

A Survey : Extracting Features from Online Reviews Corpus by Domain Relevance

Ashwini P. Lingojar¹, L. J. Sankpal²

¹Pune University, SAE, Kondhwa, Pune, Maharashtra, India

²Assistant Professor, SAE, Kondhwa, Pune, Maharashtra, India

Abstract: Large amount of user generated data is present on web in the form of blogs, reviews tweets, comments etc. This type of data involves user's opinion, view, attitude, sentiment towards particular product, topic, event, news etc. Capturing public opinion about social events and product preferences is increasing interest from the customer and from the business world. Nowadays, if one wants to buy a product, consumer is no longer limited to asking friends and family for opinions because there are many user reviews and discussions in public forums on the Web about the product. But consumers are not satisfied with overall reviews of the product where as they want to understand which is positive and negative attribute of the product. Opinions and sentiments are expressed in the text reviews. Most of the opinion features are extracted from the online review corpus by using mining patterns. As Features are extracted from single review corpus it ignoring the important difference in word distributional characteristics of opinion features across different corpora. As a result efficient opinion features are not extract. In order to obtain valid feature it is necessary to consider two different online review according to the domain relevance.

Keyword: Opinion mining, natural language processing, opinion feature, information search and retrieval

1. Introduction

There were millions of internet users are present, as a result most of the social media has accumulated massive amount of valuable peer review on almost everything. With the rapid growth of e-commerce, more products are sold on the Web as well as more people are buying products on the Web. In order to enhance customer satisfaction and their shopping experiences, it has become a common practice for online merchants to enable their customers to review or to express opinions on the products that they buy. As a result of this, the product present on the web receives large number of the customer reviews. Most of the time the reviews are in the textual form. This makes it very hard for a potential customer to read them and to make a decision on whether to buy the product. In order to make easy purchase decision, Mining this pool of review and detect opinion feature has become useful[1].

Analysis of opinions from the online review corpus is known as opinion mining or sentiment analysis. Opinion mining is also known as sentiment analysis which is used to analyze people's opinions, sentiments, attitudes and emotions towards entities such as products, services, organizations and their attribute[2]. In opinion mining, opinion feature indicates an attribute or entity on which user express their opinion. In opinion mining the opinion that are express in textual form are analysed at various level such as document level opinion mining and sentence level opinion mining. As, the reviews that are express in sentence level can derive better opinion as compare to the document level. Nowadays customer are not satisfied with the overall rating of the product but they want to know the positive and negative attribute of the product. So it is very important to extract valid opinion feature from the text reviews and associate them to opinion.

In order to extract opinion feature in opinion mining various approach have been proposed. Supervised learning model, there is problem related to the domain adaptation. It work good in given domain but not in different domain if it is applied to [3][4]. Unsupervised natural language processing approach [5][6][7], do not work well on colloquial real life review due to the lack of the formal structure. As it identify the feature by defining the syntactic rule. Topic modelling approach can mine the aspects of specific feature that are commented on explicitly in review [8][9]. Corpus statistic approaches try to extract opinion feature only from the given review corpus without considering their distributional characteristics another different corpus[1].

As the structure of opinion feature is distributional in a given domain independent and domain depended corpus. So here the survey is based on the approaches that identify the opinion feature by considering their distributional characteristics across different domain. And this approach is an IDER (Intrinsic and Extrinsic domain relevance) approach that evaluated the domain relevance of an opinion feature across two corpora and extract valid opinion feature.

2. Related Work

Extracting opinion feature is the subproblem of opinion mining. Large amount of work is done in the product review domain. There are various approaches are present in order to extract opinion feature which are mainly classified into two categories, namely, Supervised and Unsupervised.

Supervised learning models are including hidden Markov models. Supervised models perform well on a given domain, but it need to training again when it applied to a different domain. For a wide variety of products and services in different domains, supervised methods are not efficient because it is very expensive to construct labelled data for

each product or service. In addition, this model requires a decent-sized set of labelled data for model learning on every domain [11]. Supervised learning models that require labelled data have been successfully used to build sentiment classifiers for a given domain. In supervised learning model there is problem in domain adaption [10].

Unsupervised approach can be applied to any domain because it does not require any predefined aspects or sentiment lexicons. In an unsupervised approach, Syntactic relationships between features and opinions can be used to locate opinion features in a sentence by using carefully generated syntactic rules[6]. Syntactic relations which are identified by this methods help to locate features associated with opinion words. Unsupervised NLP approaches extract opinion features by mining syntactic patterns of features present in review sentences[5]. An unsupervised approach to capturing sentiment-oriented aspects for online reviews is challenging, since the reviews are short and usually each aspect is mentioned only once in a review. This approach could extract large number of invalid features due to the colloquial nature of online reviews.

Unsupervised corpus statistic approach resists the colloquial nature of son line review in order give a large review corpus to extract valid feature. To understand the distributional characteristic of the opinion feature this approach use the results of statistical analysis. Hu and Liu [1] proposed an association rule mining (ARM) approach to mine frequent itemsets and generate opinion features. These opinion features are nouns and noun phrases with high sentence-level frequency. As this ARM depends on the frequency of itemsets, it has some limitations for the task of feature identification, i.e., frequent but invalid features are extracted incorrectly and rare but valid features maybe overlooked.

To address the problem of feature-based opinion mining, Suet al. [14] introduced a mutual reinforcement clustering (MRC) approach to mine the associations between product feature categories and opinion word groups. MRC approach utilizes multisource knowledge including semantic and textual structure. Based on a co-occurrence weight matrix generated from the given review corpus. The mutual reinforcement clustering approach fully exploit the relationship between product features and opinion words. Unlike many other corpus statistics methods, MRC approach is able to extract infrequent features. The infrequent features are extracted only when, there are the mutual relationships between feature and opinion groups are found during the clustering phase is accurate. MRC's precision is low due to the difficulty in obtaining good clusters on real-life reviews.

Yu et al. [13] proposed an aspect ranking algorithm which is based on the probabilistic regression model. This ranking algorithm is use to automatically identify important product aspects from online consumer reviews. Aspect ranking algorithm identify the important aspects by simultaneously taking the aspect frequency and the influence of consumers 'opinions given to each aspect on their overall opinions. Moreover, this algorithm is not focus on extracting feature terms that are commented on explicitly in reviews,

but on ranking product aspects that are actually coarse-grained clusters of specific features.

Another approach of unsupervised topic modelling, such as latent Dirichlet allocation (LDA) [8], approach does not consider the relationships among sentences, thus ignoring that the same aspect may have quite different word usages in different sentences. This approach is a generative three-way (term-topic-document) probabilistic model that have been used to solve aspect-based opinion mining tasks. The models which are developed primarily for mining latent topics or aspects, are actually use to identify coarse-grained topics or aspects that correspond to distinguishing properties or concepts of the commented entities. This approach may not expressed opinion feature explicitly in the reviews, but rather user-defined clusters of specific opinion features.[8], [9], [12], [10]. The approaches are effective in discovering latent structures of review data but they may be less successful in dealing with identifying specific feature terms commented on explicitly in reviews. Most of these existing approaches to feature extraction typically only use the knowledge or patterns mined from a given single review corpus, while it completely ignoring the possible variations present in a different domain independent corpus.

Another IEDR approach utilizes the fact that word distribution characteristics vary across two different types of corpora, in particularly domain-specific versus domain-independent. This approach is use to derive powerful hints that help to discriminate valid features from the invalid ones. This approach derive candidate feature according to domain relevance with respect to intrinsic domain relevance score and extrinsic domain relevance score. In this the first step is to extract candidate features which is similar to NLP approach. To extract candidate feature define syntactic dependence rule. Then the next step is, for each candidate feature identify its domain relevance score with respect to domain specific and domain independent corpora. In the final step candidate with low IDR core and high EDR scores are pruned. The main difference of IEDR compared to all other existing methods is that this approach contains the smart fusion of domain-dependent and domain-independent information sources.

3. Conclusion

In this paper we have presented a literature survey on Extracting features from online reviews corpus. Domain relevance IDER approaches mainly focus on extracting valid candidate feature from online review corpus. Various features extraction approaches are listed this survey.

References

- [1] M.Hu and B. Liu, "Mining and Summarizing Customer Reviews," Proc. 10th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 168-177, 2004
- [2] B. Liu, "Sentiment Analysis and Opinion Mining," Synthesis Lectures on Human Language Technologies, vol. 5, no. 1, pp. 1-167, May 2012.

- [3] W. Jin and H.H. Ho, "A Novel Lexicalized HMM-Based Learning Framework for Web Opinion Mining," Proc. 26th Ann. Int'l Conf. Machine Learning, pp. 465-472, 2009.
- [4] N. Jakob and I. Gurevych, "Extracting Opinion Targets in a Single and Cross-Domain Setting with Conditional Random Fields," Proc. Conf. Empirical Methods in Natural Language Processing, pp. 1035-1045, 2010.
- [5] S.-M. Kim and E. Hovy, "Extracting Opinions, Opinion Holders, and Topics Expressed in Online News Media Text," Proc. ACL/COLING Workshop Sentiment and Subjectivity in Text, 2006.
- [6] G. Qiu, C. Wang, J. Bu, K. Liu, and C. Chen, "Incorporate the Syntactic Knowledge in Opinion Mining in User-Generated Content," Proc. WWW 2008 Workshop NLP Challenges in the Information Explosion Era, 2008.
- [7] G. Qiu, B. Liu, J. Bu, and C. Chen, "Opinion Word Expansion and Target Extraction through Double Propagation," Computational Linguistics, vol. 37, pp. 9-27, 2011.
- [8] D.M. Blei, A.Y. Ng, and M.I. Jordan, "Latent Dirichlet Allocation," J. Machine Learning Research, vol. 3, pp. 993-1022, Mar. 2003.
- [9] I. Titov and R. McDonald, "Modelling Online Reviews with Multi-Grain Topic Models," Proc. 17th Int'l Conf. World Wide Web, pp. 111-120, 2008.
- [10] Y. Jo and A.H. Oh, "Aspect and Sentiment Unification Model for Online Review Analysis," Proc. Fourth ACM Int'l Conf. Web Search and Data Mining, pp. 815-824, 2011.
- [11] F. Li, C. Han, M. Huang, X. Zhu, Y.-J. Xia, S. Zhang, and H. Yu, "Structure-Aware Review Mining and Summarization," Proc. 23rd Int'l Conf. Computational Linguistics, pp. 653-661, 2010.
- [12] J. Yu, Z.-J. Zha, M. Wang, and T.-S. Chua, "Aspect Ranking: Identifying Important Product Aspects from Online Consumer Reviews," Proc. 49th Ann. Meeting of the Assoc. for Computational Linguistics: Human Language Technologies, pp. 1496-1505, 2011.
- [13] W.X. Zhao, J. Jiang, H. Yan, and X. Li, "Jointly Modeling Aspects and Opinions with a Maxent-Lda Hybrid," Proc. Conf. Empirical Methods in Natural Language Processing, pp. 56-65, 2010.
- [14] Q. Su, X. Xu, H. Guo, Z. Guo, X. Wu, X. Zhang, B. Swen, and Z. Su, "Hidden Sentiment Association in Chinese Web Opinion Mining," Proc. 17th Int'l Conf. World Wide Web, pp. 959-968, 2008.
- [15] M. Hu and B. Liu, "Mining and Summarizing Customer Reviews," Proc. 10th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining, pp. 342-351, 2004.

Ms. L.J. Sankpal is working as Assistant Professor in Computer Engineering Department, SAE, Kondhwa, Pune, Pune University. India.

Author Profile

Ms. Ashwini P. Lingojwar is Student of ME Computer Engineering, SAE, Kondhwa, Pune, Pune University. India.