

Sentiment Classification: A Lexical Similarity Based Approach for Extracting Subjectivity in Documents

Kiran Sarvabhotla · Prasad Pingali · Vasudeva Varma.

the date of receipt and acceptance should be inserted later

Abstract With the growth of social media, document sentiment classification has become an active area of research in this decade. It can be viewed as a special case of topical classification applied only to subjective portions of a document (sources of sentiment). Hence, the key task in document sentiment classification is extracting subjectivity. Existing approaches to extract subjectivity rely heavily on linguistic resources like sentiment lexicons and complex supervised patterns based on POS information. It makes the task of subjective feature extraction complex and resource dependent. In this work, we try to minimize the dependency of linguistic resources in sentiment classification. We propose a simple and statistical methodology called RSUMM and use it in combination with well known feature selection methods to extract subjectivity. Our experimental results on a movie review dataset prove the effectiveness of the proposed methodology.

Keywords Social Media · Sentiment Classification · Subjectivity · Linguistic Resources · RSUMM

1 Introduction

The textual information on the web can be broadly categorized into two: *facts* and *opinions* (Liu 2010; Pang and Lee 2008). Until the early part of this decade, most of the research work in the areas of natural language processing (NLP), text mining and information retrieval (IR) focused on factual information. With the advent of customer reviews, blogs and the growth of e-commerce in this decade, *user-generated-content* has grown rapidly on the Web (Liu 2010). It has an inherent property called *sentiment*. These sentiments are playing a prominent role in decision-making process of the people (Gretzel and Yoo 2008). Analyzing and predicting their polarity has received much attention among research community and market analysts for its potential business applications. Hence, sentiment classification has become one of the hot topics of research in this decade.

Kiran Sarvabhotla · Prasad Pingali · Vasudeva Varma
Search and Information Extraction Lab, International Institute of Information Technology, Hyderabad, India.
Tel.: +91-40-6653 1157/1137
Fax: +91-40-6653 1413
E-mail: kiransarv@research.iiit.ac.in

Sentiment analysis or classification predicts the polarity of a given text unit. The text unit can be a word, phrase, sentence or a document. The polarity is predicted on either binary¹ or multi-variant scale². The problem of sentiment classification³ can be viewed as a special case of topical classification applied only to subjective portions of a document. In other words, topical classification focus on keywords where as sentiment classification focuses on subjectivity in the text (Pang et al. 2002). Hence, the key task in sentiment classification is extracting the subjective portions in a document. In this work, we apply supervised learning approaches to classify the overall sentiment of a document⁴ on binary scale. We focus more on the aspect of extracting subjective features, existing approaches for doing it, problems with them and present solutions.

Existing approaches in *subjective feature extraction* rely heavily on linguistic resources. Popular among them are lexicons like SentiWordNet, General Inquirer and part-of-speech (POS) tagger (Baccianella et al. 2009; Hu and Liu 2004; Matsumoto et al. 2005; Turney 2002). Lexicons are very generic and cannot capture subtle variations in expressing sentiment from context to context and from domain to domain. Using POS tagger, researchers frame patterns that are assumed to be subjective based on POS information. The POS patterns vary from simple noun phrases (NP), verb phrases (VP) to very complex patterns⁵. The text units that match these patterns in a document are extracted. There are also other techniques like using WordNet, appraisal adjectives and dependency parsing for extracting subjectivity (Mullen and Collier 2004; Whitelaw et al. 2005; Matsumoto et al. 2005).

With the usage of linguistic resources and complex patterns, the task of subjective feature extraction has become more complex and resource dependent. As regional language content is growing on the web gradually, extending the resource based approaches for analyzing sentiments across several languages is a tedious job. It requires a lot of human effort to build such linguistic resources in each language. Also, inadequate availability of resources in a language should not prevent researchers to conduct experiments on analyzing sentiments. To make the task of sentiment classification more feasible, we need approaches that minimize the use of linguistic resources in subjective feature extraction.

We attempt to address the problem of *resource dependency* in sentiment classification by assuming that entire document does not have subjective information. The basis for our claim is manual observation of documents from different domains on the Web. We have noticed many documents with less subjective content compared to total content. Hence, we can say that a sentiment bearing document is a mixture of subjective and objective information where the latter does not convey the feelings of the author. It has been proved in the literature that discarding the objective information enhances the performance of the sentiment classifier (Pang and Lee 2004). We propose a simple and statistical two-step *filtering methodology* for extracting subjective features from a document.

In the first step, we estimate the subjectivity of each sentence in a document. We use techniques very similar to vector space model in IR. We call this method of estimating subjectivity as RSUMM. We define two metrics Average Document Frequency (ADF) and Average Subjective Measure (ASM) and use them in RSUMM for estimating subjectivity of each sentence. Our RSUMM is based on lexical similarity model between two term vectors. It retains the most subjective sentences in a document discarding the objective information.

¹ positive or negative

² grading reviews typically on a scale of one to five (starred rating(*))

³ By sentiment classification, we mean document sentiment classification.

⁴ A document can be a review or a blog post

⁵ NP, VP, JJ NN, RB JJ not NN, JJ JJ not NN, RB VB, NN JJ not NN etc. where NN is a noun and RB stands for adverb, JJ stands for adjective and VB stands for verb.

Thus, we have an *excerpt* that preserves the subjectivity at a level comparable or better than the original document.

Later, we represent the document and its *subjective excerpt* as a feature vector to the classifier using n-gram models. In the next step, we apply two well known feature selection methods, mutual information (MI) and fisher discriminant ratio (FDR) on n-gram feature vectors to obtain the final feature set. We conduct experiments on a movie review dataset to validate our proposed two-step filtering methodology. Through our experimental results, we conclude that subjective feature extraction is possible with minimum usage of linguistic resources.

RSUMM was inspired from the work by Pang and Lee (Pang and Lee 2004). They established a strong relationship between sentence level subjectivity detection and document level sentiment classification. In their work, sentence level subjectivity detection had increased the performance of the sentiment classifier. They used a min-cut graph based classification framework using individual and contextual penalties for each sentence to obtain subjective excerpt of a document.

The rest of the paper is organized as follows. In section 2, we discuss the work related to document sentiment classification, subjective feature extraction and existing techniques in it. In section 3, we describe RSUMM that estimates subjectivity of each sentence in a document and optimize the sentiment. In section 4, we discuss on how to represent the text units as feature vectors to the classifier. We also discuss the feature selection methods for obtaining the final feature set. In section 5, we describe our experimental activity, that includes description of the datasets, evaluation measures and results. In section 6, we discuss the performance of the sentiment classifier and present our observations. Finally, we conclude the paper by giving possible future directions in section 7.

2 Related Work

Sentiment classification dates back to late 1990s (Hatzivassiloglou and McKeown 1997; Argamon et al. 1998; Kessler et al. 1997), but in the early part of this decade, it has become an important discipline in the areas of NLP, text mining and information retrieval. The classification is done at several levels of text units. The text unit can be a word, phrase, sentence or a document. In this section, we discuss more on the work related to document sentiment classification, subjective feature extraction using linguistic resources and resource or language independent approaches in it. Among the existing approaches in document level classification, supervised learning methods⁶ are popular among researchers to predict the polarity. Movie review domain is popular among them (Pang et al. 2002; Mullen and Collier 2004; Baccianella et al. 2009; Matsumoto et al. 2005; Pang and Lee 2004; Beineke et al. 2004; Whitelaw et al. 2005). It can be due to abundant availability of movie reviews on the Web and their challenging nature (Turney 2002).

Document sentiment classification is typically composed of two steps: 1) Extracting subjective features from training data and converting them as feature vectors. 2) Training the classifier on the feature vectors and applying classification on an unseen sample. Raw documents are preprocessed before extracting the subjective features. The preprocessing stage include removing HTML tags from a document, tokenization etc..

At word level, Hatzivassiloglou and McKeown used conjunctive expressions to extract semantic polarities of words (Hatzivassiloglou and McKeown 1997). The approach was

⁶ Support vector machines (SVM), Naive Bayes, Maximum entropy based classification

based on linguistic constraints that 'and' always conjuncts two words with same orientation whereas 'but' contradicts them. Turney proposed an unsupervised approach to predict the overall sentiment of a document (Turney 2002). The approach was based on point wise mutual information (PMI) between a given phrase and the words "excellent" and "poor" to predict the semantic orientation. The phrases were extracted using Brills tagger based on a set of POS patterns as mentioned in Section 1 and a semantic orientation predicted.

Pang et al. in 2002 experimented with several machine learning algorithms using unigrams, bigrams, POS information and sentence position as features (Pang et al. 2002). They reported that support vector machines (SVM) with unigrams as features yield best results. They conducted experiments on a movie review dataset and concluded that machine learning techniques outperform human produced baselines. They predicted the polarity of a review on a binary scale and reported an accuracy of 82.9%. Mullen and Collier used diverse information scores that assign value to each word or phrase using WordNet, topic proximity and syntactic relations (Mullen and Collier 2004). They reported an accuracy of 86% on the same dataset.

Matsumoto et al. used word sub-sequence mining and dependency parsing for extracting subjective features (Matsumoto et al. 2005). Their motive was to preserve the word order and syntactic relations between words in sentences or clauses. They extracted clauses using a clausal extraction tool. They imposed some linguistic constraints on the extracted clauses based on POS information. They used n-gram features to represent each review by setting intuitive support thresholds for selecting a feature. They reported a maximum accuracy of 88.3% on the same movie review dataset.

Pang et al. in 2004 examined the relation between sentence level subjectivity detection and document level polarity classification (Pang and Lee 2004). They did not use any linguistic resource, instead trained a naive-bayes classifier on an annotated subjective/objective collection. In addition, they incorporated contextual information and used min-cut based classification framework to obtain subjective extract of a document. They reported an increase in accuracy of about 4.5% using subjective extract than full review on a movie review dataset. Matsumoto et al. applied sub-sequence mining, dependency parsing techniques on the same dataset and reported an increase in accuracy of about 7% from Pang et al. accuracy values (Matsumoto et al. 2005). However, their approach used linguistic resources like clausal extraction tools, POS tagger and dependency parsing.

Thet et al. (Thet et al. 2008) conducted experiments on aspect level sentiment classification of movie reviews. They used information extraction techniques like entity extraction, co-referencing and pronoun resolution to segment the text into sections, where a particular section focuses on particular aspect of a product (movie). They predicted the sentiments of people towards crew (casts,directors) and overall sentiment of a review. Baccianella et al. used lexicons like General Inquirer and POS patterns for extracting subjective features in a document (Baccianella et al. 2009). They used ϵ -SVR regression method to predict the polarity. They proposed a new feature selection method called minimum variance (MV) to select relevant features. They also implemented a round robin way of selecting features to minimize of the impact of skewed dataset in their experiments. They conducted experiments on hotel reviews and graded them on an ordinal scale of one to five (starred rating).

There were approaches in the literature that minimize or did not use any resource in sentiment classification. Cui et al. used n-gram model to represent each document as a feature vector and compared the performance of different classifiers (Cui et al. 2006). Their claim was that the existing work in sentiment classification focused on small set of documents. They wanted to test the performance of different sentiment classifiers on a large dataset with statistical n-gram models. Hu et al. used information retrieval techniques based

on language model to predict the orientation (Hu et al. 2007). They used a combination of Kullback-Leibler divergence and different smoothing techniques to predict the polarity.

Raychev and Nakov used a different weighting scheme based on word position and its subjectivity for naive-bayes classifier (Raychev and Nakov 2009). They conducted experiments on standard IMDB movie review dataset and reported an accuracy of 89.5% using their modified weighting scheme. In 2010, Li et al. proposed an approach that combines polarity shifting of sentences and document polarity classification (Li et al. 2010). Polarity shifting means that the polarity of a sentence is different from the polarity expressed by the sum of the content words in it⁷. Their motive was that default polarity classification techniques could not capture subtle variations in polarity of each sentence. Hence, they followed a meta classifier approach where they estimate the shift in the polarity of each sentence and then predict the overall polarity of a document. They conducted experiments on different datasets with minimal use of linguistic resources.

There are also other topics of research in sentiment analysis that have become popular recently. Generally sentiment analysis tasks are highly domain dependent. Classifier trained on one domain may not perform accurately on the other (Tan et al. 2009; Aue and Gammon 2005). Hence, researchers are carrying out experiments to adapt classifiers for cross domain sentiment classification. There are also other tasks like utility of reviews, review spam detection etc.

3 RSUMM

Our subjective feature extraction process occur in two steps: In the first step, we propose a simple and statistical methodology called RSUMM for estimating subjectivity of each sentence in a document. We obtain a subjective excerpt of a document by discarding the sentences that are objective. We represent the document and its subjective excerpt using n-gram model. In the second step, we apply well known feature selection methods on the resultant n-grams to obtain the final feature set. In this work, we focus more on estimating subjectivity and obtaining the subjective excerpt of a document.

Our RSUMM is similar to vector space modeling techniques used in information retrieval. We view each document d as a mixture of subjective and objective information where the former convey the feelings of the author on a particular topic and the latter are the facts. We represent each sentence $s \in d$ as a vector of terms, $\bar{s} := (t_{s,1}, t_{s,2}, \dots, t_{s,n})$. We define two metrics Average Document Frequency (ADF) and Average Subjective Measure (ASM) and derive the corresponding term vectors \overline{adf} and \overline{asm} for the respective collections. We then compute the lexical similarity between \bar{s} and the vectors \overline{adf} , \overline{asm} , and estimate its subjectivity using Jaccard similarity measure. We retain the top X% of sentences from each document d as its *subjective excerpt*.

3.1 Vector Space Model

Vector space model assigns weights to index terms (Baeza-Yates and Ribeiro-Neto 1999). It is widely used in information retrieval to determine the relevance of a document for a given query. Both document and the query are represented as weighted vectors of terms and these weights are used to compute the degree of similarity between the query and a document.

⁷ the presence of negation, contrast transition etc.

Higher the similarity degree, more relevant is a document to the query.

Formal Definition: Both query q and document d are represented as a weighted vector of terms. The query vector is defined as $q := (w_{1,q}, w_{2,q}, \dots, w_{t,q})$ and the document vector $d := (w_{1,d}, w_{2,d}, \dots, w_{t,d})$ where t is the total number of index terms.

Then, the degree of similarity between the document d and the query q is the correlation between the two vectors. The correlation is quantified by a variety of similarity measures, for instance by the *cosine* of the angle between the two vectors. The weighting measure used typically in vector space model is *tfidf*.

$$TFIDF(t, C) = tf(t, d) \times \log \left(\frac{N}{n} \right) \quad (1)$$

where $tf(t, d)$ denote the frequency of the term in the given document d , N denotes the total number of documents in collection C , and n denotes the number of documents containing term t in C .

$$\begin{aligned} \cos\theta &= \frac{\vec{d} \cdot \vec{q}}{\|\vec{d}\| \|\vec{q}\|} \\ &= \frac{\sum_{i=1}^t w_{i,d} \times w_{i,q}}{\sqrt{\sum_{i=1}^t w_{i,d}^2} \times \sqrt{\sum_{i=1}^t w_{i,q}^2}} \end{aligned} \quad (2)$$

3.1.1 Average Document Frequency (ADF)

Document frequency is a widely used statistical measure in information retrieval to determine the importance of a term in a given collection. We use ADF metric to represent the collection C_{pol} ⁸ as a vector of terms, $adf := (t_{adf,1}, t_{adf,2}, \dots, t_{adf,n})$. Each term $t_{adf} \in adf$ has a document frequency greater than the average document frequency of the collection. This metric intuitively selects the most important features from a given collection and features specific to its domain.

$$ADF(C_{pol}) = \frac{\sum_{i=1}^{|V_{pol}|} df(t_i, c_{pol})}{|V_{pol}|} \quad (3)$$

where $ADF(C_{pol})$ denotes the average document frequency of the collection C_{pol} , $|V_{pol}|$ denote the total number of terms present in C_{pol} and $df(t_i, c_{pol})$ denotes the document frequency of term t_i in the collection C_{pol} .

3.1.2 Average Subjective Measure (ASM)

We use ASM metric to represent a collection C_{sub} ⁹ as a vector of terms. Each term in asm has subjective measure greater than the average subjective measure of the collection, $asm := (t_{asm,1}, t_{asm,2}, \dots, t_{asm,n})$. This metric intuitively selects the most subjective features in a given annotated collection.

⁸ An annotated collection of documents

⁹ An annotated collection of subjective and objective sentences

The subjective measure of a term $t_i \in C_{sub}$, is calculated as follows:

$$\Phi(t_i, C_{sub}) = \frac{subj(t_i, C_{sub})}{obj(t_i, C_{sub}) + tot(C_{sub})} \quad (4)$$

where $\Phi(t_i, C_{sub})$ denotes the subjective measure of $t_i \in C_{sub}$, $subj(t_i, C_{sub})$ denotes the frequency of term t_i in subjective instances of annotated collection C_{sub} , $obj(t_i, C_{sub})$ denotes the frequency of term t_i in objective instances of annotated collection C_{sub} (penalizing factor) and $tot(C_{sub})$ denote the total number of instances in C_{sub} (normalizing factor).

The Average Subjective Measure of a collection is calculated as follows:

$$ASM(C_{sub}) = \frac{\sum_{i=1}^{|V_{sub}|} \Phi(t_i, C_{sub})}{|V_{sub}|} \quad (5)$$

where $ASM(C_{sub})$ is the average subjective measure of the collection C_{sub} , $|V_{sub}|$ denotes the total number of terms in C_{sub} and $\Phi(t_i, C_{sub})$ denotes the subjective measure of term $t_i \in C_{sub}$ calculated using eqn.(4).

3.1.3 Final Scoring

After obtaining both *adf* and *asm* vectors using eqn.(3) and eqn.(5) respectively, we compute lexical similarity between each sentence vector s and the vectors *adf* and *asm*. In our work, we use raw terms rather than their weights. We use Jaccard similarity measure to compute the similarity between two given vectors as shown in eqn.(6)

$$\sigma(\bar{a}, \bar{b}) = \frac{n(\bar{a} \cap \bar{b})}{n(\bar{a} \cup \bar{b})} \quad (6)$$

where $\sigma(\bar{a}, \bar{b})$ denotes the similarity score between two vectors a and b , $n(\bar{a} \cap \bar{b})$ denotes the number of terms overlapping between vectors a and b and $n(\bar{a} \cup \bar{b})$ denotes the total number of terms in both vectors.

The final score of a sentence s , ($FS(s)$) is a combination of lexical similarity scores as shown in eqn.(7). We rank sentences in document d in decreasing order of their $FS(s)$ and retain top X% of them. The first part of eqn.(7) computes lexical similarity between a sentence and the most important features in a collection. Hence, the most informative sentence in a document gets high priority. The latter part of it computes the lexical similarity between a sentence and the most subjective terms. So, sentences that are more subjective are ranked higher.

$$FS(s) = \sigma(\overline{adf}, \bar{s}) + \sigma(\overline{asm}, \bar{s}) \quad (7)$$

Our RSUMM ensures that sentiment in a document is preserved to the maximum extent by optimizing 'X'. Thus, we have a *subjective excerpt* of d discarding objective sentences for effective sentiment classification. We use n-gram model to represent the subjective excerpt of a document as a feature vector to the classifier.

4 Feature Selection

RSUMM extracts most subjective sentences from a document. As sentences are relatively larger text units compared to words or phrases, using n-gram model for converting them as feature vectors lead to a very high dimensional feature set. For faster learning and better classification accuracies, we have to reduce the dimensionality by selecting features that are more relevant and capable of discriminating the class variable. Hence, a feature selection phase is essential in our case. We apply two state-of-the-art feature selection methods that are proven effective in text categorization and sentiment classification: mutual information (MI) and fisher discriminant ratio (FDR) to select final subjective features from a document (Yang and Pedersen 1997; Wang et al. 2009).

4.1 Mutual Information

Mutual Information (MI) is a widely used feature selection method in text categorization. It computes the mutual dependence between a feature f and class C . It measures the amount of information the presence/absence of a feature contributes in making the correct classification decision on C . In our case, the feature f is an n-gram and class C is either positive or negative.

$$MI(f; C) = P(f, C) \log \left(\frac{P(f, C)}{P(f)P(C)} \right) \quad (8)$$

where $MI(f; C)$ denotes the mutual information between the feature f and class C , $P(f, C)$ denotes the conditional probability of feature occurring in class C , $P(f)$ denotes the probability of feature in the entire collection and $P(C)$ denotes the class probability.

4.2 Fisher Discriminant Ratio

Fisher discriminant ratio (FDR) is one of the effective approaches for dimensionality reduction in pattern recognition. The main idea of FDR is that the points in a D dimensional space are projected in a way that there is a maximum difference between the means and minimum variance within each class. For a two class classification problem, FDR can be computed as follows.

$$J(w) = \frac{|m_1 - m_2|^2}{S_1^2 + S_2^2} \quad (9)$$

where m_i denotes a mean, S_i denotes the within-class variance and $i = 1, 2$.

We modify the above equation in a way that it computes the discriminating power of a feature f . Let $d_{P,i}$ ($i = 1, 2, \dots, m$) and $d_{N,j}$ ($j = 1, 2, \dots, n$) denote the i^{th} positive document and j^{th} negative document respectively. We define two random variables $d_{P,i}(f)$ and $d_{N,j}(f)$ as in (Wang et al. 2009).

$$d_{P,i}(f) = \begin{cases} 1 & \text{if } f \text{ occurs in } d_{P,i} \\ 0 & \text{otherwise} \end{cases}$$

$$d_{N,j}(f) = \begin{cases} 1 & \text{if } f \text{ occurs in } d_{N,j} \\ 0 & \text{otherwise} \end{cases}$$

The modified version of equation 9 is written as follows.

$$FDR(f) = \frac{\left(\frac{m_1}{m} - \frac{n_1}{n}\right)^2}{\frac{\sum_{i=1}^m (d_{P,i}(f) - \frac{m_1}{m})^2}{m} + \frac{\sum_{j=1}^n (d_{N,j}(f) - \frac{n_1}{n})^2}{n}} \quad (10)$$

where m and n denote the total number of documents in class P and N respectively, m_1 and n_1 denote the number of instances of feature f in class P and N , $d_{P,i}(f)$ denotes the presence or absence of feature f in review i of class P and $d_{N,j}(f)$ denotes the presence or absence of feature f in review j of class N

4.3 Final Subset Selection

In each case, high MI and FDR values of f implies that it has more discriminative power. We sort features in decreasing order of their corresponding MI and FDR values. In MI, we use a round robin way of selecting features as we compute the mutual information of f with each class variable C . In FDR, the estimate is for entire collection and not per each class. We retain top $Y\%$ of n -grams for each method in the final feature set.

5 Experiments

We conduct experiments on a movie review dataset collected from IMDB archive. We use n -gram model to represent the total review and its subjective excerpt as a feature vector to the classifier. Each n -gram is weighted using *tfidf* score. We predict the overall polarity of a review as positive or negative.

5.1 Experimental Setup

5.1.1 Datasets

We have downloaded the available IMDB archive¹⁰ of the rec.arts.movies.reviews newsgroup¹¹. It contains 27,886 unprocessed and unlabeled html files that convey opinions of different authors on different movies. Predominantly, it has reviews rated on three different scales: zero to four, zero to five and grade F to A+. A set of rules are framed to mine the rating patterns from the unprocessed set. Following these rules, we annotate each review as positive or negative.

With respect to rating scale of zero to five, reviews having a rating of three and a half (**1/2) or above are annotated as positive. The reviews with a rating of two and a half (*1/2) and below are annotated as negative. With respect to rating scale of zero to four, we annotate reviews with rating of three(**) or above as positive and reviews with rating of one and a half(*1/2) and below as negative. We annotate movie reviews that are graded B or above as positive (B, B+, A-, A, A+) and C- or below (C-, D, D-, F) as negative. Using this annotation scheme, we have annotated 7700 reviews as positive and 3600 reviews as negative.

¹⁰ http://www.cs.cornell.edu/people/pabo/movie-review-data/polarity_html.zip

¹¹ <http://reviews.imdb.com/Reviews>

We have selected about 11,300 reviews from the archive for our experiments. It is comparatively a larger movie review dataset to others used in sentiment classification (Pang et al. 2002; Mullen and Collier 2004; Pang and Lee 2004; Matsumoto et al. 2005). We call this dataset as PDS in the remainder of this paper. We split PDS into ten equal folds and perform ten-fold cross validation to test the statistical significance of our results. We maintain uniform class distribution across each fold.

We use a movie review dataset as a validation set to estimate parameters. It contains 1400 reviews with equal class distribution and each review labeled as positive or negative. We call it as VDS¹² in the remainder of this paper (Pang et al. 2002). To populate the knowledge of subjectivity into our experiments, we use an annotated collection of 5000 subjective and objective sentences respectively. We call this dataset as SDS¹³. (Pang and Lee 2004)

5.1.2 Preprocessing

We extract body text from unprocessed html files in PDS. We have framed a set of rules to extract the body text. Since, all unprocessed files are from the same source (IMDB), framing rule based patterns for extracting it is not that difficult although we may have a little noise. We use the sentence breaker and the tokenizer implemented in openNLP¹⁴ tool to split the extracted body text into sentences and sentences into words respectively. Except for preprocessing, we do not use any linguistic resource in our experiments. In subjective feature extraction, we rely on statistical measures as described above. Hence, we are reducing resource dependency or usage of complex patterns in subjective feature extraction.

5.1.3 Classifier and Evaluation

We use SVM classifier implemented in SVMLight package¹⁵ in our experiments with parameters set to their default values. We focus more on extracting subjective features and representing them as feature vectors to the classifier rather than tuning its parameters.

We use accuracy of the classifier on a test set as the evaluation metric. The accuracy of the classifier is calculated as shown in eqn. 11.

$$ACC = \frac{t}{n} \times 100 \quad (11)$$

where t denote the number of samples correctly classified and n denote the total number of test samples.

Since PDS is skewed towards the positive class, we modify the above equation slightly to include per class accuracy to give best generalization for our results. The modified equation for a two class problem is shown in eqn. 12 (Baccianella et al. 2009).

$$ACC = \frac{\frac{t_P}{n_P} \times 100 + \frac{t_N}{n_N} \times 100}{n_C} \quad (12)$$

where t_P denote the number of test samples correctly classified as positive (P), n_P denote the total number of test samples with label positive, t_N denote the number of test samples correctly classified as negative (N), n_N denote the total number of test samples with label negative and n_C denote the number of classes present in the dataset (in our case n_C is two).

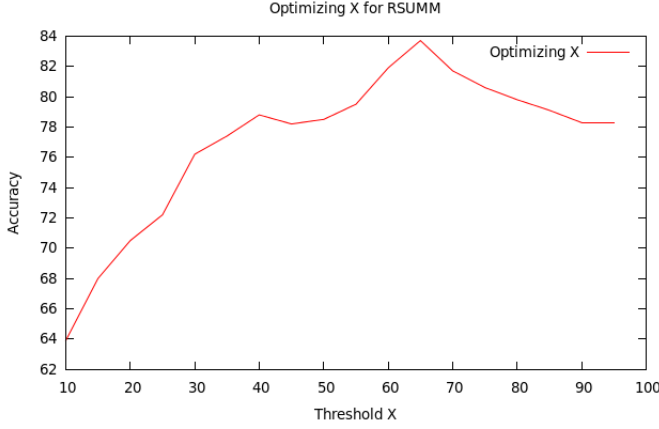
¹² http://www.cs.cornell.edu/people/pabo/movie-review-data/mix20_rand700_tokens_cleaned.zip

¹³ http://www.cs.cornell.edu/people/pabo/movie-review-data/rotten_imdb.tar.gz

¹⁴ <http://opennlp.sourceforge.net/>

¹⁵ <http://svmlight.joachims.org/>

Fig. 1: Plot showing the effect of accuracy with increase in 'X' for RSUMM on VDS



5.2 Estimating parameters

We estimate optimal values for parameters X and Y on VDS and later apply them on PDS. We use the dataset VDS as C_{pol} in eqn. 3 and SDS as C_{sub} in eqn. 4. We derive the corresponding term vectors \overline{adf} and \overline{asm} and use RSUMM to estimate the subjectivity of each sentence \bar{s} in a review as described in Section. 3.

We start with a very low value of X=10, (i.e.) retaining top 10% of sentences in each review as its subjective excerpt. We then incrementally add 5% in each iteration and examine the increase or decrease in the performance of the classifier. We use unigram representation of the subjective excerpt as a feature vector to the classifier. We perform ten-fold cross validation to test statistical significance of our results maintaining class uniform distribution. We choose the value of X that produces best accuracy on VDS as its optimal value. The effect of accuracy with increase in X is depicted pictorially in Figure. 1.

From the figure, it is clear that using an excerpt of a review we are able to achieve better accuracies than using the total content. Using the total review, we are able to achieve an accuracy of 76.1% with unigrams as features. But, retaining only 65% of it, the increase in the performance of the sentiment classifier is around 7%. It validates our assumption that entire review cannot be subjective and it is a mixture of subjective and objective information. We use the same value of 65 for X on PDS in our subsequent experiments.

For estimating Y, we start with a very high value of Y being 95, (i.e) placing 95% of features in the final feature set for each method. We then, decrement Y by 5% for every iteration. The optimal value of Y for MI and FDR is the one that produces the best accuracy value with the respective n-gram model. The experimental results with unigrams (n=1) as features for MI and FDR on the total review are shown in Figure 2. The effect on accuracy with change in Y for MI and FDR with bigrams (n=2) as features is shown in Figure 3.

From the figures 2 and 3, it is evident that very low threshold values for Y produce accurate results in each method for a given n-gram model. The optimal value of Y in each case is 5. It conveys that very few n-grams in the entire text are indeed subjective. There is

Fig. 2: Tuning the parameter γ for MI and FDR on RSUMM with unigrams as features

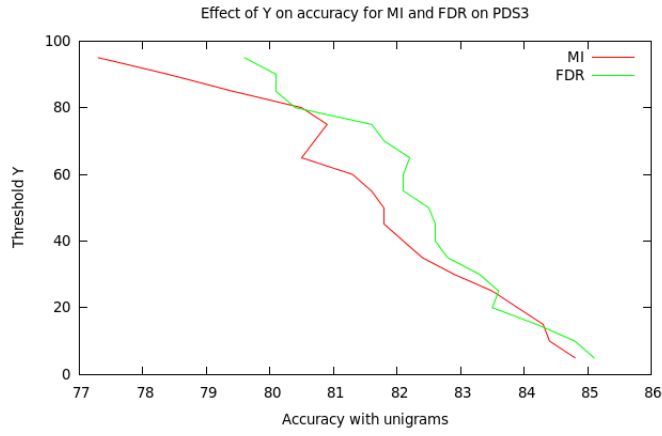
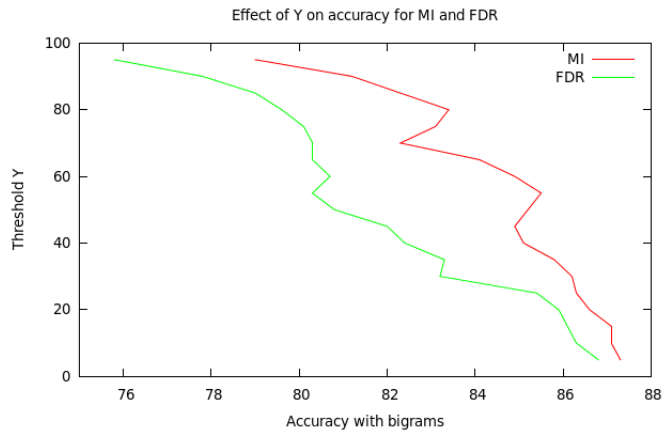


Fig. 3: Tuning the parameter γ for MI and FDR on RSUMM with bigrams as features



a lot of noise or misleading text surrounding the subjective features. We have to discard it for efficient sentiment classification as proved in the literature (Pang and Lee 2004). In this work, we limit ourselves to using unigrams, bigrams and their combination as feature vector representation of the text.

Table 1: Results showing CV accuracies for BL, TH and BH on PDS

| Features | BL | TH | BH |
|----------|------|------|------|
| Uni | 77.9 | 74.4 | 77.1 |
| Bi | 63.5 | 58.4 | 60.5 |
| Uni+Bi | 74.8 | 70.9 | 74.5 |

Table 2: Results showing CV accuracies for BL+MI and BL+FDR on PDS

| Features | BL+MI | BL+FDR |
|----------|-------|--------|
| Uni | 80.4 | 81.7 |
| Bi | 75.6 | 74.6 |
| Uni+Bi | 79.9 | 80.4 |

Table 3: Results showing CV accuracies for BL+DSE, BL+DSE+MI and BL+DSE+FDR on PDS

| Features | BL+DSE | BL+DSE+MI | BL+DSE+FDR |
|----------|--------|-----------|------------|
| Uni | 79.4 | 80.8 | 81.5 |
| Bi | 64.2 | 80.1 | 79.0 |
| Uni+Bi | 77.3 | 81.4 | 81.2 |

5.3 Experimental Results

The baseline (BL) in our experiments is using the total review with unigrams (Uni), bigrams (Bi) and their combination (Uni+Bi) as features. We split each review into two equal halves and carry out experiments using the top half (TH) and bottom half (BH) of the review. This is done to test the general pattern followed by authors in expressing their sentiment. The general pattern in movie review domain is that authors discuss objective information such as plot, characters and other aspects of a movie at the beginning. They convey their sentiments at the end of a review (Pang and Lee 2004). We report the accuracy values of BL, TH, and BH on PDS respectively in Table. 1.

We carry two more experiments to test the relative significance of RSUMM. Firstly, we train a naive-bayes classifier on SDS and test it on each sentence of a review in PDS. We discard sentences labeled as objective by the naive-bayes classifier and retain the subjective sentences in subjective excerpt. We call this method of extracting subjectivity as default subjective excerpt (DSE). We represent the subjective excerpt using n-grams and then apply MI and FDR on the resultant n-grams to obtain the final feature set. Table 3 reports the accuracy values using DSE and the effect of applying MI and FDR on it. Secondly, we apply MI and FDR on the baseline to verify whether the performance of the sentiment classifier is more sensitive to Y than X. The results for this experiment are reported in Table 2.

We use the collection PDS as C_{pol} in eqn. 3 and derive the respective \overline{adf} term vector. The \overline{asm} vector is derived using the collection SDS as C_{sub} in eqn. 5. After obtaining the corresponding \overline{adf} and \overline{asm} vectors, we compute the lexical similarity score between each sentence \bar{s} of a review and \overline{asm} , \overline{adf} vectors as per eqn. 6. We use a combination of two similarity scores to obtain the final subjective score of a sentence as shown in eqn. 7. We

Table 4: Results showing CV accuracies for BL+RSUMM, BL+RSUMM+MI and BL+RSUMM+FDR on PDS

| Features | BL+RSUMM | BL+RSUMM+MI | BL+RSUMM+FDR |
|----------|----------|-------------|--------------|
| Uni | 83.9 | 86.1 | 86.7 |
| Bi | 69.3 | 83.9 | 84.5 |
| Uni+Bi | 80.1 | 85.7 | 84.9 |

retain top X% of sentences of a document in its subjective excerpt. We choose the value of X to be 65. We use n-gram model ($n \leq 2$) to represent the subjective excerpt as a feature vector to the classifier. We then apply MI and FDR on the resultant unigrams and bigrams with Y being 5 in each case to obtain the final feature set. The accuracy values for RSUMM, MI and FDR are reported in Table. 4.

6 Discussion

The baseline accuracy values reported in Table 1 indicate that unigram feature vector representation of a review yield better results compared to bigrams or combination of unigrams and bigrams. Using the total content with unigrams as features, the baseline accuracy reported is 77.9%. In case of bigrams and combination of unigrams and bigrams, the accuracy values are 63.5% and 74.8% respectively. The bottom half of review in movie review domain is more subjective compared to the top half with each feature vector representation of the text. It conveys that authors model subjective expressions after discussing the plot, establishing the characters etc. in movie review domain. Using only bottom half a review with unigram representation of feature vectors, we are able to achieve a comparable accuracy to the baseline of the sentiment classifier (77.9% to 77.1%). In case of bigrams and there is a drop of about 3% in accuracy value from the baseline. But, the drop in accuracy is very small when the combination of unigrams and bigrams are used as features. It adds robustness to our assumption that entire review does not contain subjective information.

Using only top half of a review, there is a drop of about 3.5%, 5% and 5% from the baseline in case of unigrams, bigrams and their combination respectively. There is a drop in the performance of the classifier when bigrams and combination of unigrams and bigrams are used as features than unigrams. It conveys that too many features will degrade the performance of the sentiment classifier and there has to be selection criteria. Hence, researchers have defined some support thresholds in their experiments for selecting features (Pang et al. 2002; Matsumoto et al. 2005). From the results in Table 1 we conclude that unigram representation yield better results compared to bigrams and combination of unigrams and bigrams.

To test the relative significance of RSUMM, we use DSE method to obtain the subjective excerpt of a document. The increase in the performance of the classifier is around 1.5%, 0.7% and 2.5% in case of unigrams, bigrams and their combination respectively. Using RSUMM, there is a significant increase in the performance of the classifier from the baseline with unigrams as features (around 6%). In case of bigrams and combination of unigrams and bigrams, the increase is 5.8% and 5.3% respectively. It conveys that RSUMM is able to optimize the essence of sentiment in a document in threshold X. In comparison to DSE, RSUMM performed better although they use the same SDS to populate the knowledge of subjectivity. It can be attributed to the fact that we use an additional metric called ADF in

RSUMM that gives importance to most informative sentences in a document. We don't use any linguistic resource in our experiments yet obtain an increase in the performance of the sentiment classifier. These experimental results prove that subjective feature extraction is possible with minimum use of linguistic resources and no complex patterns.

Feature selection techniques have proved vital in the performance of several text categorization tasks as they enhance the performance of the classification system considerably (Yang and Pedersen 1997). Even in sentiment classification tasks, selecting features based on techniques like document frequency, term frequency, minimum variance etc. is done for obtaining good performance (Pang et al. 2002; Matsumoto et al. 2005; Baccianella et al. 2009). In our experiments, we employ two state-of-the-art feature selection methods MI and FDR. From tables 3 and 4, we can infer that there is a little to choose between the two as they both enhance the performance of the sentiment classifier. In case of unigrams, the increase in accuracy of RSUMM+FDR is less comparable to RSUMM (2.8% to 6.0%). In case of bigrams, the performance of the classifier is highly sensitive to parameter γ . It infers the presence of large number of irrelevant features in documents.

Using MI as the feature selection method, there is a marginal increase of about 2% from the baseline in case of unigrams whereas there is a significant increase of about 12% in accuracy value for bigrams. After applying MI with RSUMM, there is an increase of 2.2% and 14% with unigrams and bigrams respectively. Using the combination of DSE and MI, there is an increase of 1.4%, 16% and 4% in case of unigrams, bigrams and their combination respectively. Using RSUMM and MI, the accuracy values obtained for unigrams, bigrams and their combination are 86.1%, 83.9% and 85.7% respectively. The increase in their values from the baseline is about 8%, 20% and 11% with unigrams, bigrams and their combination as features respectively which is significant.

FDR as the feature selection technique has performed very similar to MI. In case of unigrams, it has increased the performance of the classifier compared to MI (81.7% to 80.4% and 86.7% and 84.5%). Using the combination of RSUMM and FDR, we obtain the highest accuracy of 86.7% and 84.5% with unigrams and bigrams respectively. There is an increase in the accuracy values of about 9%, 21%, 10% for each feature representation respectively from the baseline. There is a slight drop in the accuracy values when the combination unigrams and bigrams are used as features compared to unigrams and bigrams in isolation.

7 Conclusion and Future Work

We focused on subjective feature extraction, the key component in document sentiment classification. We followed a two-step filtering methodology to mine subjective portions in a document. We used RSUMM to obtain subjective excerpt of a document by estimating subjectivity at sentence level. We then applied well known feature selection methods on the subjective excerpt to obtain the final feature set. The major contributions of this work are:

1. We explored various subjective feature extraction methodologies in sentiment classification and their limitations.
2. We attempted to minimize resource dependency or complex patterns in subjective feature extraction.
3. We developed a simple and statistical methodology called RSUMM to extract subjectivity in a document.

We explored several resource independent or language independent subjective feature extraction methods in the literature. We used techniques similar to vector space model in

RSUMM to estimate the subjectivity of each sentence. To the best of our knowledge, usage of techniques similar to vector space model for extracting subjective features was not proposed yet in the literature. Our experimental results concluded that subjective feature extraction is possible with minimum use of linguistic resources. We explored two frequency based metrics ADF, ASM and used them in RSUMM to estimate subjectivity. We then applied two well known feature selection methods MI and FDR on the subjective excerpt of RSUMM.

Our work can be considered as a building block for analyzing sentiments with minimal usage of linguistic resources and no complex patterns. In future, we need to explore different metrics to extract subjectivity and conduct experiments. As feature selection methods are proved critical in the performance of the classification, we need to explore more novel methods for selecting features.

Acknowledgements We thank Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan for providing URL in their paper to download the IMDB movie review dataset. We thank the Department of Computer Science, Cornell University for providing us link to download the dump of IMDB archive.

References

- Shlomo Argamon, Moshe Koppel, and Galit Avneri. Routing documents according to style. In *Proceedings of First International Workshop on Innovative Information Systems*, 1998.
- A. Aue and M. Gamon. Customizing sentiment classifiers to new domains: A case study. In *Proceedings of the International Conference RANLP-2005.*, 2005.
- Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. Multi-facet rating of product reviews. In *Proceedings of the 31th European Conference on IR Research on Advances in Information Retrieval*, ECIR '09, pages 461–472. Springer-Verlag, 2009.
- Ricardo A. Baeza-Yates and Berthier Ribeiro-Neto. *Modern Information Retrieval*. Addison-Wesley Longman Publishing Co., Inc., 1999.
- Philip Beineke, Trevor Hastie, and Shivakumar Vaithyanathan. The sentimental factor: improving review classification via human-provided information. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, ACL '04. Association for Computational Linguistics, 2004.
- Hang Cui, Vibhu Mittal, and Mayur Datar. Comparative experiments on sentiment classification for online product reviews. In *proceedings of the 21st national conference on Artificial intelligence - Volume 2*, pages 1265–1270. AAAI Press, 2006.
- Urike Gretzel and Kyung Hyan Yoo. Use and impact of online travel reviews. *Information and Communication Technologies in Tourism*, pages 35–46, 2008.
- Vasileios Hatzivassiloglou and Kathleen R. McKeown. Predicting the semantic orientation of adjectives. In *Proceedings of the eighth conference on European chapter of the Association for Computational Linguistics*, pages 174–181. Association for Computational Linguistics, 1997.
- Minqing Hu and Bing Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '04, pages 168–177, 2004.
- Yi Hu, Ruzhan Lu, Xuening Li, Yuquan Chen, and Jianyong Duan. A language modeling approach to sentiment analysis. In *Proceedings of the 7th international conference on Computational Science, Part II*, ICCS '07, pages 1186–1193. Springer-Verlag, 2007.

- Brett Kessler, Geoffrey Numberg, and Hinrich Schütze. Automatic detection of text genre. In *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics and Eighth Conference of the European Chapter of the Association for Computational Linguistics*, ACL-35, pages 32–38. Association for Computational Linguistics, 1997.
- Shoushan Li, Sophia Y. M. Lee, Ying Chen, Chu-Ren Huang, and Guodong Zhou. Sentiment classification and polarity shifting. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 635–643, 2010.
- Bing Liu. Sentiment analysis and subjectivity. In *Handbook of Natural Language Processing, Second Edition*. CRC Press, Taylor and Francis Group, Boca Raton, FL, 2010.
- Shotaro Matsumoto, Hiroya Takamura, and Manabu Okumura. Sentiment classification using word sub-sequences and dependency sub-trees. In *Proceedings of PAKDD*, pages 301–311, 2005.
- Tony Mullen and Nigel Collier. Sentiment analysis using support vector machines with diverse information sources. In *Proceedings of EMNLP*, pages 412–418, 2004.
- Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the ACL*, pages 271–278, 2004.
- Bo Pang and Lillian Lee. Opinion mining and sentiment analysis. *Foundation and Trends in Information Retrieval*, 2:1–135, 2008. ISSN 1554-0669.
- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Thumbs up?: sentiment classification using machine learning techniques. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing - Volume 10*, pages 79–86, 2002.
- Veselin Raychev and Preslav Nakov. Language-independent sentiment analysis using subjectivity and positional information. In *Proceedings of the International Conference RANLP-2009*, pages 360–364. Association for Computational Linguistics, 2009.
- Songbo Tan, Xueqi Cheng, Yuefen Wang, and Hongbo Xu. Adapting naive bayes to domain adaptation for sentiment analysis. In *Proceedings of the 31th European Conference on IR Research on Advances in Information Retrieval, ECIR '09*, pages 337–349. Springer-Verlag, 2009.
- Tun Thura Thet, Jin-Cheon Na, and Christopher S. Khoo. Sentiment classification of movie reviews using multiple perspectives. In *Proceedings of the 11th International Conference on Asian Digital Libraries: Universal and Ubiquitous Access to Information, ICADL 08*, pages 184–193. Springer-Verlag, 2008.
- Peter D. Turney. Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, pages 417–424. Association for Computational Linguistics, 2002.
- Suge Wang, Deyu Li, Yingjie Wei, and Hongxia Li. A feature selection method based on fisher’s discriminant ratio for text sentiment classification. In *Proceedings of the International Conference on Web Information Systems and Mining, WISM '09*, pages 88–97. Springer-Verlag, 2009.
- Casey Whitelaw, Navendu Garg, and Shlomo Argamon. Using appraisal groups for sentiment analysis. In *Proceedings of the 14th ACM international conference on Information and knowledge management, CIKM '05*, pages 625–631. ACM, 2005.
- Yiming Yang and Jan O. Pedersen. A comparative study on feature selection in text categorization. In *Proceedings of the Fourteenth International Conference on Machine Learning, ICML '97*, pages 412–420. Morgan Kaufmann Publishers Inc., 1997.