

# Sentiment Classification based Product Reviews and its Application

H.Sultana parveen

*M.SC(I.T.), M.Phil. Scholar, Department of Computer Science, Mother Teresa Women's University Chennai, India*

Dr. G. Rasitha Banu

*MCA., M.Phil., Ph.D.*

*Assistant Professor, Department of Health Information Management and Technology,  
Faculty of Public Health and Tropical Medicine, \*Jazan University, \*KSA*

**Abstract** - Generally, a product may have hundreds of aspects. Consumer reviews contain rich and valuable knowledge for both firms and users. However, the reviews are often disorganized, leading to difficulties in information navigation and knowledge acquisition. This article proposes a product aspect ranking framework, which automatically identifies the important aspects of products from online consumer reviews, aiming at improving the usability of the numerous reviews. The important product aspects are identified based on two observations: (a) the important aspects are usually commented by a large number of consumers; and (b) consumer opinions on the important aspects greatly influence their overall opinions on the product. In particular, given the consumer reviews of a product, we first identify product aspects by a shallow dependency parser and determine consumer opinions on these aspects via a sentiment classifier. We then develop a probabilistic aspect ranking algorithm to infer the importance of aspects by simultaneously considering aspect frequency and the influence of consumer opinions given to each aspect over their overall opinions.

**Keywords** -- Probabilistic aspect ranking algorithm, Sentiment Classification, Product Aspects, Aspect Identification,, Consumer review, Aspect ranking, Extractive review.

## I. INTRODUCTION

We proposed a framework to analyze and report customers overview comments into particular part of the customized product. A product aspect ranking framework to automatically identify the important aspects of products from numerous consumer reviews. Years have witnessed the rapidly expanding e-commerce. A recent study from ComScore reports that online retail spending reached \$37.5 billion in Q2 2011 U.S. Millions of products from various merchants have been offered online. For example, Bing Shopping1 has indexed more than five million products. Amazon.com archives a total of more than 36 million products. Shopper.com records more than five million products from over 3,000 merchants. Most retail Websites encourage consumers to write reviews to express their opinions on various aspects of the products.

Here, an aspect, also called feature in literatures, refers to a component or an attribute of a certain product. A sample review "The battery life of Nokia N95 is amazing." reveals positive opinion on the aspect "battery life" of product Nokia N95. Besides the retail Websites, many forum Websites also provide a platform for consumers to post reviews on millions of products.

For example, CNet.com involves more than seven million product reviews; whereas Pricegrabber.com contains millions of reviews on more than 32 million products in 20 distinct categories over 11,000 merchants. Such numerous consumer reviews contain rich and valuable knowledge and have become an important resource for both consumers and firms.

Consumers commonly seek quality information from online reviews prior to purchasing a product, while many firms use online reviews as important feedbacks in their product development, marketing, and consumer relationship management. Generally, a product may have hundreds of aspects. For example, iPhone 3GS has more than three hundred aspects (see Fig. 1), such as "usability," "design," "application," "3G network." We argue that some aspects are more important than the others, and have greater impact on the eventual consumers' decision making as well as firms' product development strategies. For example, some aspects of iPhone 3GS, e.g., "usability" and "battery," are concerned by most consumers, and are more important than the others such as "usb" and "button."

## II. PROPOSED ALGORITHM

**PROBABILISTIC ASPECT RANKING ALGORITHM**

A probabilistic aspect ranking algorithm to identify the important aspects of a product from consumer reviews. Generally, important aspects have the following characteristics: (a) they are frequently commented in consumer reviews; and (b) consumers opinions on these aspects greatly influence their overall opinions on the product. The overall opinion in a review is an aggregation of the opinions given to specific aspects in the review, and various aspects have different contributions in the aggregation. That is, the opinions on (un)important aspects have strong (weak) impacts on the generation of overall opinion.

---

Probabilistic Aspect Ranking Algorithm :

---

**Input:** Consumer review corpus  $R$ , each review  $r \in R$  is associated with an overall rating  $Or$ , and a vector of opinions  $or$  on specific aspects.

**Output:** Importance scores  $_k/mk=1$  for all the  $m$  aspects.

**while not converged do**

Update  $\{or\}/R/r=1$  according to Eq. (9);

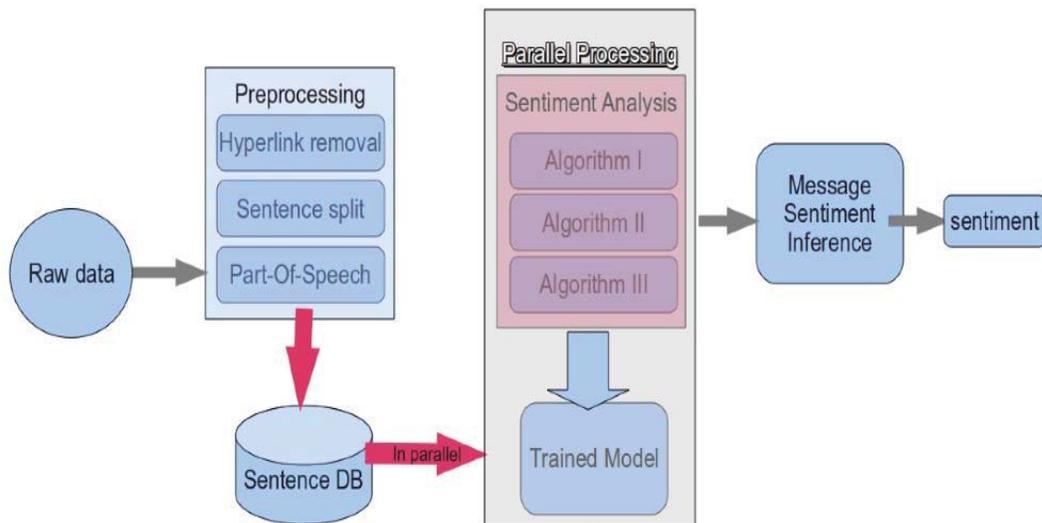
Update  $\{\mu, \Sigma, \sigma^2\}$  according to Eq. (13);

**end while**

Compute aspect importance scores  $\{_k\}mk=1$

---

## B. SENTIMENT CLASSIFICATION



In the general review extracted from the aspect identification have to classify the review whether which is shared for the good rating or bad rating. Ex the word headphone and earphone targets about the same product but in different words. Thus the review might be containing good rating without exclamation words are differentiated. The task of analyzing the sentiments expressed on aspects is called aspect-level sentiment classification in literature. Existing techniques include the supervised learning approaches and the lexicon-based approaches, which are typically unsupervised. The lexicon-based methods utilize a sentiment lexicon consisting of a list of sentiment words, phrases and idioms, to determine the sentiment orientation on each aspect.

While these methods are easily to implement, their performance relies heavily on the quality of the sentiment lexicon. On the other hand, the supervised learning methods train a sentiment classifier based on training corpus. The classifier is then used to predict the sentiment on each aspect. Many learning-based classification models are applicable, for example, Support Vector Machine (SVM), Naive Bays, and Maximum Entropy (ME) model etc. Supervised learning is dependent on the training data and cannot perform well without sufficient training samples. However, labeling training data is labor-intensive and time-consuming. In this work, the Pros and Cons reviews

have explicitly categorized positive and negative opinions on the aspects. These reviews are valuable training samples for learning a sentiment classifier. We thus exploit Pros and Cons reviews to train a sentiment classifier, which is in turn used to determine consumer opinions (positive or negative) on the aspects in free text reviews. Specifically, we first collect the sentiment terms in Pros and Cons reviews based on the sentiment lexicon provided by MPQA project. These terms are used as features, and each review is represented as a feature vector. A sentiment classifier is then learned from the Pros reviews (i.e., positive samples) and Cons reviews (i.e., negative samples).

Given a free text review that may cover multiple aspects, we first locate the opinionated expression that modifies the corresponding aspect, e.g. locating the expression “well” in the review “The battery of Nokia N95 works well.” for the aspect “battery.”

Generally, an opinionated expression is associated with the aspect if it contains at least one sentiment term in the sentiment lexicon, and it is the closest one to the aspect in the parsing tree within the context distance of 5. The learned sentiment classifier is then leveraged to determine the opinion of the opinionated expression, i.e. the opinion on the aspect.

### C. PRODUCT ASPECT



Software product line engineering aims to reduce development time, effort, cost, and complexity by taking advantage of the commonality within a portfolio of similar products. The effectiveness of a software product line approach directly depends on how well feature variability within the portfolio is implemented and managed throughout the development lifecycle, from early analysis through maintenance and evolution. This paper presents an approach that facilitates variability implementation, management and tracing by integrating model-driven and aspect-oriented software development. Features are separated in models and composed by aspect-oriented composition techniques on model level. Model transformations support the transition from problem to solution domain. Aspect-oriented techniques enable the explicit expression and modularization of variability on model, code, and template level. The presented concepts are illustrated with a case study of a home automation system.

### D. ASPECT IDENTIFICATION

With the general review have to find out the product name or parts of the product name executed in the first phase. Ex in mobile phone product the review is specially for the camera or the battery is identified in the first phase. The Websites such as CNet.com require consumers to give an overall rating on the product, describe concise positive and negative opinions (i.e. Pros and Cons) on some product aspects, as well as write a paragraph of detailed review in free text. Some Websites, e.g., Viewpoints.com, only ask for an overall rating and a paragraph of free-text review. The others such as Reevo.com just require an overall rating and some concise positive and negative opinions on certain aspects. In summary, besides an overall rating, a consumer review consists of Pros and Cons reviews, free text review, or both.

For the Pros and Cons reviews, we identify the aspects by extracting the frequent noun terms in the reviews. Previous studies have shown that aspects are usually nouns or noun phrases, and we can obtain highly accurate aspects by extracting frequent noun terms from the Pros and Cons reviews. For identifying aspects in the free text reviews, a straightforward solution is to employ an existing aspect identification approach. It first identifies the nouns and noun phrases in the documents. The occurrence frequencies of the nouns and noun phrases are counted, and only the frequent ones are kept as aspects.

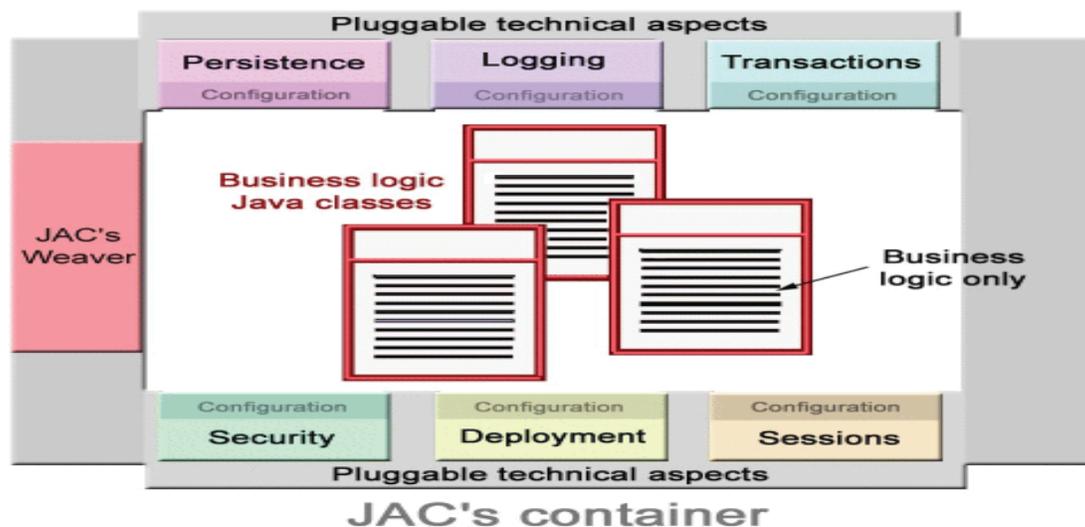
Although this simple method is effective in some cases, its well-known limitation is that the identified aspects usually contain noises. Recently, a phrase dependency parser to extract noun phrases, which form candidate aspects. To filter out the noises, they used a language model by an intuition that the more likely a candidate to be an aspect, the more closely it related to the reviews. The language model was built on product reviews, and used to predict the related scores of the candidate aspects.

The candidates with low scores were then filtered out. However, such language model might be biased to the frequent terms in the reviews and cannot precisely sense the related scores of the aspect terms, as a result cannot filter out the noises effectively. In order to obtain more precise identification of aspects, we here propose to exploit the Pros and Cons reviews as auxiliary knowledge to assist identifies aspects in the free text reviews. In particular, we first split the free text reviews into sentences, and parse each sentence using Stanford parser2. The frequent noun phrases are then extracted from the sentence parsing trees as candidate aspects. Since these candidates may contain noises, we further leverage the Pros and Cons reviews to assist identify aspects from the candidates.

We collect all the frequent noun terms extracted from the Pros and Cons reviews to form a vocabulary. We then represent each aspect in the Pros and Cons reviews into a unigram feature, and utilize all the aspects to learn a one-class Support Vector Machine (SVM) classifier. The resultant classifier is in turn used to identify aspects in the candidates extracted from the free text reviews. As the identified aspects may contain some synonym terms, such as “earphone” and “headphone,” we perform synonym clustering to obtain unique aspects. In particular, we collect the synonym terms of the aspects as features.

The synonym terms are collected from the synonym dictionary. We represent each aspect into a feature vector and use the Cosine similarity for clustering. The ISODATA (Iterative Self-Organizing Data Analysis Technique) clustering algorithm is employed for synonym clustering. ISODATA does not need to fix the number of clusters and can learn the number automatically from the data distribution. It iteratively refines clustering by splitting and merging of clusters. Clusters are merged if the centers of two clusters are closer than a certain threshold.

One cluster is split into two different clusters if the cluster standard deviation exceeds a predefined threshold. The values of these two thresholds were empirically set to 0.2 and 0.4 in our experiments.



#### E. CONSUMER REVIEW

In a product review website users commonly share their thoughts about a product used. That review generally contains both good and bad about the things. If new user search about a product review it displays in a same form. Hence it's difficult to find out the original issue with the product. For ex a mobile have good camera resolution but worst battery features, means finding the reality is difficult in the existing system. Hence, identifying important product aspects will improve the usability of numerous reviews and is beneficial to both consumers and firms. Consumers can conveniently make wise purchasing decision by paying more attentions to the important aspects,

while firms can focus on improving the quality of these aspects and thus enhance product reputation effectively. However, it is impractical for people to manually identify the important aspects of products from numerous reviews. Therefore, an approach to automatically identify the important aspects is highly demanded

#### F. ASPECT RANKING

With the help of both aspect identification and the sentiment word classification the product overall and individual parts review is classified and ranked separately. This framework may differ from one website to other website. In this section, we propose a probabilistic aspect ranking algorithm to identify the important aspects of a product from consumer reviews. Generally, important aspects have the following characteristics:

- (a) They are frequently commented in consumer reviews; and
- (b) Consumers opinions on these aspects greatly influence their overall opinions on the product. The overall opinion in a review is an aggregation of the opinions given to specific aspects in the review, and various aspects have different contributions in the aggregation.

That is, the opinions on (UN) important aspects have strong (weak) impacts on the generation of overall opinion.

#### G. EXTRACTIVE REVIEW

Our method aims to find what customers like and dislike about a given product. However, due to the difficulty of natural language understanding, some types of sentences are hard to deal with. The next sentences were taken from the reviews of a mobile phone. The first two can be considered easy sentences and the last sentence a hard one to handle with: "It has a nice color screen." "T-mobile was a pretty good server." "When you put this phone in your pocket you forget it is there ; it is unbelievably small and oh , so light." In the first two sentences, it is easy to note that the user is talking about color screen and T-mobile server respectively because these aspects are explicitly mentioned. However, some aspects are implicit and hard to find, like in the third sentence, where the customer is talking about size and weight. Semantic understanding is needed to find these implicit aspects, but this is out of the purpose of this paper. This work is focused on finding explicit aspects. In general, most product aspects indicating words are nouns or noun phrases. Therefore, after parsing the sentence, the next step is to identify noun phrases as potential product aspects. In this sense, we apply the linguistic filtering patterns shown in Table 1. Each pattern is defined in terms of an extended regular expressions over the POSTagging labels: JJ (adjective), NN (common noun), NNP (proper noun), VBG (gerund verb), VBN (past participle verb), and DT (general determiner). These definitions allow the extraction of both simple and compound noun phrases as potential aspects.

Name	Pattern	Examples
NP1	(JJ NN NNP)+	battery life lcd screen
NP2	NP1 (V BG V BN) NP1	battery charging system

### III.COMPARISON WITH EXISTING SYSTEM

Document-level sentiment classification aims to classify an opinion document as expressing a positive or negative opinion. Existing works use unsupervised, supervised or semi-supervised learning techniques to build document level sentiment classifiers. Unsupervised method usually relies on a sentiment lexicon containing a collection of positive and negative sentiment words. It determines the overall opinion of a review document based on the number of positive and negative terms in the review. Supervised method applies existing supervised learning models, such as SVM and Maximum entropy (ME) etc., while semi supervised approach exploits abundant unlabeled reviews together with labeled reviews to improve classification performance. The other related topic is extractive review summarization, which aims to condense the source reviews into a shorter version preserving its information content and overall meaning. Extractive summarization method forms the summary using the most informative sentences and paragraphs etc. selected from the original reviews. The most informative content generally refers to the "most frequent" or the "most favorably positioned" content in exiting works.

The two widely used methods are the sentence ranking and graph-based methods. In these works, a scoring function was first defined to compute the informativeness of each sentence. Sentence ranking method ranked the sentences according to their informativeness scores and then selected the top ranked sentences to form a summary. Graph-based method represented the sentences in a graph, where each node corresponds to a sentence and each edge characterizes

the relation between two sentences. A random walk was then performed over the graph to discover the most informative sentences, which were in turn used to compose a summary.

#### IV. ADVANTAGES OF PROPOSED SYSTEM

- a. We propose a product aspect ranking framework to automatically identify the important aspects of products from numerous consumer reviews.
- b. We develop a probabilistic aspect ranking algorithm to infer the importance of various aspects by simultaneously exploiting aspect frequency and the influence of consumers' opinions given to each aspect over their overall opinions on the product.
- c. We demonstrate the potential of aspect ranking in real world applications. Significant performance improvements are obtained on the applications of document level sentiment classification and extractive review summarization by making use of aspect ranking.
- d. We perform extensive experiments to evaluate the efficacy of aspect ranking in these two applications and achieve significant performance improvements.
- e. The occurrence frequencies of the nouns and noun phrases are counted, and only the frequent ones are kept as aspects.

#### V. CONCLUSION

We have proposed a product aspect ranking framework to identify the important aspects of products from numerous consumer reviews. The framework contains three main components, i.e., product aspect identification, aspect sentiment classification, and aspect ranking. First, we exploited the Pros and Cons reviews to improve aspect identification and sentiment classification on free-text reviews. We then developed a probabilistic aspect ranking algorithm to infer the importance of various aspects of a product from numerous reviews. The algorithm simultaneously explores aspect frequency and the influence of consumer opinions given to each aspect over the overall opinions. The product aspects are finally ranked according to their importance scores. We have conducted extensive experiments to systematically evaluate the proposed framework. The experimental corpus contains 94,560 consumer reviews of 21 popular products in eight domains. This corpus is publicly available by request. Experimental results have demonstrated the effectiveness of the proposed approaches. Moreover, we applied product aspect ranking to facilitate two real-world applications, i.e., document level sentiment classification and extractive review summarization. Significant performance improvements have been obtained with the help of product aspect ranking.

#### VI. FUTURE ENHANCEMENT

The project designed with reevoo([www.reevoo.com/shopping](http://www.reevoo.com/shopping)) website comments for a particular product. In future the framework can be implemented for a particular product with a help of multiple website comments. Hence the overall product aspect ratio in the open market can be analyzed.

#### REFERENCES

- [1] J. C. Bezdek and R. J. Hathaway.: Convergence of alternating optimization. in *Journal of Neural, Parallel & Scientific Computations*, vol. 11, pp. 351-368. USA. 2003.
- [2] C. C. Chang and C. J. Lin.: Libsvm: a Library for Support Vector Machines. <http://www.csie.ntu.edu.tw/~simjlin/libsvm/>, 2004.
- [3] G. Carenini, R. T. Ng, and E. Zwart.: Multi-document Summarization of Evaluative Text. in *Proc. of ACL*, pp. 3-7. Sydney, Australia. 2006.
- [4] China Unicom 100 Customers iPhone User Feedback Report, 2009.
- [5] ComScore Reports [http://www.comscore.com/Press Events/Press Releases](http://www.comscore.com/Press%20Events/Press%20Releases), 2011.
- [6] X. Ding, B. Liu, and P. S. Yu.: A Holistic Lexicon-based Approach to Opinion Mining. in *Proc. of WSDM*, pp. 231-240. USA. 2008.
- [7] G. Erkan and D. R. Radev.: LexRank: Graph-based Lexical Centrality as Saliency in Text Summarization. in *Journal of Artificial Intelligence Research*, vol. 22, pp. 457-479. 2004.
- [8] O. Etzioni, M. Cafarella, D. Downey, A. Popescu, T. Shaked, S. Soderland, D. Weld, and A. Yates.: Unsupervised Named-entity Extraction from the Web: An Experimental Study. in *Journal of Artificial Intelligence*, vol. 165, pp. 91-134. 2005.
- [9] A. Ghose and P. G. Ipeirotis.: Estimating the Helpfulness and Economic Impact of Product Reviews: Mining Text and Review Characteristics. in *IEEE Trans. on Knowledge and Data Engineering*, vol. 23, pp. 1498-1512. 2010.
- [10] V. Gupta and G. S. Lehal.: A Survey of Text Summarization Extractive Techniques. in *Journal of Emerging Technologies in Web Intelligence*, vol. 2, pp. 258-268. 2010.