



AUTOMATIC PRODUCT ASPECT IDENTIFICATION FOR OPINION MINING

N.KIRUTHIKA¹, M.SIVA KUMAR²

¹PG Student ²Assistant Professor.

Department Of CSE, K.Ramakrishnan College Of Technology, Samayapuram, Trichy.
Kiruthika035@gmail.com, Siva.krct@gmail.com

ABSTRACT

The growth of web 2.0 application, consumer feedback about product is analyzed to improve the quality of the product. The consumer feedback or reviews are extracted from the social media and then determine the polarity (positive, negative or objective) is called sentiment analysis. It is also known as opinion mining or appraisal extraction or review mining. The sentiment lexicon plays an important role in sentiment analysis. General purpose sentiment lexicons are not suitable for aspect level or identifying particular domain topics. It is also difficult to identify the polarity of same opinion word is used for different aspects for example, in the laptop review, “large” is negative for battery aspect while being positive for the screen aspect. The need for automatic tool to extract, analyze and also understand the consumer suggestions about a product in the aspect level. The novel method is used the context dependent sentiment lexicon. The lexicon is context dependent, domain dependent and also aspect dependent. Initially, the product aspects are extracted from the consumer feedback and then determine the opinion word related to the product aspect. Both aspect and the opinion word are given to the context dependent sentiment lexicon is used to determine the polarity of the opinion word. The novel method is more accurate than other state of the art methods.

Keywords: Sentiment Analysis, Opinion Mining, Context dependent Sentiment Lexicon

I. INTRODUCTION

The explosion of web application, Sentiment analysis is used to determine the polarity (positive, negative and objective) of the each opinion word from the customer reviews about a particular product. It is also known as opinion mining. It is classified into four levels. Those are document level, feature level, entity-aspect level and sentence level. Document level, determining the overall opinion of a document. Sentence level, examining the sentiment expressed in sentences. Feature level, determining the feature from the reviews and then determining the opinions that are linked with each feature. Entity-aspect level, this is granular analysis takes into consideration each opinion expressed in the content. General sentiment lexicons are not suitable for aspect level or identifying particular domain topics. It is also difficult to identify the polarity of same opinion word is used for different aspects for example, General sentiment lexicon is considering the word “unpredictable” in negative orientation. But the word “unpredictable” is negative in the electronics domain and also positive in the

movie domain. The need for automatic tool to extract, analyze and also understand the consumer suggestions about a product in the aspect level. The novel method is used the context dependent sentiment lexicon. The context information about a product is given to the lexicon. The lexicon is context dependent, domain dependent and also aspect dependent.

In the related area of product aspect level sentiment analysis from product reviews, several research efforts have focused on the identification of product aspect and the opinion about the particular product aspect (e.g. Popescu and Etzioni (2005), Hu and Liu (2004)). Popescu and Etzioni (2005) introduced OPINE, an unsupervised information extraction system. It used to extract the product features from the reviews and determine the opinion about each feature. It focused on three subtasks for opinion mining: In first step the opinion mining used the relaxation labeling to find the semantic orientation of words for the given feature and sentences. Relaxation labeling was iterative procedure. It considered the neighborhood features. Second step to find the precision

and recall values and the third step used to determine the polarity of each features [1].

Traditional lexicon based method used to perform the following tasks. First, to express the positive or negative. If it had more positive words, the final opinion of the feature was positive otherwise negative. Xiaowen Ding et.al (2008) proposed novel Holistic Lexicon Based Method had not only considering the opinions of the current sentence it also considered the opinion of the future sentences. It used to determine the opinion about the product more accurately [3]. Andrea Esuli et.al (2006) proposed a lexical resource for opinion mining. It based on the WordNet. It assigned the score for each synset of WordNet (positive, negative or objective). It used to identify the opinion about particular product. It used for determining text SO polarity, determining text PN polarity and determining the strength of text PN-polarity. The score range from 0.0 to 1.0 and their sum 1.0 for each synset [4].

Peter D. Turney (2002), proposed a simple classification learning algorithm. Semantic orientation had estimated using PMI-IR algorithm Pointwise Mutual Information (PMI) and Information Retrieval (IR) to measure the similarity of pairs of words or phrases. The semantic orientation of a phrase calculated as the mutual information between the given phrase and the word "excellent" minus the mutual information between the given phrase and the word "poor" [12]. Mei et.al also applied the topic modeling to capture the topic or aspect from the reviews. Topic modeling used to determine or mine the topic from document set [9]. Veselin Stoyanov et.al (2008), introduced A Coreference Approach to Topic identification[13]. Vasileios Hatzivassiloglou et.al , proposed log-linear regression model used to identify and validate from large corpus constraints from conjunctions on the positive or negative semantic orientation of the conjoined adjectives [5].

Lei Zhang et.al, introduced a method to identify the opinion for the particular product features. The existing technique double propagation method used for medium size corpora. The novel method used whole pattern and part pattern to identify the feature and opinion word. The method utilized Bipartite and Hits algorithm for feature ranking [14].

Raymond Y.K. Lau et.al, proposed a method to extract, analyze and summarize consumer reviews about a products with references to the specific product context automatically. Ontology Based Opinion Mining System used to extract opinion for product aspect. In this paper used the ontology extraction and context sensitive opinion mining. KullBack-Leibler(KL) divergence algorithm used for contextual information. It used for context sensitive opinion mining. Fuzzy domain ontology captured the taxonomic information for example "iphone" (product) is a mobile phone (product category) and non taxonomic relationship such as "screen" (product feature) is "associated with" "iphone" (product). It used for the context sensitive sentiment orientation ("excellent") of a product feature ("screen") captured in the fuzzy domain

ontology. Opinion mining used the sentiment lexicon and the NLP rules for determine the polarity [7].

Yue Lu et.al, focused on the problem of sentiment lexicon in domain specific and the domain dependent and the aspect in context. It provided the optimization framework that provided the way to combine different source of information for learning such as context dependent sentiment lexicon. The context dependent sentiment lexicon is also aspect dependent. It used to determine the opinion about particular aspect more accurately [6].

II. SYSTEM ARCHITECTURE

The aspect level sentiment analysis is used to extract the product aspect or feature from the user product reviews and to determine the opinion about each product aspect. The aspect identification is based on the similarity measures. Point wise mutual information term frequency and inverse document frequency (PMI-TFIDF) is used for similarity measures. The domain vector creation is based on the association between the candidate feature and the domain entity. The domain vector is used to extract the domain specific features. It is used to improve the performance compared with the other state of the art methods.

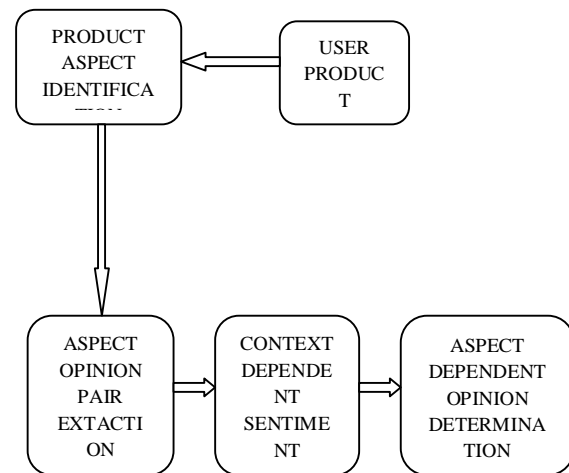


Figure 1.1 System Architecture

Figure 1.1 describes, User product reviews are used for aspect identification. In the aspect identification is used the comparative domain corpora. It is set of different product reviews. It is used to create domain vector. The domain vector is derived the association between the candidate feature and the domain entity.

In the domain vector creation is used the preprocessinBstep and the candidate features list. Preprocessing step is used the tokenizer, POS Tagger, Stemming and the stopword elimination. The term similarity measure is used to evaluate the association between candidate feature and domain entity. It is used to derive the domain vector. Domain vectors are used to extract the domain specific features.

The distance between the domain vectors is derived based on association between the candidate feature (CF) and the domain entity (DE) is used to extract the domain specific feature or aspect. This proposed method is used the PMI-TFIDF. PMI-TFIDF is a measure of information theory. It is used to measure the association between the words. In the problem of PMI for measuring the association of words. The traditional PMI method is not suitable for small corpus. In small size corpus PMI values of different words are similar values. So the PMI is combined to the TFIDF.

Opinion mining includes three subtasks: (1) Aspect opinion pair extraction, (2) Context dependent sentiment lexicon and (3) Aspect dependent opinion determination.

Set of aspects are given as the input in the aspect opinion pair extraction. It is used to extract the set of opinion words for each aspect. Context dependent sentiment lexicon is combined general sentiment lexicon and the domain dependent corpus, WordNet and language clues. It is used to extract the polarity (positive, negative or neutral) for each aspect. Aspect dependent opinion determination is used to combine the reviews about a particular product feature. It is used to determine the overall opinion about a product aspect.

III. PRODUCT ASPECT IDENTIFICATION

Aspect identification is using the comparative domain corpora. Comparative domain corpora are several product review sets. Extract domain product features through the evaluation of their weights in different related domains. Similarity measures are used to evaluate the association candidate feature and entity. It is used to derive the domain vector. Domain vector is used to determine the domain specific features or aspects.

DOMAIN VECTOR CREATION

Domain vector creation includes the preprocessing step. The preprocessing step is used to extract the aspect of the particular domain that are nouns and noun phrases, nouns and multiple noun string are extracted in this step. The candidate feature extraction includes noun phrase extraction, name entity recognition and stop word elimination.

Term similarity measure to evaluate the association of candidate features and domain words. Based on these similarities, a domain vector for each candidate feature can be derived. The similarity is evaluated by a new measure of PMI-TFIDF, which combines point wise mutual information (PMI) and term frequency-inverse document frequency (TF-IDF). A desired Domain Entity (DE) should meet two conditions: a noun and with high frequency in a domain corpus. Term similarity measure is used new measure of PMI-TFIDF which combines the pointwise mutual information (PMI) and term frequency inverse document frequency (TFIDF).

TERM FREQUENCY

Term frequency (TF) denotes the relative importance of a term to a document, which means that the number of times a word appears in a document.

Where n_{ij} is the frequency of word i in document j , Σ is the number of all words in document j .

INVERSE DOCUMENT FREQUENCY

$$IDF_i = \log \frac{|\Sigma|}{| \{j : w_i \in j\} |}$$

Where $|\Sigma|$ the frequency of documents in a corpus is, $| \{j : w_i \in j\} |$ is the number of documents containing word w_i . In the small corpus, different words to get same PMI so the TFIDF is combine to the PMI. To evaluate the association of candidate features and domain words.

POINT WISE MUTUAL INFORMATION

$$PMI(w_i, w_j) = \log_2 \frac{Pr(w_i, w_j)}{Pr(w_i) \cdot Pr(w_j)}$$

Where $Pr(w_i, w_j)$ is the probability of a sentence containing word w_i and w_j in a corpus, $Pr(w_i)$ is the probability of a sentence containing word w_i in a corpus.

PMI and TF-IDF

$$PMI - TFIDF(w_i, w_j) = \log \frac{(n_{ij} \cdot \Sigma) \cdot (\Sigma \cdot \Sigma)}{(n_{i.} \cdot \Sigma) \cdot (n_{.j} \cdot \Sigma)}$$

As domain entity (DE) is one of the candidate features (Cf) as well, and satisfies $DE \in \{Cf_1, \dots, Cf_m\}$, evaluate each candidate feature by measuring the distance between its domain vector and the domain vector of a domain entity.

$$Sim(DV_{cf}, DV_{de}) = \frac{\sum_i (DV_{cf} \cdot DV_{de})}{\sum_i DV_{cf} \cdot \sum_i DV_{de}}$$

Where DV_{cf} is the domain vector of candidate feature, DV_{de} is the domain vector of domain entity. If the value of $Sim(DV_{cf}, DV_{de})$ above the threshold value, cf would be recognized as a domain-specific feature dsf.

$$Cf \in \{ dsf \}$$

IV. OPINION MINING

It includes three subtasks. That are (1) Aspect Opinion Pair Extraction, (2) Context dependent sentiment lexicon and (3) Aspect dependent Opinion Determination.

ASPECT OPINION PAIR EXTRACTION

Aspect opinion pair extraction is used the given set of aspects of a certain product, the goal is to extract the aspect-opinion pair. It is used to identify the set of opinion words for each product. Dependency distance (dd) is used to extract aspect opinion pair. It is extract the relationship between the aspect and opinion.

$$= \frac{1}{\sum_{h=1}^n \frac{1}{h}} \left(\frac{1}{h} \right)$$

CONTEXT DEPENDENT SENTIMENT LEXICON

Context dependent is used to combine the general sentiment lexicon is used to identify the positive and negative polarity for each aspect. It is the combination of general sentiment lexicon, WordNet, language clues and the domain dependent corpora.

$$\begin{aligned} & \left(\begin{array}{c} - \\ ; \\ ; \\ ; \end{array} \right) = \\ & \begin{array}{c} \in - \\ \in - \\ h \end{array} \end{aligned}$$

ASPECT DEPENDENT OPINION DETERMINATION

It is used to determine the overall opinion about the particular aspect. It is classified the different reviews about the particular product aspect. Support vector machine is used for classification of reviews. Finally overall opinion about the product aspect is determined in the aspect dependent opinion determination.

V. CONCLUSION

The novel context dependent sentiment lexicon is used to determine the opinion in the aspect level. It is also have the context information about the product. The automatic tool is used to extract, analyze and understand the consumer feedback in the aspect level. The sentiment lexicon is domain dependent, context dependent and also aspect dependent. It is used to identify the opinion about the product aspect more accurately compare to the other state of the art methods.

REFERENCES

- [1] I. Chetviorkin, N. Loukachevitch, Extraction of domain-specific opinion words for similar domains, in: Proc. of the Workshop on Information Extraction and Knowledge Acquisition, 2011, pp. 7–12.
- [2] X. Ding, B. Liu, The utility of linguistic rules in opinion mining, in: Proc. of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2007, pp. 811–812.
- [3] X. Ding, B. Liu, P. S. Yu, A holistic lexicon based approach to opinion mining, in: Proc. of the International Conference on Web Search and Web DataMining (WSDM'08), 2008, pp. 231239.
- [4] A. Esuli, F. Sebastiani, SentiWordNet: a publicly available lexical resource for opinion mining, in: Proc. of the Fifth International Conference on Language Resources and Evaluation (LREC 2006), 2006, pp. 417–422.
- [5] V. Hatzivassiloglou, K. R. McKeown, Predicting the semantic orientation of adjectives, in: Proc. of the 40th Annual Meeting on Association for Computational Linguistics (ACL 1997), 1997, pp. 174–181.
- [6] Y. Lu, M. Castellanous, U. Zhai, Automatic Construction of a Context-Aware Sentiment Lexicon: An Optimization Approach, in: 20th International Conference on the World Wide Web (WWW), Hyderabad, India, 2011, pp. 347–356.
- [7] R.Y.K. Lau, C.C.L. Lai, J. Ma, Y.F. Li, Automatic domain ontology extraction for context-sensitive opinion mining, in: Proc. of the International Conference on Information Systems celebrates (ICIS 2009), 2009, pp. 35–53.
- [8] T. Ma, X. Wan, Opinion target extraction in Chinese news comments, in: Proc. of the International Conference on Computational Linguistics (COLING 2010), 2010, pp. 782–790.
- [9] Q. Mei, X. Ling, M. Wondra, H. Su, C. Zhai, Topic sentiment mixture: modeling facets and opinions in weblogs, in: Proc. of the World Wide Web Conference (WWW'07), 2007, pp. 171–180.
- [10] A. Popescu, O. Etzioni, Extracting product features and opinions from reviews, in: Proc. of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing (HLT-EMNLP), 2005, pp. 339–346.
- [11] D. Putthividhya, J. Hu, Bootstrapped named entity recognition for product attribute extraction, in: Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP), 2011, pp. 1557–1567.
- [12] P. Turney, Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews, in: Proc. of the 40th Annual Meeting on Association for Computational Linguistics (ACL 2002), 2002, pp. 417–424.
- [13] S. Veselin, C. Cardie, Topic identification for fine-grained opinion analysis, in: Proc. of the International Conference on Computational Linguistics (COLING 2008), 2008, pp. 817–824.
- [14] L. Zhang, B. Liu, Extracting and ranking product features in opinion documents, in: Proc. of the International Conference on Computational Linguistics (COLING 2010), 2010, pp. 1462–1470.