# Ontology Based Sentiment Clustering Of Movie Review

**C.Sivagami[1], S.C.Punitha[2]**

[1] Research Scholar, Department of Computer Science,
PSGR Krishnammal College for Women, India

[2] Head, Department of Computer Science,
PSGR Krishnammal College for Women, India

***Abstract:*** *Sentiment analysis is the mining the sentiment or opinion words and identification and analysis of the opinion and arguments in the text. Text document clustering contains opinions or sentiments about the objects, such as product reviews, movie reviews, and book reviews etc. This paper presents a method of ontology-based sentiment clustering to cluster and analyse the movie reviews. In this paper, we proposed the domain ontology to extract the related generic category, in order to find the class of the movie based on the theme (Comedy, Action, Thriller, Tragedy, Horror, Sci-Fi, Family) the domain ontology and the related adjectives are analyzed and applied to the FCM clustering process the attributes with a fuzzy score is calculated and used semi supervised learning for clustering techniques.*

**Keywords:** Sentiment Analysis, Ontology using Formal Concept Analysis Design, Opinion Mining, Bag-Of-Words, Fuzzy C-Means Clustering.

## 1. INTRODUCTION

In recent years, Web 2.0 has exploded with user-generated different platforms such as Internet forums, blogs, discussion boards and forums, and social network sites, customer review sites etc. Opinionated consumers have at their disposal unprecedented reach and power by which to share their brand experiences and opinions. Major companies are now beginning to realize that their consumers wield more influences than ever before and as results are beginning to look at reviews of their products and services more closely. In fact, this phenomenon is not only related to companies that make products and services, but also public figures and celebrities. These entities can then respond to consumer insights and public sentiments by monitoring the places that generate them, and ultimately gain an advantage that has so far only come about with an extensive and usually expensive consumer research campaign. According to two surveys of more than 2000 American adults 81% of Internet users have done online research on a product at least once and among73% and 87% report that reviews had a significant influence on their purchases. However, the sheer volume of user-generated opinion content has become so enormous that most companies and public figures have to spend a lot of time engaged in it to find an overall sentiment. Indeed, even fellow consumers spend a lot of time sifting through many reviews to find that one review that discusses features and matters that is important to them. One way for organizing such data is known text classification, which involves mapping documents into topical categories based on the occurrences of particular features. Sentiment Analysis (SA) can be framed as a text clustering task where the categories are polarities such as positive and negative. However, the similarities end here. Whereas general text clustering is concerned with features that distinguish different topics, sentiment analysis deals with features about subjectivity, affect, emotion, and points-of-view that describe or change the related entities. Since user-generated review documents contain both kinds of features, SA solutions ultimately face the challenge of separating the objective content from the subjective content describing it. For example, taking a segment from a randomly chosen document in Pang et al.'s [1] movie review corpus, we see how entities and modifiers are related to one other. In this paper, we describe the sentiment clustering based on the combination approach of Natural Language Processing (NLP), Formal Concept Analysis (FCA) based ontology on movie review domain and Fuzzy C-means (FCM) cluster.

## 2. LITERATURE REVIEW

Sentiment analysis of reviews has been the focus of recent research. Several techniques are used for the opinion mining and sentiment analysis tasks. It has been attempted in different domains such as product reviews, movie reviews, and customer feedback reviews [1], [2]. Much of the research until now has focused on training Machine Learning algorithms such as Support Vector Machines (SVMs) to classify reviews. Research has also been done on positive/negative term-counting methods and automatically determining if a term is positive or negative.

### 2.1 Taxonomy Based Similarity Extraction in words

Resnik proposed a similarity measure using information content. Rensik defined the similarity of two concepts among C1and C2 in the taxonomy as the most of the information content of all concepts C that subsume both C1 and C2. Then, the similarity among two words are

*International Journal of Emerging Trends & Technology in Computer Science (IJETTCS)*
**Web Site: www.ijettcs.org Email: editor@ijettcs.org, editorijettcs@gmail.com**
**Volume 2, Issue 4, July – August 2013**                                        **ISSN 2278-6856**

defined as the greatest of the similarity among any concepts that the words belong with. Rensik used Word Net as the taxonomy; information content is calculated using the Brown corpus.

Li et al. [4] combined structural semantic information from a lexical taxonomy and information content from a corpus in a nonlinear model. They found a similarity measure that uses shortest, depth, and local density, in taxonomy. Their proposed system reported a high Pearson correlation coefficient of 0.8914 on the Miller and Charles [5].They did not compare their method in terms of similarities among named entities. Lin defined the similarity between two concepts as the information that is in common to both concepts and the information contained in each person concept.

### 2.2 A Web-Based Kernel Function for Measuring the Similarity of Short Text Snippets

Cilibrasi and Vitanyi proposed a distance metric between words using only page counts retrieved from a web search engine. The planned metric is named Normalized Google Distance (NGD) and is given by Sahami and Heilman measured semantic similarity between two queries using snippets returned for those queries by a search engine. For each query, they collected snippets from a search engine and represent each snippet as a TF-IDF-biased the term vector. Each vector is L2 standardized and the centric of the set of vectors is calculated. Semantic similarities among two queries are then defined as the inner product among the corresponding centric vectors. They did not test their similarity measure with taxonomy-based similarity measures.

### 2.3 Web-based information Extraction using mining

Web-based information Extraction can work its way through a website, downloading pages, objects, images, PDF files and streaming media, and other dynamic content that it discoveries so that they can be viewed locally, without needing to be attached to the internet. The sites can be saved locally as a fully brows able website which can be viewed with any browser (such as Internet Explorer, Netscape, Opera, Mozilla Firefox, etc), or they can be saved into the Internet Explorer cache and viewed using IE's offline mode as if the you had surfed the sites 'by hand'.

Web based information Extractor, simply enter a starting URL, and strike the Go button. The program will then take the page to that URL, parsing the HTML as it goes, looking for links to other pages and objects etc. It will extract the list of sub-links and downloads it. This process continues recursively until either no more links fulfill the system's filter criteria or computer's hard disk becomes full - which ever happen first. The locally saved files will have their HTML links adjusted so that they can be browsed as if they were being read directly from the internet.

### 2.4 Machine learning approach in classifying relationship of heterogeneous web pages

Hinton et al. [19] describe a new machine learning approach that creates expert-like rules for field matching, a key sub process in record linkage. In this approach, they defined the affiliation among two field values using a set of heterogeneous transformations, and could produce more accurate results by modeling more advanced relationships. In other paper, Hinton [20] had discussed, there have been many others published in various journals and conferences; while all these papers attempt to improve the accuracy of record linkage technique, it has not been any work that attempts to make the online record linkage process more efficient by reducing the communication overhead in a distributed environment.
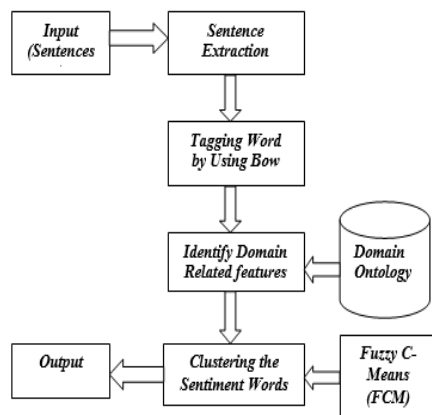
## 3. MOTIVATION

The Aim of the system is to extract a set of new feature selection schemes that use a Content and Syntax model to learn automatically a set of features in a review document by separating the entities that are being reviewed from the subjective expressions that describe those entities in terms of divisions, by spotlighting only on the personal terminology and ignoring the entities, the system needs to choose more salient features for document level sentiment analysis and review clustering. In this approach we used domain ontology and use semi supervised learning technique to extract the features and opinions from the movie reviews or comments to enhance the existing sentiment clustering tasks. We show the direct opinions on the feature level to get the opinions of each feature.

As our contribution we developed the domain ontology of the problem domain based on the FCA design, [15], and [16]. And then use the FCM clustering to train the features based on the semi-supervised learning which classified the opinions or sentiments in the reviews or comments as positive or negative on each feature.

## 4. PROPOSED SYSTEM

In this system we focus on the feature level sentiment clustering. There are three main parts in our approach: assigning the Bag-of-words, identifying domain related features and clustering the sentiment words. We use Bag-of-words to assign Bag-of-words tags to words in a sentence (such as: tags for nouns, verbs, and adjective). After that we use domain ontology to extract the related concepts and attributes based on the generic group, to find the class of the movie based on the theme (Comedy, Action, Thriller, Tragedy, Horror, Sci-Fi, Family) the domain ontology and the related adjectives and the representations are analyzed and applied to the clustering process which uses FCM for labeling positive or negative or neutral on the related concepts and attributes with a fuzzy score is calculated.

Our proposed system is as shown in the following figure.

Proposed system Architecture

## 5. METHODOLOGY

### 5.1 Assigning the Bag-of-words tags

To make this process, first we need to do sentence extracting. It generally consider '.', '!' and '?' as sentence delimiters for splitting process, although allowances are given for the occurrence of '.' in abbreviations, number, URLs etc. For identifying sentiment phrases we use BOW tagging. To make this process we use Microsoft NLP Library. This tagger is based on transformation based error ambitious training, a method which is efficient in several natural language applications which include the bag of words and word sense tagging, prepositional phrase attachment, and semantic parsing. The BOW tagger assigns BOW tags to words in a sentence (such as tags for nouns, verbs, and adjective). These tags are used to extract the certain features related to the proposed movie review domain.

### 5.2 Identifying the domain related features

We use the domain ontology to get the domain related features. Ontology aims to give knowledge about specific domain that is understandable by both developers and Input computers and necessary for knowledge representation and knowledge exchange. Using existing taxonomical hierarchies are not enough for knowledge exchange or for informational retrieval. By using the taxonomical ordering, the concepts have no other differentiating attributes. It is not easy to change the frames and their slots once they are defined. These may cause the problems in knowledge sharing. Therefore we need a better way to describe the concepts and relation. We present a method that is based on Formal Concept Analysis (FCA), used for analyzing data and forming semantic structures that are formal abstraction of concepts of human thoughts and find conceptual structures between data sets. It also allows the analysis of complex structures and the discovery of dependencies within the data. The adjective "formal" is used to emphasize that these are formal notions. "Formal objects" need not be "objects" in any kind of common sense meaning of

"object." But the use of "objects" and "attributes" is indicative because in many applications it may be useful to choose object-like items as formal concepts and to choose their features or characteristics as formal attributes. The main characteristics of FCA,
• Concepts are described by properties
• The properties find the hierarchy of the concepts
• When the properties of different concepts are the same, then the concepts are the same Contexts in FCA are triples (O, A, R)

Where O=finite set of object,
A=finite set of attributes,
R=binary relation on O and A. We developed the domain ontology in OWL (Web Ontology Language) based on FCA design. Implementation of domain ontology is used protégé 2000. Features that are parts of the domain is extracted to be clustered by FCM.

### 5.3 Clustering the sentiment words using FCM

Sentimental words play a major role in how we describe and understand opinions. In this paper, we presented an experiment that used interval type fuzzy C Means from Sentimental words. Our method featured that prompted users with Sentimental words, which refer to good or bad. Each Sentimental word was rated by users with intervals on valence, activation and dominance dimensions' data are collected by user. It is essential to collect data about how people describe or know Sentimental words naturally because the aim of this research is to make communication easier about chosen specific Sentimental words among human and computer. We plan to confirm the results of this paper by experiments on the survey and the natural language corpora, which will be analyzed in more detail to consider carrying out Interval Approach for sentence-level. In addition, for both sentence and word level, we will calculate fuzzy C Means similarity metrics between fuzzy C Means membership functions that projects the related distances for visualization using multidimensional scaling. Beside we also plan to look at the inter and intra subject variability, and create a vocabulary of the most specific Sentimental words since it should be large enough so that a human will feel comfortable in interacting with a computer. The critical question we will try to answer in future study is: Is the fuzzy C Means approach useful to represent Sentimental words correlated with sentences including these sentimental words. The answer for this question will lead us to focus on this aspect.

## 6. EXPERIMENTAL RESULTS

To test sentiment clustering system, we use the customer review of a few movies form IMDB corpus which was grouped into positive, negative and neutral categories for content analysis and tested and compared with the manually tagged set of around 300 movie reviews. The results and comparisons are shown in the following Table.1.

**Table.1.** Results and comparisons of movie reviews

|  | Without the use of domain ontology | Proposed approach |
|---|---|---|
| Positive Sentences | Accuracy- 72 % Recall- 90 % | Accuracy- 90% Recall- 95% |
| Negative Sentences | Accuracy- 64% Recall- 85% | Accuracy- 85% Recall 93% |
| Ambiguous Sentences | Accuracy- 65% Recall-90% | Accuracy- 88% Recall- 97% |
| Neutral Sentences | Accuracy-60% Recall-90% | Accuracy- 89% Recall- 96% |

## 7. CONCLUSION AND FUTURE WORK

We proposed the combination approach of Bag-of-words, FCA-based domain ontology and FCM clustering intend to enhance the sentiment clustering. By using this approach we can view the strength or weakness of the movie more detail and we hope will be useful for further development and improvement of the development and improvement of the movie. As the future work we need to tested with a large number of data sets and require further training and clustering to solve the problem of the comparative sentences. The domain ontology to extract the related generic category, to find the class of the movie based on the theme (Comedy, Action, Thriller, Tragedy, Horror, Sci-Fi, Family) the domain ontology and the related adjectives and the representations are analyzed and applied to the FCM clustering process which applies positive or negative or neutral on the related concepts and attributes with a fuzzy score is calculated..

## References

[1] Pang, B. and L. Lee, (2008) "Opinion Mining and Sentiment Analysis", Foundations and Trends in Information Retrieval, 2(1-2), 1-135.

[2] Bo Pang and Lillian Lee, "Using very simple statistics for review search: An exploration". In Proceedings of the International Conference on Computational Linguistics (COLING), 2008. Poster paper.

[3] Peter Turney, "Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews". In Proceedings of the Association for Computational Linguistics (ACL), pages417–424, 2002.

[4] Liu, Y., X. Huang, A. An, and X. Yu, (2007) "ARSA: A Sentiment-Aware Model for Predicting Sales Performance Using Blogs", Proc. 30th Annual Int. ACM SIGIR Conf. on Research and Development in Information Retrieval, pp. 607-614, Amsterdam, the Netherlands, July 23-27.

[5] Miller, G. A. (1995) "WorldNet: A Lexical Database for English", Communications of the ACM, 38(11), 39-41.

[6] Sanjiv R. Das, Mike Y. Chen, "Yahoo! for Amazon: sentiment extraction from small talk on the Web", Management Science 53 (9) (2007) 1375–1388.

[7] Xiaowen Ding, Bing Liu, Philip S. Yu, "A holistic lexicon-based approach to opinion mining", Proceedings of the Conference on Web Search and Web Data Mining (WSDM), 2008.

[8] Thomas L. Griffiths, Mark Steyvers, David M. Blei, Joshua B. Tenenbaum, "Integrating topics and syntax", In Advances in Neural Information Processing Systems, 17, MIT Press, 2005, pp. 537–544.

[9] Vasileios Hatzivassiloglou, Kathleen McKeown, "Predicting the semantic orientation of adjectives", Proceedings of the Joint ACL/EACL Conference, 1997, pp. 174–181.

[10] Vasileios Hatzivassiloglou, Janyce Wiebe, "Effects of adjective orientation and grad ability on sentence subjectivity", Proceedings of the International Conference on Computational Linguistics (COLING).

[11] Abbasi, A. (2010). "Intelligent feature selection for opinion classification". IEEE Intelligent Systems, 25, 75–79.

[12] Abbasi, A., Chen, H., & Salem, A. (2008). "Sentiment analysis in multiple languages: Feature selection for opinion classification in web forums". ACM Transactions on Information Systems, 26, 12:1–12:34.

[13] Abbasi, A., France, S., Zhang, Z., & Chen, H. (2011). Selecting attributes for sentiment classification using feature relation networks. IEEE Transactions on Knowledge and Data Engineering, 23, 447–462.

[14] Alpaydin, E. (2004). "Introduction to machine learning". In Proceedings of the Australian joint conference on advances in artificial intelligence (pp. 100–109).

[15] Baccianella, A. E. S., & Sebastiani, F. (2010). "Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining". In Proceedings of the 7th conference on international language resources and evaluation (LREC'10). Valletta, Malta: European Language Resources Association (ELRA).

[16] Bai, X. (2011). "Predicting consumer sentiments from online text". Decision Support Systems, 50, 732–742.

[17] Genter B, Wille, R. "Formal Concept Analysis, Mathematical Foundation", Berlin, Springer Verlag 1999.

[18] Marek Obitko et al. "Ontology Design with Formal Concept Analysis". In CLA 2004, pp, 111-119, ISBN 80-248-0597-9

[19] Hinton, G. E. and Salakhutdinov, R. (2006)." Reducing the dimensionality of data with neural networks". Science, 313(5786), 504-507.

[20] Hinton, G. E., Osindero, S., and Teh, Y. (2006). "A fast learning algorithm for deep belief nets". Neural Compu- tation, 18, 1527-1554.