

# Discovering the Features in Opinion Mining via Domain Dependent and Domain Independent Relevance

Pallavi D. Jawalkar<sup>1</sup>, G. J. Chhajed<sup>2</sup>

<sup>1</sup>ME-II Student, Vidya Pratishthan's College of Engineering, Baramati, Savitribai Phule Pune University, Pune, Maharashtra, India.

<sup>2</sup>Assistant Professor, Vidya Pratishthan's College of Engineering, Baramati, Savitribai Phule Pune University, Pune, Maharashtra, India.

**Abstract:** *Opinion feature mining is also known as aspect mining used to take out users opinions, and attitudes towards a specific product, services and their characteristics. The most of the existing approaches to opinion feature extraction on mining patterns is only by using a single review corpus. This paper presents the new method to discover the opinion features from online reviews by taking out the difference in opinion feature statistics across two different corpora, one domain specific corpus and another is domain independent corpus (i.e. the contrasting corpus). Domain relevance is the measure which is used to capture the disparity. The domain relevance characterizes the relevant term from the text collection. Firstly, the sentences are extracted from the reviews. Then the POS Tagger is applied to separate out the nouns, noun phrases and adjectives. Next the candidate features are extracted by applying the syntactic rules designed for Standard English. For every candidate feature, the Intrinsic Domain Relevance (IDR) and Extrinsic Domain Relevance (EDR) scores are calculated by using Domain dependent and domain independent corpus respectively. a The interval threshold approach, called as IEDR Criteria is applied to confirm the final Opinion Feature in which the candidate feature having IDR score greater than IDR threshold, and EDR scores less than EDR threshold is checked.*

**Keywords:** Opinion mining, Domain relevance, part-of-speech tagging, Opinion Feature

## 1. Introduction

What the user think about has an important issue of information for us while taking the decision.. Mostly we asks to our friend to recommend or give opinion about any product, service, regarding job, either positive or negative. What the user think about has an important issue of information for us while taking the decision. Mostly we asks to our friend to recommend or give opinion about any product, service, regarding job, either positive or negative.

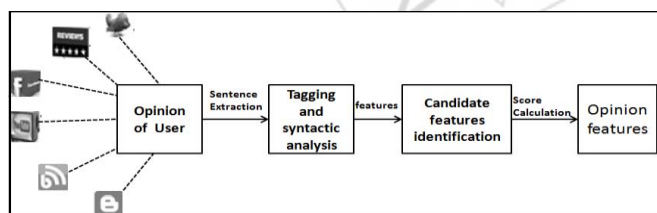


Figure 1: Process of Opinion Mining

Now a days people are planning to develop a system which can be identify and classify opinion or sentiment as represented in an electronic text i.e. (e-text). By analyzing every text the opinion mining system see the part which contains opinionated word, which is to be opinionated and who is written the opinion. Opinions expressed in textual reviews are generally pruned on different dimensions. Opinion mining is works at document level and feature level. With the marvellous growth of social media such as reviews, remarks, comments and postings in social media web sites on the web are used by personal, private and company, for taking the decision.

Opinion features are the attributes of an entity on which the opinions are suggested. The polarity of opinion is referred by

the point of reference. A major research area in this domain is of opinion feature recognition and extraction which has already been considered and various techniques such as Natural Language Processing techniques and modelling techniques are proposed. In real life reviews, the syntactic rules which are used in NLP are not working properly as these reviews having lack formal structure, modelling techniques which creates semantic rules are used for coarse grained analysis.

### 1.1 Problem Definition

Identification of opinion features are useful for decision making regarding product, service selection etc. In this work opinion features are collected by checking the difference in opinion feature statistics across two corpora. This work not only uses domain-specific corpus but also refers domain independent corpus for better opinion feature selection. By using Part-Of- Speech Tagger probable opinion features are extracted, which are further processed using IEDR criterion.

## 2. Related Work

- Vasileios Hatzivassiloglou and Jance Wiebe [2] is stated Supervised classification method for an effect of adjectives on prediction of the subjectivity of opinions.
- Yessenalina and Cardie [4] stated a compositional matrix space model whose mechanism is of phrase-level sentiment analysis.
- Advantage of the this model is that the model can manage invisible word composition by learning matrices for that word.
- Zen Hai and C Yang [5] stated a model to unsupervised natural language processing to discover the candidate

feature.

- The authors Pang and Lee [6] used Machine learning algorithms. The algorithms Nave Bayes, Maximum entropy classification and support vector machine (SVM) which are used for text classification Naive Bayes classification method classifies the on the whole document by Bayes rule.
- Ryan McDonald and Kerry Hannan [7] have stated a structured model for classifying sentiments at a variety of different levels of granularity in that document level, sentence level or word level which is also called as fine to coarse sentiment analysis.
- Lizhen Qu and Georgiana Ifrim [8] have proposed a model based on regression approach. From sparse text pattern, this method is used to forecasting the review ratings. The regression method proposes an algorithm for calculating opinion scores from regression method.
- W.jin and H. H. Ho,[9] have planned the supervised machine learning framework that use lexicalized Hidden Markov Model. The linguistic features integrated into automatic learning, supported by this framework model.
- Hanshi Wang Lizhen Liu presents the new method that uses the similar type of opinion words to take out the features it filters the noises according to mutual support scores and confidence score. The method implicit the features and clusters the features which are based on the knowledge of the context dependent information .
- Latent Dirichlet Allocation (LDA) [10]: This is the probabilistic model for collections of distinct data as text corpora. LDA is also called as three-level hierarchical Bayesian model. In this model which every item of a collection is on sidered as a finite combination over the given set of topics.
- Association Rule Mining (ARM)[11] : Consumers have commented on the mining product features. To identify opinion comments on that review and make confirm that is it each and every opinion sentence is positive or negative.
- Mutual Reinforcement Clustering (MRC) [12] : The research on feature-level opinion mining basically believe on recognising the explicit related among the product feature words and opinion words in reviews.

### 3. Proposed System

#### 3.1 Proposed system Architecture

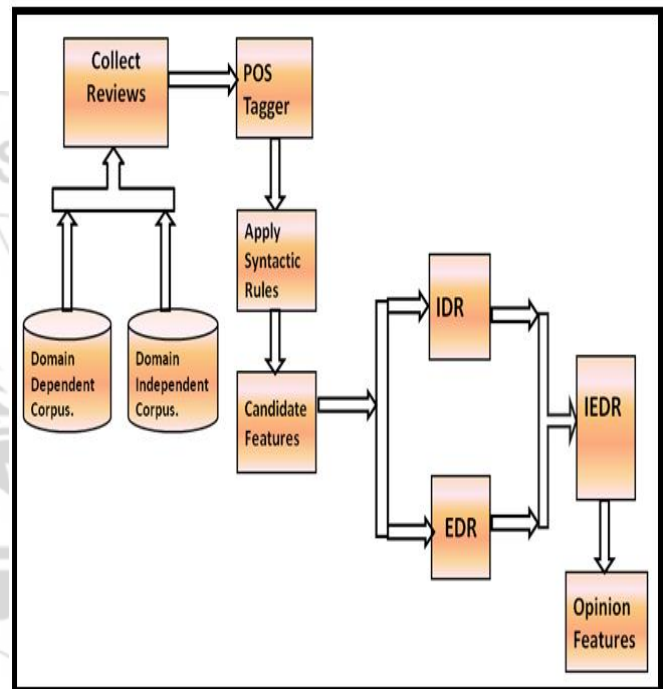
As shown in figure 2, this proposed system focuses on opinion feature identification. In this, firstly load the Review corpus of Domain dependent corpus and domain independent corpus. Sentences are extracted from reviews which usually provided with XML version. These sentences should loaded into array of data type String.

A Part-Of-Speech Tagging also called grammatical tagging. This is the process of marking up each word in a text as analogous to a particular part of speech which is based on not only its definition but its contents also.

By applying Syntactic rules which are in standard English language form gives the candidate features. Dispersion and deviation are calculated using the term frequency and inverse document frequency (TF-IDF). Intrinsic Domain Relevance score is calculated for each extracted candidate feature of domain dependent corpus.

Extrinsic Domain Relevance score is calculated for each extracted candidate feature of domain independent corpus.

The Intrinsic Extrinsic Domain Relevance is a criterion to identify and exact opinion feature. The candidate feature is considered as opinion feature only when IDR score is greater than the threshold value and EDR score is less than another threshold value. Finally opinion features are confirmed.



**Figure 2:** Proposed system Architecture

#### 3.2 Algorithm

**Algorithm 1:**  
 Calculating Intrinsic and Extrinsic Domain

##### Relevance IDR/EDR

**Input:** A domain dependent / independent corpus.

**Output:** Domain relevance scores(IDR or EDR).

##### STEPS

- 1) For each candidate feature **do**
- 2) For each document in the corpus **do**
- 3) Calculate weight ;
- 4) Calculate standard deviation ;
- 5) Calculate dispersion ;
- 6) For each document in the corpus **do**
- 7) Calculate deviation ;
- 8) Compute domain relevance;
- 9) return A list of domain relevance (IDR/EDR) scores for all candidate features;

**Algorithm 2:  
 Identifying Opinion Features using IEDR criteria.**

**Input:** Domain Review corpus and Domain independent corpus

**Output:** A validated list of opinion features.

**STEPS**

- 1) Extract candidates from the review corpus.
- 2) for each candidate feature **do**
- 3) Compute IDR score  $idr_i$  via Algorithm 1 on the review corpus
- 4) Compute EDR score  $edr_i$  via Algorithm 1 on the domain independent corpus
- 5) **if** ( $idr_i \geq i^{th}$ ) AND ( $edr_i \leq e^{th}$ ) **then**
- 6) Confirm candidate as a feature
- 7) **return** A validated set of opinion features.

**4. Implementation Details**

The proposed framework is organized into following phase.

- 1) POS Tagging
- 2) Apply Syntactic Rules
- 3) Measure IDR/EDR score
- 4) Apply IEDR
- 5) Extract Opinion feature

**4.1 POS Tagging**

A Part-Of-Speech Tagging also called grammatical tagging. this is the process of marking up each word in a text as analogous to a particular part of speech which is based on not only its definition but its contents also. The relationship with adjacent and related words in a phrase. This work use open source Stanford NLP parser for POST. The parser is instantiated with English Model.

**4.2 Apply Syntactic Rules**

By applying the following syntactic rules, can identify noun, noun phrase, noun plural, and adjectives. The output of this process is taken as candidate features.

The table shows the following rules which are used in Standard English language.

**Table 1:** List of Syntactic Rules

Rules	Interpretation
NN+JJ→CF NN+VB→CF JJ+NN→CF VB+NN→CF	Extract NN as CF
NNP+JJ→CF NNP+VB→CF JJ+NNP→CF VB+NNP→CF	Extract NNP as CF
NNS+JJ→CF NNS+VB→CF JJ+NNS→CF VB+NNS→CF	Extract NNS as CF

**4.3 Measure IDR/EDR score**

After getting candidate feature, the IDR/EDR scores are calculated by using the domain relevance which is applied on the domain dependent and domain independent corpus respectively.

**4.4 Apply IEDR**

The Intrinsic Extrinsic Domain Relevance is a criterion which is used to identify and extract opinion features. The value of IDR of that candidate feature is more than threshold value and value of EDR is less than another threshold value, then that candidate feature is considered as an Opinion feature.

**4.5 Extract Opinion feature**

After applying IEDR over the candidate feature lastly we get the final feature called as Opinion Feature.

**5. Result Analysis**

**5.1 GUI of System**

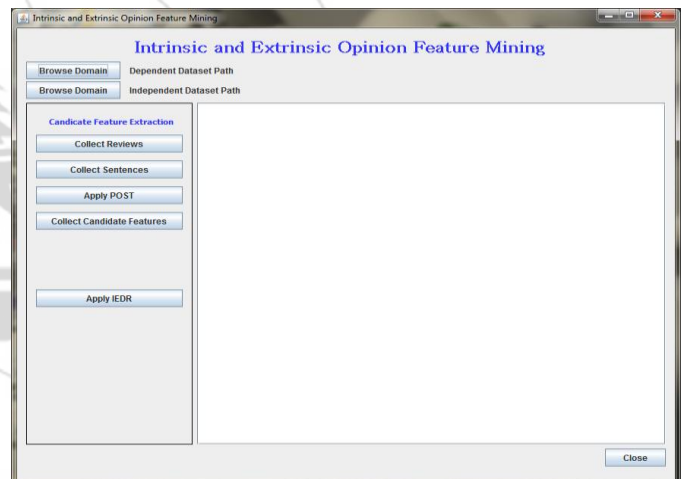
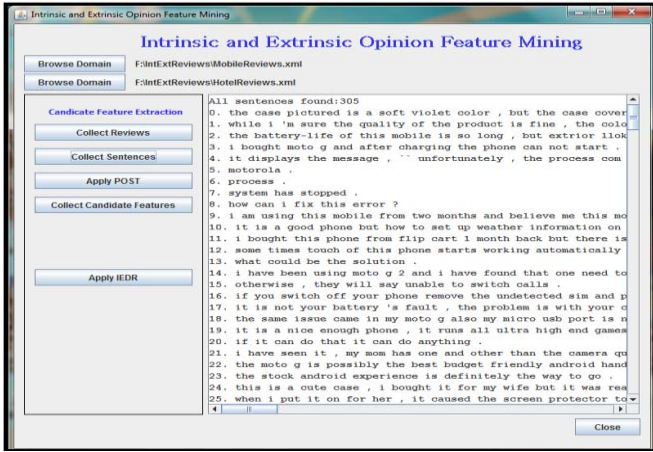


Figure shows GUI of the system. This work requires two domain corpora, one is domain dependent corpus and another is domain independent corpus. In figure, mobile domains is domain dependent corpus and Hotel is domain independent corpus.

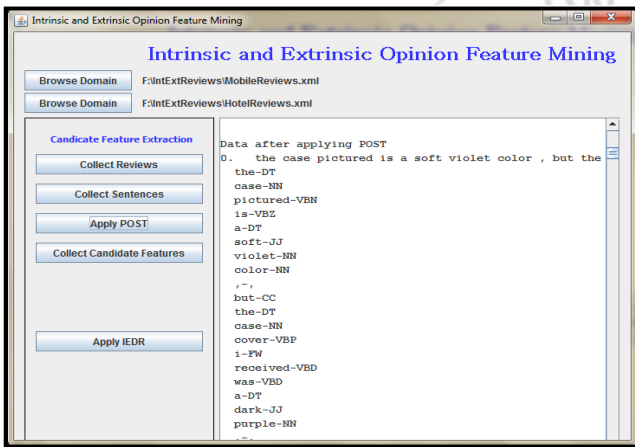
**5.2 Loading Dataset**

A user has to browse both the domains from the source location of the storage system as shown in figure.



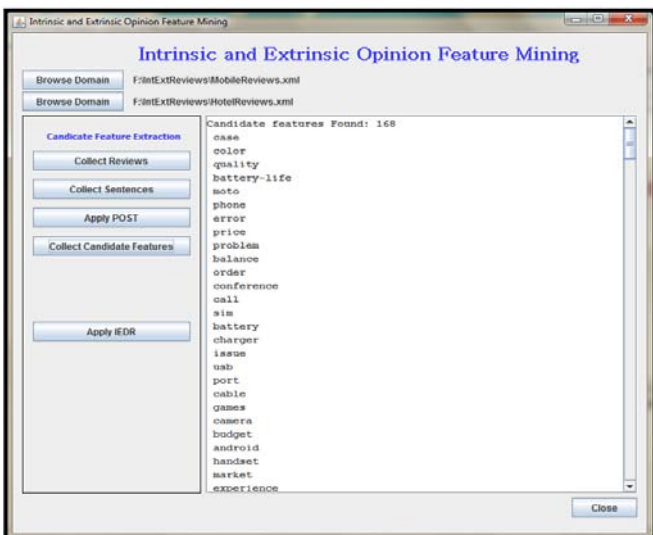
### 5.3 Apply Part-Of-Speech

Stanford NLP is used for Part-of-speech tagging which is applied on the extracted sentences to tag every word of sentence which is shown in figure.



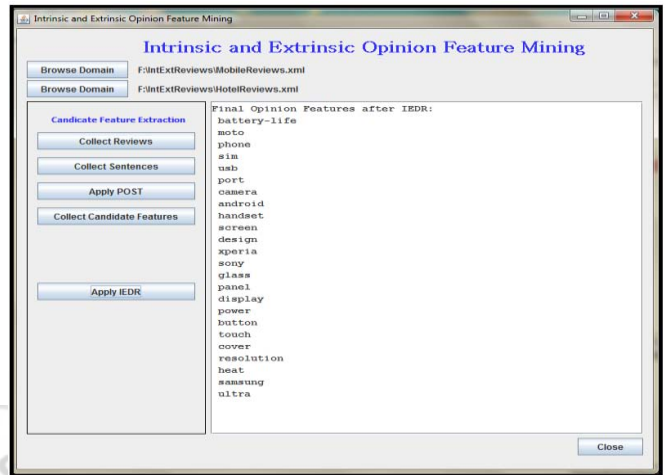
### 5.4 Collect Candidate features

The part of speech tagging is applied on the collected reviews and a set of syntactic rules are applied to identify candidate features from the review corpuses. These syntactic rules identify candidate features properly. It is shown in figure.



### 5.5 Final Set of Opinion Features

By applying IEDR criterion on candidate features it will retrieve a set of opinion features. This set of Final Opinion Features as shown in figure



### 5.6 Result Tables

Table 2: Dataset Description

Domain Dependent Corpus	#Reviews	#Sentences	#Features
Mobile	110	305	29
Hotel	108	364	32
Laptop	121	337	41

### 6. Precision and Recall

Table 3: Result table (Precision and Recall)

Sr. No.	Intrinsic Domain	Extrinsic Domain	#Features in Domain	#Retrieved Features	# Correct Features	Precision	Recall
1	Mobile	Hotel	29	24	19	0.79	0.82
2	Hotel	Mobile	32	27	22	0.81	0.84
3	Laptop	Hotel	41	33	26	0.78	0.80

Table 3. Shows precision and recall values for different Corpus. By considering this values graph is constructed.

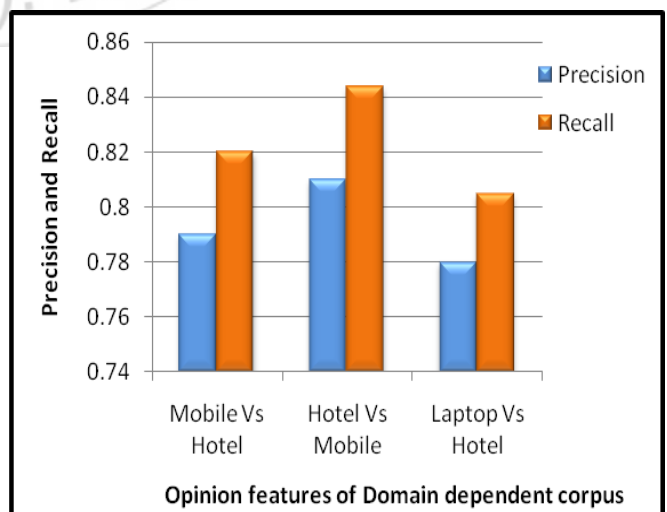


Figure 3: Graph for Precision and Recall

## 7. Conclusion

In this work, a novel approach of opinion feature extraction which is based on the IEDR feature filtering criteria is adopted. It utilizes the variation or inequality in distributional characteristics of features across two corpora, the domain-specific and domain independent.

IEDR identifies candidate features that are specific to the given review domain. POS Tagger is used to separate the noun, verb, adjectives etc from the sentence. By calculating scores will extract the features called as a opinion feature. This opinion features are helpful to Purchase a product and other decision making tasks.

## References

- [1] Hai, Kuiyu Chang, Jung-Jae Kim, and Christopher C. Yang, "Identifying Features in Opinion Mining via Intrinsic and Extrinsic Domain Relevance", IEEE Transaction on knowledge and data engineering, vol. 26, no. 3, March 2014.
- [2] V. Hatzivassiloglou and J.M. Wiebe, "Effects of Adjective Orientation and Gradability on Sentence Subjectivity", Proc 18th Conf. Computational Linguistics, pp. 299-305, 2000.
- [3] B. Pang, L. Lee, and S. Vaithyanathan, "up?: Sentiment Classification Using Machine Learning Techniques," Proc. Conf. Empirical Methods in Natural Language Processing, pp. 79-86, 2002.
- [4] A. Yessenalina and C. Cardie, "Matrix-Space Models for Sentiment Analysis,". Conf. Empirical Methods in Natural Language Processing, pp. 172-182, 2011.
- [5] Z. Hai, K. Chang, Q. Song, and J.-J. Kim, "Statistical Nlp Approach for Feature and Sentiment Identification from Chinese Reviews,". CIPS-SIGHAN Joint Conf. Chinese Language Processing, pp. 105-112, 2010
- [6] B. Pang and L. Lee, "Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts,". 42nd Ann. Meeting on Assoc. for Computational Linguistics, 2004.
- [7] R. McDonald, K. Hannan, T. Neylon, M. Wells, and J. Reynar, "Structured Models for Fine-to-Coarse Sentiment Analysis,". 45th Ann. Meeting of the Assoc. of Computational Linguistics, pp. 432-439, 2007.
- [8] L. Qu, G. Ifrim, and G. Weikum, "Bag-of-Opinions Method for Review Rating Prediction from Sparse Text Patterns," Proc 23rd Intl Conf. Computational Linguistics, pp. 913-921, 2010.
- [9] W. Jin and H.H. Ho, "Novel Lexicalized HMM-Based Learning Framework for Web Opinion Mining," Proc. 26<sup>th</sup> Ann. Intl Conf. Machine Learning, pp. 465-472, 2009.
- [10] D.M. Blei, A.Y. Ng, and M.I. Jordan, "Dirichlet Allocation,". Machine Learning Research, vol. 3, pp. 993-1022, Mar. 2003.

## Author Profile



**Ms. Pallavi D. Jawalkar** received her B.E. Degree in Computer Science and Engineering in 2012 from B.M.I.T. Solapur, Solapur University, Solapur. and persuing Master of Computer Engineering in VPCOE, baramati, Savitribai Phule Pune University, Pune



**Prof. Mrs. G. J. Chhajed** obtained her B.E. Degree in Computer Science and Engineering in 1991-95 from S.G.G.S.I.E.T, Nanded and Postgraduate Degree in Computer Engineering from College of Engineering, Pune (COEP) 2005-2007 both with distinction. She is approved Undergraduate and Postgraduate teacher of Pune University and has 19 years of experience. she has total 32 Publications in the National, International level Journals and proceedings of Conferences. She has membership of IEEE, LMCSI, IET, LMISTE and International Association LMIACSIT.