



Mining Opinion Features in Customer Reviews

¹Havinash P.H, ²Jeril Johnson N, ³Glen Thomas, ⁴Emily Stephen P

Department of Computer Science and Engineering
Sahrdaya college of Engineering and Technology
Thrissur, Kerala, India

[¹havinashh@gmail.com](mailto:havinashh@gmail.com), [²jeril95.n@gmail.com](mailto:jeril95.n@gmail.com), [³glenkt2@gmail.com](mailto:glenkt2@gmail.com), [⁴emilystephenp@gmail.com](mailto:emilystephenp@gmail.com)

Abstract:-Now days, E-commerce systems have become extremely important. Large numbers of customers are choosing online shopping because of its convenience, reliability, and cost. Client generated information and especially item reviews are significant sources of data for consumers to make informed buy choices and for makers to keep track of customer's opinions. It is difficult for customers to make purchasing decisions based on only pictures and short product descriptions. On the other hand, mining product reviews has become a hot research topic and prior researches are mostly based on pre-specified product features to analyse the opinions. Natural Language Processing (NLP) techniques such as NLTK for Python can be applied to raw customer reviews and keywords can be extracted. This paper presents a survey on the techniques used for designing software to mine opinion features in reviews. Eleven IEEE papers are selected and a comparison is made between them. These papers are representative of the significant improvements in opinion mining in the past decade.

Keywords:-Python, POS, CRF, NLTK, LBSN

1. INTRODUCTION

Online product reviews currently have substantial effect on acquiring choice of potential clients. However, the overpowering number of those audits upset individuals to discover helpful data and use sound judgment on the buys. Subsequently; sentiment mining has turned into a theme that got much consideration as of late. Sentiment classification, which aims at classifying sentiment data into polarity categories (e.g., positive or negative), is widely studied because many users do not explicitly indicate their sentiment polarity [1]. Sentiment analysis or Opinion mining uses natural language processing. It analyses text and computational linguistics to identify and extract subjective information in source materials [3].

The essence of performing such classification is the recognition of sentiment carrying words in a sentence. One method of clustering similar sentiment carriers is through the use of Word Net. Semantic similarity between nouns and verbs can be easily evaluated from Word Net utilizing the ordered relations. The early data extraction focused on adjective and adverb, and then stretched out to verb and nouns. The sentiment classifiers frequently expect that every document has stand out subject, and the point of every report is known. The burden is that these suspicions are dependably not the case, particularly for web document [6].

Content mining is an interdisciplinary field which draws on data recovery, information mining, machine learning,

Insights and computational etymology. POI recommendation is one of the important applications of Location Based Social Network (LBSN) which brings great uses to both customers and merchants. Existing point-of-interests (POIs) recommendation systems mainly uses collaborative filtering (CF), which only exploits user given ratings about a merchant regardless of the user preference across multiple aspects, which exists in real scenarios. MARS (Multi-Aspect Recommender System) is a novel POI recommender system based on multi-aspect user preference learning from reviews by using utility theory [2].

Another approach is Graph-based opinion entity ranking framework which is used to mine opinion data from former customers and rank their entities or aspects in accordance with those opinions. Aspects and sentiment words for each entity is extracted. Relationships between reviewers and pairs of entity aspect are then represented by a weighted bipartite graph and an algorithm to compute the ranking scores is proposed [8]. Genetic algorithm also can to analyse the opinion sentences and determines the orientations of the opinions, provide an exact summary to the user and identify the polarity of the text. It is an evolutionary algorithm which is used in different problem domains from astronomy to sports, medical science and optimization of different computer science applications [3].

This paper proposes a number of techniques based on data mining and natural language processing methods to mine opinion/product features.

2. BACKGROUND

A. Natural language processing

The term Natural Language Processing envelops a wide arrangement of systems for computerized era, control and examination of normal or human languages. Although most NLP strategies acquire to a great extent from Linguistics and Artificial Intelligence, they are additionally affected by generally more up to date zones, for example, Machine Learning, Computational Statistics and Cognitive Science. Mallet is a package mainly used for NLP. Mallet finds applications in a number areas including document classification, topic modeling etc [4].

B. Basic Terminologies:

Tokenization: The way of splitting a sentence into its constituent tokens. For sectioned languages, for example, English, the presence of whitespace makes tokenization generally less demanding.

Corpus: A body of text, typically containing an expansive number of sentences.

Part-of-speech (POS) Tag: A word can be characterized into one or more of an arrangement of lexical or part-of-speech categories, for example, Nouns, Verbs, Descriptors and Articles. A POS tag is a symbol representing such a lexical classification - NN(Noun), VB(Verb), JJ(Adjective), AT(Article). One of the oldest and most ordinarily utilized label sets is the Brown Corpus label set [5][4][9].

Parse Tree: A tree characterized over a given sentence that speaks to the syntactic structure of the sentence as defined by a formal grammar.

C. Some common NLP tasks:

POS Tagging: Given a sentence and an arrangement of POS labels, a typical language processing task is to consequently assign POS labels to each word in the sentences. For instance, given the sentence The ball is red, the yield of a POS tagger would be The/AT ball/NN is/VB red/JJ. Best in class POS taggers can accomplish exactness as high as 96%. Labeling content with parts-of-speech ends up being greatly valuable for more complicated NLP errands, for example, parsing and machine interpretation[5][4][9].

Computational Morphology: Natural languages consist of a very large number of words that are built upon basic building blocks known as morphemes (or stems), the smallest linguistic units possessing meaning. Computational morphology is concerned with the discovery and analysis of the internal structure of words using computers.

Parsing: In the parsing task, a parser develops the parse tree given a sentence. A few parsers accept the presence of an arrangement of language structure rules so as to parse however late parsers are sufficiently keen to conclude the parse trees straightforwardly from the given information utilizing complex measurable models . Most parsers likewise work in a managed setting and require the sentence to be POS-labeled before it can be parsed. Statistical parsing is a range of dynamic exploration in NLP.

Machine Translation (MT): In machine translation, the goal is to have the computer translate the given text in one natural language to fluent text in another language without any human in the loop. This is one of the most difficult tasks in NLP and has been tackled in a lot of different ways over the years. Almost all MT approaches use POS tagging and parsing as preliminary steps.

D. Python

The Python programming language is a powerfully written, object-oriented interpreted dialect. Despite the fact that, its essential quality lies in the simplicity with which it permits a software engineer to quickly model a venture, its effective furthermore, develop set of standard libraries make it an extraordinary fit for extensive scale generation level programming designing undertakings also. Python has an exceptionally shallow expectation to absorb information and a online internet learning asset.

E. Natural Language Toolkit (NLTK)

In order to deal with and manipulate the text, particular toolkits are expected to sort out the content into sentences then split them into words, to encourage semantic and meaning extraction. One of these toolkit is the broadly utilized NLTK which is a free module for Python.

The Natural Language ToolKit (NLTK) is an arrangement of modules, instructional exercises and activities which are open source and cover Natural Language Processing typically and measurably. NLTK was created at the University of Pennsylvania in 2001 allowing computational phonetics considering three instructive applications: projects, assignments and exhibitions. It can be found inside the Python

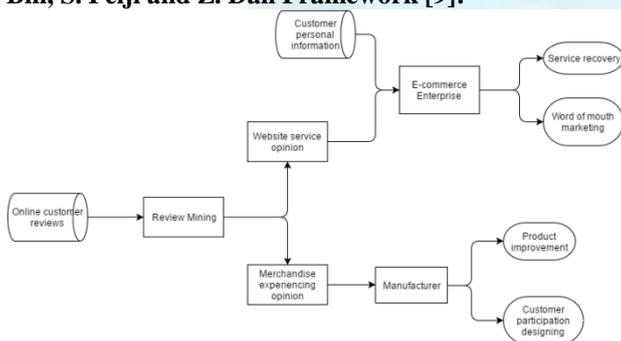
Libraries for Graph manipulation GPL open permit. NLTK is utilized to part words in a string of content and separate the content into parts of discourse by labeling word names as per their positions and capacities in the sentence. The subsequent labeled words are then prepared to extricate the importance and produce a reaction as discourse or activity as required. Distinctive linguistic use principles are utilized to order the labeled words in the content into gatherings or expressions identifying with their neighbors and positions. This type of grouping is

called chunking into phrases, such as noun phrases and verb phrases.

3. RELATED WORKS

Fine Grained opinion mining using CRF [12]: A supervised learning approach using Conditional Random Fields (CRF) to not only identify product aspects and corresponding opinions, but also determining the usage of each aspect. The underlying idea of CRFs is to define a conditional probability distribution over label sequences given a particular observation sequence, rather than defining a joint distribution over both label and observation sequences as in HMMs. CRFs relax the strong independence assumptions made in HMMs. Specifically, the tag of a word may depend to a large extent on its neighboring words and/or their POS tags. CRFs allow us to freely define features capturing this property, unlike in case of using HMMs. Furthermore, field segmentation techniques based on HMMs are not suitable for the problem either, as they usually perform well when there is some structure for fields in the sentences (according to the domain), such as bibliographical citations, house advertisements, and so on. On the other hand, the product aspects, aspect usages and opinions can appear in the sentences in many different forms, and also the majority of the words are background words which do not belong to any categories.

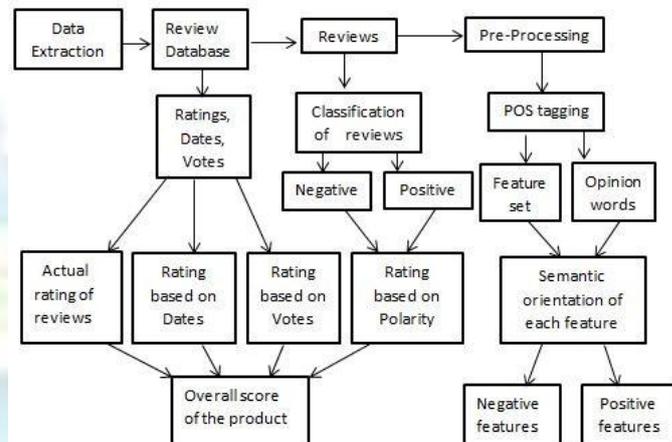
Bin, S. Peiji and Z. Dan Framework [9]:



Recommending products to customers using opinion mining of online product reviews and features [4]:

A prototype web based system for recommending and comparing products sold online. Methodology used here for this prototype contains natural language processing to automatically read reviews and used Naive Bayes classification to determine the polarity of reviews. This extracts 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 and the polarity of reviews. The availability of robust machine learning algorithms and tools, companies and individuals are able to create platforms that can help to compare products based on reviews compare E-Commerce sites, recommend Products to customers, make decisions on pricing and promotion of products

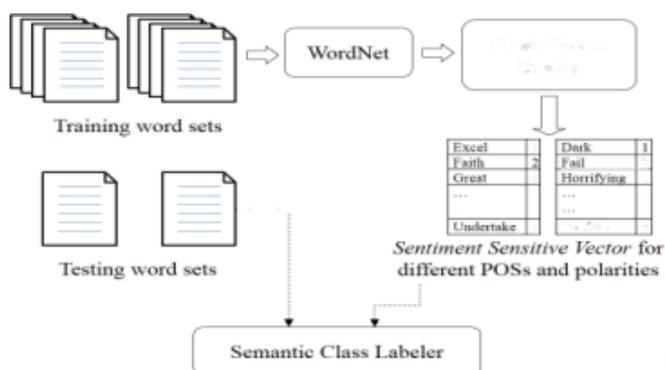


Learning-based aspect identification in customer review products [5]:

Learning-based approach using decision tree and rule learning to generate pattern set based on sequence labelling. The patterns will be used to identify and extract aspect in customer product review combined with opinion lexicon. Our experiment results based on some generated pattern using Decision Tree and Rule Learning, show that the generated pattern can produce better performance than baseline model. However, there is significant increase in the number of patterns generated from learning-based aspect extraction compared with previous pattern.

Constructing Sentiment Sensitive Vectors for Word Polarity Classification [1]:

Sentiment classification has been applied to scenarios such as mining customer reviews on merchandise sold online and film reviews of movies. Word polarity classification can be extended to perform classification of sentences and paragraphs. Initially adjectives were considered as primary sentiment carriers in a document. Eventually, other parts of speech were such as adverbs, adjective phrases, non-affective adjectives etc were considered useful for sentiment mining. A sentiment can be either positive or negative. Polarity can be determined by creating and referring a dictionary with words annotated with their semantic orientation and incorporating the relative intensification and negation. Due to the high cost of constructing a sentiment dictionary, a second approach in which sentiment of words is identified through machine learning techniques was introduced.



Even though this method is straight forward and effective, more experiments can be conducted since the proportion of words with different parts of speeches are not similar in the current data source.

Graph-Based Opinion Entity Ranking Customer Reviews [8]: Graph-based opinion entity ranking framework is used to mine opinion data from former customers and rank their entities or aspects in accordance with those opinions. Aspects and sentiment words for each entity is extracted. Relationships between reviewers and pairs of entity aspect are then represented by a weighted bipartite graph and an algorithm to compute the ranking scores is proposed.

This method applies econometric analysis and text mining technique to devise two ranking algorithms: a consumer-oriented ranking base that ranks reviews according to their expected helpfulness and a manufacturer-oriented base that ranks reviews according to their expected effects on sales. The orientation (positive, negative or neutral) of aspects is assigned by examining its corresponding sentiment word through SentiWordNet library. A Graph representation is constructed based on the relationships between reviewers and entity-aspect pairs and finally a ranking score is calculated for an individual aspect, entity, or reviewer.

User study for performance evaluation in this method found that rankings produced from this framework have higher degree of agreement with those of user study than those of the traditional baseline reported in the literature.

MARS: A Multi-Aspect Recommender System for Point-of-Interest[2]: MARS (Multi-Aspect Recommender System) is a novel POI recommender system based on multi-aspect user preference learning from reviews by using utility theory.

POI recommendation is one of the important applications of Location Based Social Network (LBSN) which brings great uses to both customers and merchants. Customers gain better user experiences through easy exploration of their interested merchants referred by other customers and the merchants may attract more customer visits and

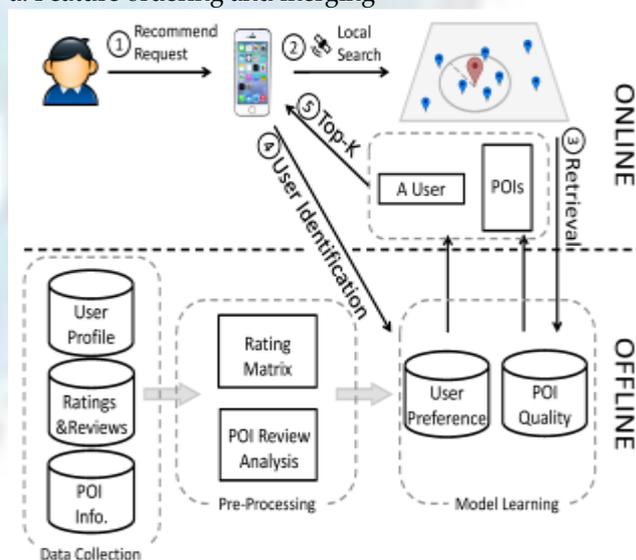
increase business turnover. MARS encompasses the following benefits:

- _ It provides a multi-aspect recommendation framework in combination of ratings and reviews.
- _ It supports a multi-mode visualization for merchants to quantitatively understand their qualities across multiple aspects.
- _ It implements POI recommendations to users based on their multi-aspect preferences.

FEROM (Feature Extraction and Refinement) Using Genetic Algorithm[3]: FEROM-GA (Genetic Algorithm) showed satisfactory results on (www.testfreak.com) through a series of experiments conducted on real product (Electronic, clothing, accessories etc) review data. The system gives good performance by using genetic algorithm in a virtual/real opinion mining framework. It increases the efficiency of search, based on the experimental result we conclude that the proposed system is a proper method for finding opinion of product.

Necessary steps in Genetic Algorithm

1. Enter the product name.
2. Fetch all the reviews of the product.
3. Clean the review.
4. Preprocessing Module
 - a. Morphological analysis and Sentence splitting
5. Extraction Module
 - a. Feature selection method
 - b. Opinion information extraction method
 - c. Opinion phrase conversion method
6. Refiner Module
 - a. Feature ordering and merging



4. CONCLUSIONS

This survey has covered a number of selected papers that have focused specifically on Opinion mining techniques. A study of eleven selected IEEE papers has

been presented, and the contribution of each study has been identified. Pros and Cons of each technique were analyzed. From the survey above, it can be said that the development and improvement of opinion mining software has not grow at a predictable rate due to a variety of methods used. The techniques of Opinion mining are still a matter for debate and no common approach has yet been identified. Researchers have so far worked in isolated environments with reluctance to divulge any improved techniques they have found, consequently, slowing down the improvements.

There are many improvements which can be made to the opinion mining application in terms of using further linguistic and contextual clues: the development of the application described here is a first stage towards a more complete system, and also contextualizes the work within a wider framework.

REFERENCES

- [1] C. H. Chu, A. H. Roopa, Y. C. Chang and W. L. Hsu, "Constructing sentiment sensitive vectors for word polarity classification," 2015 Conference on Technologies and Applications of Artificial Intelligence (TAAI), Tainan, 2015, pp. 252-259.
- [2] X. Li, G. Xu, E. Chen and L. Li, "MARS: A multi-aspect Recommender system for Point-of-Interest," 2015 IEEE 31st International Conference on Data Engineering, Seoul, 2015, pp. 1436-1439.
- [3] R. Mishra, "FEROM (Feature extraction and refinement) using genetic algorithm," 2015 International Conference on Applied and Theoretical Computing and Communication Technology (iCATccT), Davangere, 2015, pp. 344-350.
- [4] P. V. Rajeev and V. S. Rekha, "Recommending products to customers using opinion mining of online product reviews and features," Circuit, Power and Computing Technologies (ICCPCT), 2015 International Conference on, Nagercoil, 2015, pp. 1-5.
- [5] W. Maharani, D. H. Widyantoro and M. L. Khodra, "Learning-based aspect identification in customer review products," Electrical Engineering and Informatics (ICEEI), 2015 International Conference on, Denpasar, 2015, pp. 71-76.
- [6] Y. Luo and W. Huang, "Product Review Information Extraction Based on Adjective Opinion Words," Computational Sciences and Optimization (CSO), 2011 Fourth International Joint Conference on, Yunnan, 2011, pp. 1309-1313.
- [7] L. Liu, W. Wang and H. Wang, "Summarizing customer reviews based on product features," Image and Signal Processing (CISP), 2012 5th International Congress on, Chongqing, 2012, pp. 1615-1619.
- [8] K. Chutmongkolporn, B. Manaskasemsak and A. Rungsawang, "Graph-based opinion entity ranking in customer reviews," 2015 15th International Symposium on Communications and Information Technologies (ISCIT), Nara, 2015, pp. 161-164.
- [9] D. Bin, S. Peiji and Z. Dan, "E-Commerce Reviews Management System Based on Online Customer Reviews Mining," Innovative Computing & Communication, 2010 Intl Conf on and Information Technology & Ocean Engineering, 2010 Asia-Pacific Conf on (CICC-ITOE), Macao, 2010, pp. 374-377.
- [10] D. D. Chaudhari, R. A. Deshmukh, A. B. Bagwan and P. K. Deshmukh, "Feature based approach for review mining using appraisal words," Emerging Trends in Communication, Control, Signal Processing & Computing Applications (C2SPCA), 2013 International Conference on, Bangalore, 2013, pp. 1-5.
- [11] J. Bross, "A Distant Supervision Method for Product Aspect Extraction from Customer Reviews," Semantic Computing (ICSC), 2013 IEEE Seventh International Conference on, Irvine, CA, 2013, pp. 339-346.
- [12] S. Shariaty and S. Moghaddam, "Fine-Grained Opinion Mining Using Conditional Random Fields," 2011 IEEE 11th International Conference on Data Mining Workshops, Vancouver, BC, 2011, pp. 109-114.