data*, series Data-Centric Systems and Applications. New York, NY, USA: Springer, 2007.

[10] G. Qiu, B. Liu, J. Bu, and C. Che, "Expanding domain sentiment lexicon through double propagation," in Proc. 21st Int. Jont Conf. Artif. Intell., Pasadena, CA, USA, 2009, pp. 1199–1204.

[11] R. C. Moore, "A discriminative framework for bilingual word alignment," in Proc. Conf. Human Lang. Technol. Empirical Methods Natural Lang. Process., Vancouver, BC, Canada, 2005, pp. 81–88.

[12] X. Ding, B. Liu, and P. S. Yu, "A holistic lexicon-based approach to opinion mining," in Proc. Conf. Web Search Web Data Mining, 2008, pp. 231–240.

[13] F. Li, C. Han, M. Huang, X. Zhu, Y. Xia, S. Zhang, and H. Yu, "Structure-aware review mining and summarization," in Proc. 23th Int. Conf. Comput. Linguistics, Beijing, China, 2010, pp. 653–661.

[14] Y. Wu, Q. Zhang, X. Huang, and L. Wu, "Phrase dependency parsing for opinion mining," in Proc. Conf. Empirical Methods Natural Lang. Process., Singapore, 2009, pp. 1533–1541.

[15] T. Ma and X. Wan, "Opinion target extraction in chinese news comments," in Proc. 23th Int. Conf. Comput. Linguistics, Beijing, China, 2010, pp. 782–790.