

A personalized hybrid movie recommendation system for users

A dissertation submitted in the partial fulfillment for the award of Degree of

Master of Technology

In

Software Technology

Submitted by

Yogendra Singh (2K13/SWT/19)

Under the esteemed guidance of

Dr. Rajni Jindal

(Associate Professor, Department of Computer Science and Engineering)



DELHI TECHNOLOGICAL UNIVERSITY

BAWANA ROAD, DELHI

2013-2016

DECLARATION

I hereby declare that the thesis entitled “**A personalized hybrid movie recommendation system for users**” being submitted to the **Delhi Technological University**, in partial fulfillment of the requirements for the award of degree of **Master of Technology in Software Technology** is an authentic work carried out by me under the guidance of Dr. Rajni Jindal. The material contained in this thesis has not been submitted to any university or institution for the award of any degree.

Yogendra Singh

Department of Software Engineering

Delhi Technological University,

Delhi.

CERTIFICATE



Date: _____

This is to certify that the Major Project entitled “**A personalized hybrid movie recommendation system for users**” submitted by **Yogendra Singh** (2K13/SWT/19); in partial fulfillment of the requirement for the award of degree of Master of Technology in Software Technology to Delhi Technological University, Bawana Road Delhi; is a record of the candidate's own work carried out by him under my supervision.

Dr. Rajni Jindal

Associate Professor, Department of Computer Science and Engineering

Delhi Technological University

Bawana road, Delhi - 110042

ACKNOWLEDGEMENT



July 2016

I would like to take this opportunity to thank my project guide Dr. Rajni Jindal for her invaluable and consistent guidance throughout this work. I would like to thank her for giving me the opportunity to undertake this topic. I am very appreciative of her generosity with her time, advice, data, and references, to name a few of her contributions. It is her wonderful association that enabled me to achieve the objectives of this work. I humbly extend my grateful appreciation to my friends whose moral support made this study possible.

Lastly, I would like to thank all the people directly and indirectly involved in successfully completion of this project.

Yogendra Singh

Roll # 2K13/SWT/19

Master of Technology (Software Technology)

Delhi Technological University

Bawana road, Delhi – 110042

ABSTRACT

We describe a rating logical thinking approach to incorporating matter user reviews into Collaborative Filtering (CF) algorithms. The main motive of our approach is to use user preferences which is expressed in movie reviews and then convert such user's preferences into some rating that may be understood by existing CF algorithms. The linguistics score of subjective sentence is fetched from SentiWordNet Library to calculate their sentiments as +ve, -ve or neutral based on the textual review. We've used SentiWordNet library as a dataset with two completely different approaches of alternatives comprising of adverbs and verbs, adjectives and n-gram feature extraction. We have a tendency to conjointly used our SentiWordNet library to figure the document level sentiment for every movie reviewed and compared its label with results obtained victimization Alchemy API. We conjointly developed and evaluated a model of the planned framework. Preliminary results valid the effectiveness of varied tasks within the planned framework, and recommend that the framework doesn't admit an oversized coaching corpus to operate. Additional development of our rating logical thinking framework is in progress. A comprehensive analysis of the framework are administered and reported during a follow-up article.

Contents

CERTIFICATE	iii
ACKNOWLEDGEMENT	iv
ABSTRACT	v
List of Figures	viii
List of tables	viii
ABBREVIATIONS	viii
Chapter 1: INTRODUCTION.....	9
Chapter 2: LITERATURE REVIEW	13
Chapter 3: SENTIMENT CLASSIFICATION	21
3.1 Opinion Lexicons.....	22
3.2 WordNet Glosses and SentiWordNet	22
3.3 Web Crawling using HTML AGILITY PACK	24
Chapter 4: THE PROPOSED FRAMEWORK.....	26
4.1 Data Extraction	27
4.2 Analyzing reviews	27
4.2 Rating Inference:.....	32
Chapter 5: SOLUTION APPROACH.....	34
5.1 Review extraction process	34
5.2 Imdb.com Crawler	35
5.3 Sentiment Analysis	36
5.4 Web User Interface	37
Chapter 6 : DATASET AND PERFORMANCE MEASURES	38
6.1 Dataset.....	38
6.2 Analysis on the Use of Opinion Words	38

Chapter 7: CONCLUSION AND FUTURE WORK.....	42
<i>References</i>	43

List of Figures

Figure 1: Overview of framework	26
Figure 2 Pseudo code: SWN(AAC)	30
Figure 3 Pseudo code: SWN(AAAVC)	32
Figure 4 Overall process Model.....	34
Figure 5 General task of search engine.....	35
Figure 6 : Database structure (Tables ER).....	36

List of tables

Table 1 Coverage of Opinion Lexicons	24
Table 2 SENTIWORDNET DATABASE STRUCTURE	30
Table 3: Top 10 opinion words with relative frequencies.	39

ABBREVIATIONS

CF.....	Collaborative Filtering
IMDB.....	Internet Movie Database
PHOAKS.....	People Helping One Another Know Stuff
AAC.....	Adverb + Adjective Combinations
POS.....	Part of speech
PosScore.....	Positive Score
NegScore.....	Negative Score
AAAVC.....	Adverb + Adjective and Adverb + Verb Combine.
+ve.....	Positive
-ve.....	Negative
kNN.....	K Nearest Neighbor

Chapter 1: INTRODUCTION

Collaborative Filtering (CF) is a methodology in recommendation systems. It provides personalized recommendations to users supported a user preferences and similarity, from that users having similar tastes. It then recommends to a target user things liked by his, similar users [5,10].

CF-based recommendation systems may be classified into 2 major category which depends upon how these system collect user preferences:

- a) User-log based and
- b) Ratings based.

User-log based mostly CF obtains user preferences from implicit votes which are captured through users' interactions with the system[12]. Ratings based mostly CF makes use of express ratings users have given things (e.g. 5-star rating scale as in MovieLens [6]). Such ratings area unit sometimes in or will simply be reworked into numerical values (e.g. A to E). Some review hubs, like the net picture show information (IMDb), permit users to supply comments in free text format, spoken as user reviews during this work. User reviews may also be thought-about a kind of "user ratings", though they're sometimes language texts instead of numerical values. An analysis on mining user preferences from reviews, can be referred as sentiment analysis (e.g. [7,17,18, 23]), is getting more and more popular within the text mining area, its integration with CF has solely received very little analysis attention. The PHOAKS (People Helping One Another Know Stuf) system represented in [21] classifies websites suggested by users in newsgroup messages, however it doesn't involve mining user preferences from texts.

This work describes our projected framework for integration sentiment analysis and CF. we have a tendency to take a rating reasoning approach that infers numerical ratings from matter reviews, in order that user preferences delineated within the reviews will simply be fed into existing CF algorithms. The contributions of this approach area unit two-fold.

- a) Firstly, it addresses the well-known knowledge poorness drawback in CF by permitting CF algorithms to use matter reviews as a further supply of user preferences.
- b) Secondly, it allows extending CF to domains wherever numerical ratings on merchandise area unit tough to gather, or wherever preferences on domain things area unit too advanced to be expressed as scalar ratings. An example of these type domains is travel and commercial enterprise, which is most existing recommendation systems area unit engineered upon content-based or knowledge-based techniques [20]. “Reviews” written by travelers on commercial enterprise merchandise or destinations, as an example, area unit offered as travel Journals on TravelJournals.net [22]. Interpreting travelers’ preferences from their travel reviews might contribute to the event of additional advanced and personalized recommendation systems.

Sentiment Analysis is a natural language processing techniques that uses an approach to find textual content and categorize it as +ve, -ve or neutral. The unstructured matter knowledge on the net usually carries expression of opinions of users. Sentiment analysis tries to spot the expressions of opinion and mood of writers. A sentiment analysis classifies a document as '+ve', '-ve' or 'neutral', supported the opinion expressed in it. The drawback of document level sentiment analysis is basically as follows: Given a collection of documents [S], a sentiment analysis rule classifies every document [s S] into one among the 2 categories, +ve and -ve. +ve label denotes that the document d expresses a positive opinion and -ve label means d expresses a negative opinion of the user. Additional refined algorithms try and establish the sentiment at movie's feature-level,sentence-level, or entity-level.

There are mainly three types of approaches for sentiment classification of texts:

1. By using a machine learning based text classifier -such as Naïve Thomas Bayes, SVM or kNN- with appropriate feature choice theme;
2. By using the unsupervised semantic orientation scheme of extracting relevant n-grams of the text so treated them either as +ve or -ve and consequentially the document; and
3. By using the SentiWordNet open-source used based online library that gives positive, negative and neutral scores for words. A number of the relevant past works on sentiment classification is found in [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11] and [12].

Now a day's web of internet hosts an outsized volume of information created by numerous users. Users are currently co-creators of website, instead of being passive customers. The social media is currently a serious a part of the internet. The statistics shows that each four out of five users on the net use some sort of social media. The user contributions to social media vary from blog posts, tweets, reviews and photo or video uploads etc. an outsized quantity of the information on the net is unstructured text. Opinions expressed in social media in sort of reviews or posts represent a very important and attention-grabbing space price exploration and exploitation. With increase in accessibility of opinion resource product reviews, movie reviews, blog reviews, social network tweets, the new difficult task is to mine giant volume of texts and devise appropriate algorithms to know the opinion of others. This info is of very useful and information to firms that try and grasp the feedback regarding their product or services. This review helps them in taking user choices. Additionally to be helpful for firms, the reviews and opinion strip-mined from them, is useful for users in addition. reviews about hotels in the city may help a user going to a city seeking for a good hotel. Similarly, movies' reviews help other users to decide whether the movie is worth to purchase or not. Similarly, movie reviews facilitate different users choose whether or not the movies' is worth for money

or not. During this work we have got tried to explore a new SentiWordNet primarily based theme for each document-level. The document level class involves use of various linguistic options (ranging from Adverb + Adjective combination to Adverb + Adjective + Verb combination). We have got additionally devised a new domain specific heuristic for aspect-level sentiment classification of movies' reviews. This theme locates the self-opinionated text round the desired aspect feature in an exceedingly review and computes its sentiment orientation. For a movies', this is used for all the reviews.

The sentiment scores on a particular aspect from all the selected reviews are then aggregated together. Same process is used for all aspect which are under consideration. Finally a summarized sentiment profile of the users' comment on all aspects is presented in an easy to visualize and understandable pictorial form. The remaining report is organized as follow.

Chapter 2 describes the background literature survey conducted for this report. Chapter 3 describes the sentiment level classification and detail about wordnet and sentiwordnet Database. Chapter 4 and 5 describes the proposed framework and solution approach in used in this project. Chapter 6 presents details about dataset used. And In the end Chapter 7 describes the Conclusion and future work.

Chapter 2: LITERATURE REVIEW

The early work of sentiment analysis began with subject detection, qualitative analysis back to the late 1990's. After this the area of research shifts its focus towards the interpretation of metaphors, main purpose of views, narrations, affects, evidentially in text and different connected areas. Shown below is that the literature describing the first works of sound judgment and detection of affects within the text. With the rise in web usage, the online became a supply of importance as text repositories. Consequently, a switch was slowly created off from the employment of sound judgment analysis and towards the employment of sentiment analysis of the online content. Sentiment analysis is currently become one amongst the dominant approach used for extracting sentiment from text and appraisals from on-line sources like websites and blogs. Separating non self-opinionated, neutral and objective sentences and text from the subjective sentences carrying several sentiments may be a troublesome job, however, it's already been explored seriously in closely connected however separate field (J M Wiebe, 1994). It basically concentrates on creation of a distinction between "subjective" and "objective" words and texts, At one side of research, the subjective ones gives results and opinions and while on the opposite hand, the target ones square measure being employed to gift info that is factual (Wiebe, Wilson, Bruce, Bell, & Martin, 2004) (Wiebe & Riloff, 2005). this can be quite totally different than sentiment analysis with reference to the set of classes into that the language units square measure classified by every of those 2 analyses. Subjective analysis focuses on dividing the language units into 2 categories; objective and subjective, wherever as sentiment analysis tries to divide these language units into 3 categories: -ve, +ve and neutral. the realm of concentration in a number of the first studies was with sound judgment detection solely (M. Wiebe, 2000). With the passage of your time to time and a necessity for higher understanding of system and extraction, momentum slowly raised towards sentiments classification and linguistics orientations.

Like other development fields of research today, sentiments analysis is a terminology yet to be matured; moreover just attempting to define sentiment can be difficult to accomplish [14]. The words sentiment [13][15], polarity [11] [12] [17], opinion [19], [20], semantic orientation [12] [21], attitude [22] and valence [22] are used to represent similar if not the same idea. These words, more often than not, used either to make the reference to various aspects of one particular phenomenon, an example being [14] [24] where the sentiment is defined as an affective part of the opinion, or simply can be used as synonyms for each other without any true definition of their self. Furthermore, some of these sentences can be confusing because of their multiple synonyms already in linguistic tradition (ex. polarity, valence) and therefore are confusing. For our present study, the focus is on capturing expressed sentiment in a text as -ve, +ve or neutral; therefore, we will refer this domain of research as a sentiment analysis. Our preference is of using the term 'sentiment analysis' is due to the fact that: 1) the possibility of confusing this study with research in other areas is not likely because the term is not belong with any other research tradition, 2) the kind of data which was extracted from the text is accurately reflected [unlike in the case of opinions which could also possess a topical component], and lastly 3) it is parsimonious and precise.

Recently, there is a change of attitude in the area of sentiment analysis whereby the concentration is now on classification, which has been added a third category known as neutrals [16].it is no longer focused on the binary classification of only +ve/-ve [21]. Through observations, there came a realization that it is much easy to separate +ve elements from -ve ones than it is to differentiate +ves or -ves from neutrals. Majority of disagreements among the human annotators as well as the errors resulting from utilize the automatic systems are associated with attempting to separate the neutral words, sentences or text from those that are either -ve or +ve [16] . Moreover, a problem arises from the meaning attributed to the term 'neutrals'. This is because 'lack of opinion' [18] as well as 'a sentiment that lies between +ve

and -ve [13] are both meanings of 'Neutrality' used in related literature. The latter definition is favored by sentiment analysis while the field of subjectivity analysis mostly use the previous interpretation. However, it is the latter meaning of the word that will utilized in this dissertation.

A rating inference as a metric labeling problem was developed by [29]. They achieved this by first apply two n array classifiers, which can included one-vs.-all SVM and SVM regression, in order to classify these reviews in regards to multi point rating scales. After applying these classifiers, a metric labeling algorithm was utilized so that the results of the n array classifiers were completely changed in order to guarantee that the like items receive like labels. A similar function was determined from this. It is true that a typically used similar function in topic classification is the overlapping of terms however, when attempting to identify reviews having like ratings, it is not particularly effective [20]. The +ve Sentence Percentage (PSP) similarity function was subsequently introduced; which calculates the number of sentences which are considered +ve divided by the number of sentences in the review so that that are considered to be subjective. Results of experiments generally have shown an improvement in n-ary classifier performance when making use of metric labeling with PSP. Pang and Lee's work was later augmented by [30] where they used transductive semi supervised learning in their study. It is shown that classification accuracy could be improved upon with the help of reviews without user specified ratings, in other words, unlabelled reviews [18].

A kernel based regression algorithm which was introduced by [31] 2007, made use of order preferences of unlabelled data and it was successfully applied to the sentiment classification. The order preference of a pair of unlabeled data x_i and x_j indicates that x_i is preferred to x_j to some degree, even though exact preferences for x_i and x_j are unknown. For ex, in framework of sentiment analysis, when presented with the two reviews of unknown rating values, it is quite possible to determine which review is more +ve. They executed their algorithm with the

rating inference problem. As a result, it was evidenced that by utilizing order preferences the performance of rating inference was much better than standard regression.

Corpus based machine learning method or methods on compilations are able to compile lists of -ve and +ve words with a high accuracy of up to 95%. In order to reach their full potential, most of these approaches need immense annotated training datasets. Corpus based methods can overcome some of these limitations by utilizing dictionary based approaches since these approaches depend on existing lexicographical resources (such as Word-Net) to provide semantic data in regards to individual senses and words [24]. Suggested that when analyzing sentiment, semantically similar does not necessarily imply sentimental similarity. This suggestion was made on the base of statistical observations from a compilation of movies' reviews. Subsequently, a method for determining the semantic orientation of the opinion is proposed on the basis of relative frequency. An estimation of the opinion strength of a word and the semantic- orientation in regards to a sentiment class and its relative freq of appearance in that class is carried out using this methodology. For ex, if the word 'best' appeared 8 times in +ve reviews and 2 times in -ve reviews, its strength with respect to +ve semantic orientation is then $8/(8+2) = 0.8$ [24].

Introduction of the new features, that are conceptually related to the key phrase frequency were done . On the basis of candidate phrases in the input document, these new features can be generated, by issuing the queries to Web search engines. An improvement in key phrase extraction has been experienced with these new although they are neither domain specific nor training intensive. The feature values are calculated from the number of hits for queries (the number of matching Webpages). A large collection of the unlabeled data, approximately 350 million Webpages without manually assigned key phrases, has been mined for the lexical knowledge to derive these new features. Simple methods for combining individual sentiments [16] and supervised [17] statistical techniques was proposed which can measure sentiment of

the phrase or sentence level using opinion oriented words. Another popular method, proposed by [19], makes use of both lexical and syntactic features for the sentiment analysis and is a machine learning approach. This method, missed pertinent contextual information on which it indicates that the individual sentence itself is a vital when extracting semantic orientation.

An alternative method was suggested by [19] for utilizing WordNet's synonym relations for tagging words with Osgoods three semantic dimension. The shortest path of joining a particular word to the words 'good' and 'bad' was calculated through the WordNet relations in order to assign the values of +ve or -ve to the word. Dictionary based methods for sentiment classification at the word level have no need of large corpora, or search engines having special functionalities. Rather they depend on the readily available lexical resources existing today such as Word-Net. They are able to compile the comprehensive, accurate and domain independent word lists containing their sentiment and the subjectivity annotated senses. Such as a lists provide a vital resource for sentence or text sentiment classification and because of the early compilation they are able to increase the efficiency of sentiments classification at texts and sentence level. In contrast to other works this work presents sentence level lexical/dictionary knowledge base methods to tackle domain adaptability problem for different types data [9].

Dictionary based techniques that make use of the data found in references and lexicographical resources, such as Word Net and the thesaurus in which can be used for assigning sentiments to a large number of words. Majority of such methods utilize the various relationships between the words (synonymy, antonymy, hyponymy / hyperonymy) in order to find seed words and other entries as described earlier. The data exists in dictionary definitions is made use in the wordlevel sentiment orientation in some of the recent methods. For semantic orientation lexical based semantic terms are extracted using dictionaries like Senti WordNet, Concept Net etc. for sentence level classification. According to [19]. The first try at employing Word Net relations

in a word sentiment annotation was made by [16][17]. They made the suggestion about an extension to the lists of manually tagged +ve and -ve words by adding to the list the synonyms for those words. They began with just 56 verbs and 36 adjectives. The method was applied in 2 occurrences and acquired 6070 verbs and 12213 adjectives. Then on the basis of the strength of the sentiment polarity which had been assigned to each word, the words which has been acquired were ranked. This strength of-sentiment score or rank for each word was calculated by maximizing probability of the category of the word's sentiment in regards to its synonyms.

Semantic characteristics, like word sentiment, of each word are greatly acknowledged as the good indicators of semantic characteristics of a phrase or a text that contains them, e.g. in (B. Baharudin, 2010) [21]. A sentence or text level sentiment annotation system that can uses words as indicators (features) of sentiment and therefore, It requires the creation of words lists annotated with sentiment markers. The research on word level sentiment annotation has been produced a number of such lists of words that were manually or automatically tagged as sentiment or classified as related to sentiment. [40]

[20] suggested a method that would use different information occurred at the same time in order to acquire words related to opinion (ex., disapproval, accuse, commitment, belief) from the text as a way to carry out analysis of subjectivity at the word-level. Two different techniques was used. The log likelihood ratio is computed with the first technique: using the data obtained by calculating how often the words obtained from one sentence occur with seed words taken from [50]. Relative frequencies of words found in documents, either subjective or objective, are computed by using second technique.

When NLP and statistical techniques are utilized, much importance is given to sentiment analysis at the word or feature level because it is an analysis of the text with the most detail. The semantic orientation of a given phrase or a review word is determined by the techniques

proposed by [18] and [41]. Several researchers used a preset seed word to enable extraction of opinion-oriented words and features [42] [43] and form a list used for semantic orientation, extraction and classification of opinion. Determining the polarity and subjectivity of a text is not the only aim of sentiment analysis. On the contrary, what the writer of the text specifically likes or dislikes regarding an object is also of importance [44]. Our main focus here is to discuss sentence and document level sentiment analysis. Sentence level analysis decides what the primary or comprehensive semantic orientation of a sentence is while the primary or comprehensive semantic orientation of the entire document is, handled by the document level analysis [13] [43]. Document level sentiment analysis deals with a document as a whole and classifies all the sentiments which have been expressed about a certain object by the author showing whether the overall document sentiment is +ve, -ve or neutral. However, the text document or review are split down into sentences for sentiment analysis to the sentence level. These sentences then evaluated by utilizing statistical or lexical methods in order to determine their semantic orientation. Three main steps are involved, namely Pre-processing, Text Analysis and Sentiment Classification. A compilation of specific reviews are taken as input by the model and are then processed according to the above three steps to obtain results. Review classification and evaluation of sentences or expressed opinions in the reviews are the results produced by the model. The machine learning method and topic classification are similar in the sense that topics are classes of sentiment such as Negative and Positive [14]. This is how it works: a review is broken down into phrases or words, the review is then presented as a document vector (bag-of-words model), and finally, the review is classified on the basis of the document vectors.

It is apparent that classifying a sentiment can easily be formulated as a supervised learning problem which has two class labels (negative and positive). In regards to the assumption above, it is not a surprise that the reviews utilized in existing research regarding data for training and

testing are mostly product based. Data for training and testing is easily available due to any typical review site having already assigned a reviewer rating (e.g. 1-5 stars) to each review [45]. Commonly, a thumbs-up or positive review will be assigned 4-5 stars while a negative or thumbs-down review is assigned with only 1-2 stars. Studies present to date have taken unlabeled data from the domain of interest with labelled data from another domain as well as general opinion words and made use of them as features for adaptation[46] [12] [18].

In this thesis a technique for domain independent sentence level classification of sentiment is introduced [48]. Rules for all the parts of speech are applied so that this can be scored on the strength of the semantics, contextual valence shifter, and sentences structure or expressions on the basis of dynamic pattern match. However, word sense disambiguation to fetch accurate sense of the sentence has already been addressed. Opinion type, confidence level, strength and reasons are all identified using this system. Senti WordNet and Word-Net are utilized as the basic knowledge base which has the further capability of being strengthened by using these modifiers, information in the contextual valence shifter and all parts of the speech.

Chapter 3: SENTIMENT CLASSIFICATION

Sentiment characterization is an opinion mining movement concerned with figuring out what, if any, is the general feeling introduction of the sentiments contained inside of a given report. It is accepted all in all that the report being examined contains subjective data, for example, in such as in product reviews and feedback forms. Opinion introduction can be named having a place with contradicting positive or negative polarities – positive or negative criticism around an item, great or unfavorable sentiments on a point – or positioned by range of conceivable conclusions, for instance in movie from surveys with input running from one to five stars.

Supervised learning systems using different aspects of content as sources of features have been proposed in the writing. Early work seen in [13] presents a few supervised learning algos using bag of words features common in content mining research, with best execution obtained using support vector machines as a form of combination with unigrams. Grouping terms from a report into its linguistic parts, or role of speech has also been investigated: In [21] form of speech information is used as component of a feature set for performing assessments of sentiments on a dataset of news wire articles, with similar methodologies attempted in [7], [10] and [16], on distinct data sets. On [20] a strategy that identifies and scores patterns in form of speech is applied to derive features for sentiment classification, with a comparable thought that applied to review extraction for product features seen in [4]. Separation of subjective and target sentences for the purposes of enhancing reports level sentiment classification are found in [14], where significant changes were obtained over a pattern word vector classifier. Other different studies concentrate on the correlation of composing style to overall sentiment, taking into the account that the use of colloquialisms and punctuation that may pass sentiment. In [22] a lexicon of colloquial expressions and a general expression rule base is to recognize detect unique opinion terms for example unusual spellings (“greeeat”) and word combination

("supergood"). In [1] report statistics and features measuring aspects of composing style are joined with the word vectors to acquire considerable standards over a baseline classifier on a data set of movie reviews.

3.1 Opinion Lexicons

Opinion lexicons are assets that associate with words and sentiment orientation. Their utilization in review mining research originates from the theory that individual words can be considered as a unit of review information, and accordingly may give clues to reports sentiment and subjectivity. Manually created review lexicons were connected to sentiment classification as seen in [13], where a forecast of document polarity is given by count +ve and -ve terms. A similar methodology is presented in the work of Kennedy and Inkpen [10], this time utilizing an opinion lexicon based on the combination of other existing resources.

Manually created lexicons however have a tendency to be constrained to a little number of terms. By its tendency, building manual lists is a period consuming effort, and may be liable to annotator bias, To overcome from these issues lexical induction methodologies have been proposed in the writing with a view to extend the size of sentiment lexicons from a core set of seed terms, either by investigating term connections, or by calculating similarities in report corpora. Early work in this area that is seen in [9] expands a list of +ve and -ve adjectives by assessing conjunctive statements in a report corpus. Another basic approach is to get opinion terms from the WordNet database of terms and relations [12], regularly by looking the semantic connections of a term such as equivalent words and antonyms.

3.2 WordNet Glosses and SentiWordNet

As noted in [15], term connections in the WordNet form of database create a highly disconnected graph, and along these expansions of sentiment data from a core of seed words

by looking there semantic relationships such as there meanings and antonyms is bounded to be restricted to a subset of terms. To defeat this issue, data contained in term glosses – informative content going with every term – can be investigated to gather term, based on the presumption that a given term and the terms contained in its gloss are likely to demonstrate the same polarity. In [2] a strategy for lexicon expansion is proposed where terms are assigned +ve or -ve opinions based on the presence of terms known to carry opinion content found on the term gloss. The creators those argue that the glosses contains potentially low level of noise since they “are intended to match as close as possible as expected by the components of meaning of the word, have generally standard style, language and syntactic structure”; This thought is additionally seen in [5], this time by utilizing managed learning systems for extending a lexicon by investigating gloss data, yielding +ve accuracy enhancements over a gold standard in compare to some portion of techniques previously discussed in this article. This is the same methodology utilized on building the SentiWordNet opinion lexicon [6].

SentiWordNet is built in a two-stage approach: first is, WordNet term connection such as antonym, synonym and hyponymy are investigated to extend a core of seed words used in [19], and known earlier to carry +ve or -ve review bias. After a fixed number of cycles, a subset of WordNet terms is acquired with either a +ve or -ve label. These terms are then used to prepare a committee of machine learning techniques. To minimize bias, the classifiers are prepared using diverse alogs and different training sets size. The predictions from the classifier committee are then used to determine the sentiment orientation of the remainder of terms in WordNet. The table below compares the coverage of SentiWordNet in relation to other assume built opinion lexicons accessible in the literature.

Table 1 Coverage of Opinion Lexicons

Opinion Lexicon	Total Sentiment Bearing Terms
General Inquirer [17].	4216
Subjectivity Clues Lexicon [21].	7650 (out of 8221 terms)
Grefenstette et al [8].	2258
SentiWordNet [6].	28431 (out of total 86994 WordNet terms)

3.3 Web Crawling using HTML AGILITY PACK

To parse HTML from a website is otherwise called Screen Scraping. It's a process to access external website information (the information must be public – public data) and processing it as required. For instance, if we want to get the average ratings of Nokia Lumia 1020 from different websites we can scrap the ratings from all the websites and calculate the average in our code. So we can say, as a general “User” what you can have as “Public Data”, you'll be able to scrap that using HTML Agility Pack easily. Previously it was harder to scrap a website as the hold DOM elements used to be downloaded as string. So it wasn't a pleasure to work with strings and find out individual nodes by iterating through at and matching tags and attributes to specify your requirements. Gradually the way has improved and now it has become too easy using HtmlAgilityPack library. That's why this article will give you a simple demonstration on how to start with HAP.

You need to have the basics of programming and must know writing code in C# and ASP.NET.

Before HTML Agility Pack we had to use different built-in classes in .NET Framework to pull out HTML from a website. But now we don't have to use such loads of classes rather we'll use the HAP library and order it to do the task for us.

It's pretty simple. Your code will make an HTTP request to the server and parse/store the returned HTML.

First HAP creates a DOM view of the parsed HTML of a particular website. Then it's really some lines of code that will allow you to pass through the DOM, selecting nodes as you like. Using an XPath expression, HAP also can give you a specific node and its attributes. HAP also includes a class to download a remote website.

Chapter 4: THE PROPOSED FRAMEWORK

In our framework there are two main component.

- i) Preprocess user reviews and fetch rating from it about the movie.
- ii) A Collaborative Filtering Module which uses the user's rating and then generates Recommendations.

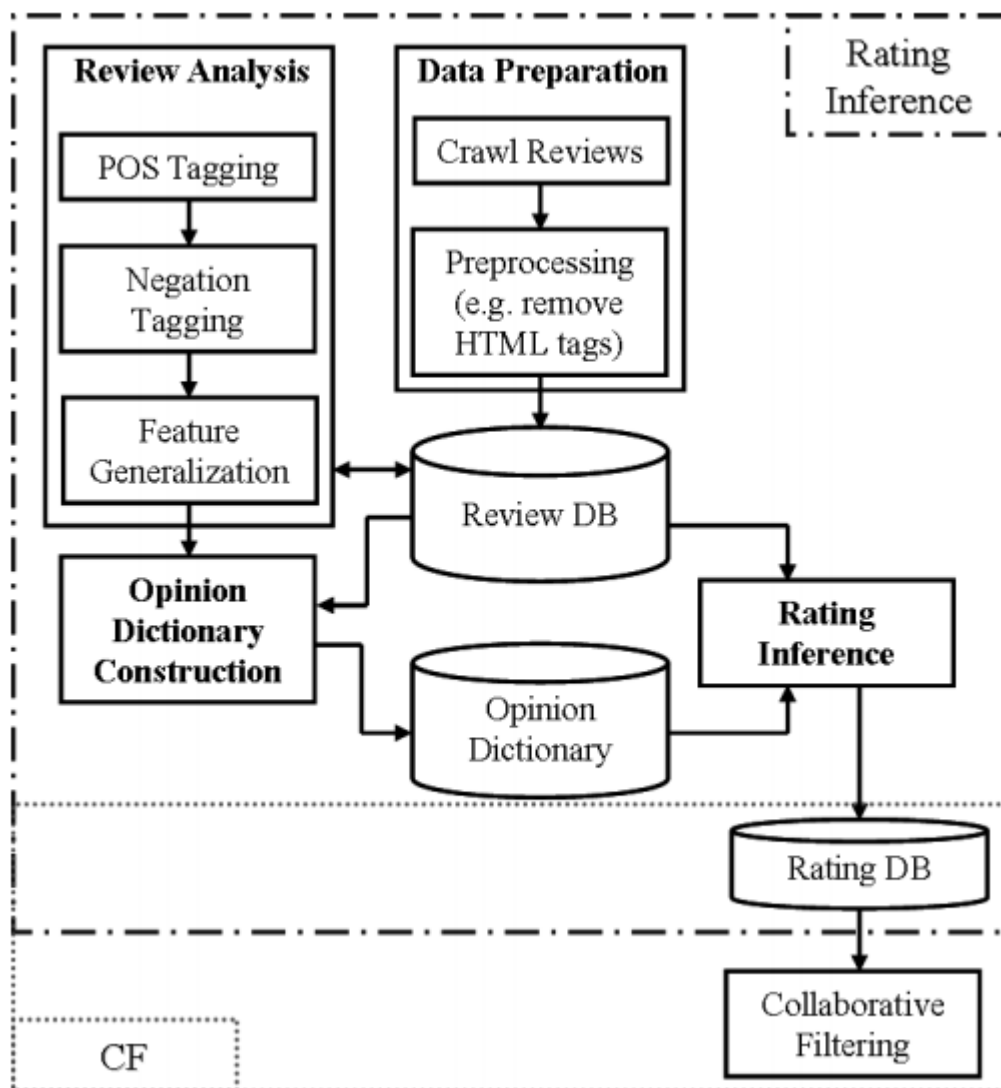


Figure 1: Overview of framework

- i) Here the review refers the text comments written by the user about the movie on the relevant subject matters. Our dataset involves the user comments. In Imdb.com different preprocessing steps are required. A user review is likely to be a semistructured document, containing some structured headers and an unstructured text body. A movie review on IMDb, for example, contains structured headers including a user (author) identity and a one-line summary, and unstructured blocks of text, which are the user's comments on the movie being reviewed, written in natural language. We are also storing the user Identity, subject matter.

4.1 Data Extraction

Data Extraction involves collecting and pre-processing user reviews for sentiment analysis. We've used different pre-processing steps looking on the information sources. In our case, we've downloaded user reviews as hypertext markup language pages, the html tags and non-textual contents they contain are removed during this step. A user review is probably going to be a semi-structured document, containing some structured headers and an unstructured text body. A movie review on IMDb.com, contains structured headers which includes a user identity, summary, and unstructured blocks of text, which are the user's comments on the movie being reviewed, written in a natural language. Sentiment analysis algorithms usually don't use data aside from the comments and the original ratings given by the users (e.g. for performance evaluation), if any. Our framework, however, extracts additionally the identities of users and also the topics being reviewed as a result of their helpful for activity CF, and such data square measure preserved to facilitate our future work. Since we tend to specialise in rating abstract thought during this work, the term "reviews" hereafter refers to the comments given by the users on the relevant topics.

4.2 Analyzing reviews

It is the first and most important part of our framework. As explained before the IMDb.com's Data are semi-structured and the review analysis step includes many task that help in

identifying the important information in reviews. The document level sentiment analysis try to analyse the full document (such as one review) into '+ve' ,'-ve' or neutral class. The methodologies based on SentiWordNet focuses the term profile of the document and concentrate terms having desired POS label (such as adjectives, adverbs or verbs). This obviously shows that before applying the SentiWordNet based formulation, the review text should be applied to a POS tagger which tags each term occurring in the review text. At that point some chose terms (with wanted POS tag) are removed and the opinion score of every extricated term is gotten from the SentiWordNet library. The scores for every removed term in a review are then accumulated utilizing some weightage and accumulation plan.

Subsequently two key issues are to choose

- (a) Which POS labels ought to be separated, and
- (b) How to choose the weightage of scores of distinct POS labels extricated while registering the total score.

We have investigated with diverse linguistic highlights and scoring plans. Computational Linguists propose that modifiers are great markers of reviews. Case in point, if a review sentence says "The movie is incredible", then utilization of modifier "incredible" lets us know that the movie was loved by the analyst and perhaps he had a nice experience by utilizing it. At times, Adverbs further adjust the sentiment communicated in audit sentences. Case in point, the sentence "The movie is extremely good" communicates a more +ve supposition about the movie than the sentence 'the movie is great'. A related past work [12] has additionally inferred that 'Adverb + Adjective' consolidate creates preferred results over utilizing modifiers alone. Subsequently we favored the 'adverb + adjective' consolidate over removing 'descriptive word' alone. The adverbs are usually used as complements or modifiers. Few more examples of this usage are:- he ran quickly; only adults; very dangerous trip; very nicely; rarely bad; rarely good

etc. In all these examples adverbs modify the adjectives. Though adverbs are of various kinds, but for sentiment classification only adjectives of degree is useful. Some past works have recommended misusing the "verb" POS labels in addition to 'adjective' for sentiment classification. Here, we have investigated with two semantic highlight determination plans. In one we only concentrate on 'adjectives' and any 'adverbs' going before the selected adjective. In the other one we separate both 'adjectives' and 'verbs', along with any 'adverbs' going before them. Since, adverbs are changing the scores of succeeding terms, it needs to be chosen as to what extent the sentiment score of an 'adverb' should change the succeeding 'adjective' or 'verb' sentiment score, to obtain higher accuracy. We have chosen the modifying weight age (scaling factor) of adverb score as 0.35, in view on the conclusions reported in [14] and [11]. The other fundamental issue that remains to be addressed is how should the sentiment scores of chosen 'adverb+adjective' and 'adverb + verb' consolidated be aggregated. For this we have attempted different factors of weight ranging from 10% to 100%, i.e. the 'adverb + verb' scores are combined to 'adverb + adjective' scores in a weighted way. In the first plan of utilizing only 'adverb + adjective' join, we have picked a scaling element of = 0.35. This is proportionate to giving just 35% weight to adverb scores. The changes in adjective scores are thus in a fixed proportion to adverb scores. Since we picked a value of scaling variable $sf = 0.35$, the adjective scores will get a higher priority in the consolidated score. The demonstrative pseudo-code of key steps for this plan i.e. Senti-WordNet (AAC) is illustrated below. Here AAC refers to Adverb + Adjective Combinations.

Table 2 SENTIWORDNET DATABASE STRUCTURE

Fields	Descriptions
POS	Part of speech linked with synset. This can take four possible values: a = adjective n = noun v = verb r = adverb
Offset	Numerical ID which associated with part of speech uniquely Identifies a synset in the database.
PosScore	Positive score for this synset. This is a numerical value ranging from 0 to 1.
NegScore	Negative score for this synset. This is a numerical value ranging from 0 to 1.
SynsetTerms	List of all terms included in this synset

For each sentence, extract adv+adj combines.

For each extracted adv+adj combine do:

- ❖ If adj score=0, ignore it.
- ❖ If adv is affirmative, then
 - If score(adj)>0
 - ◆ $fsAAC(adv,adj) = \min(1, score(adj) + sf * score(adv))$
 - If score(adj)<0
 - ◆ $fsAAC(adv,adj) = \min(1, score(adj) - sf * score(adv))$
- ❖ If adv is negative, then
 - If score(adj)>0
 - ◆ $fsAAC(adv,adj) = \max(-1, score(adj) + sf * score(adv))$
 - If score(adj)<0
 - ◆ $fsAAC(adv,adj) = \max(-1, score(adj) - sf * score(adv))$

Figure 2 Pseudo code: SWN(AAC)

Here, adj refers to adjective and adv refers to adverb. The last sentiment values [fsAAC] are scaled form of adverb and adjective Senti-WordNet scores, where the adverb score is given 35% weightage. The presence of 'Not' is taken care by subtracting the scores obtained. Firstly we picked the sentence boundaries of a review and then we process all those sentences. For every sentence we choose the adv + adj combines and then compute their sentiment scores according the scheme described in the 723 pseudo code. The final document sentiment score is an addition of these sentiment scores for every sentence occurring in it. The score value decided the polarity of the review.

The second usage that we attempted joins both 'adverb + adjective' and 'adverb + verb' sentiment scores. It is same like to the previous scheme in its method of joining the adverbs with adjectives or verbs, difference is in the logic that it counts both adjectives and verbs for choosing the overall sentiment score. We have tried it with different aggregation of weights for adjective and verb scores and conclude that 30% weight for verb score produces best precision levels. The occurrence of word 'not' has been handled in a same manner as in previous scheme. The indicative pseudo code of key step for this scheme, i.e. Senti-WordNet (AAAVC) is illustrated below. Here AAAVC refers to Adverb + Adjective and Adverb + Verb Combine.

In this scheme, we compute sentiment score for all 'adverb+adjective' and 'adverb+verb' combines in a sentence and aggregate them together. This is done for all sentences and the document-level sentiment polarity value is determined based on the aggregated sentiment score of the review document.

For each sentence, extract adv+adj and adv+verb combines.

1. For each extracted adv+adj combine do:

- ❖ If adj score=0, ignore it.
- ❖ If adv is affirmative, then
 - If score(adj)>0
 - ◆ $f(\text{adv}, \text{adj}) = \min(1, \text{score}(\text{adj}) + sf * \text{score}(\text{adv}))$
 - If score(adj)<0
 - ◆ $f(\text{adv}, \text{adj}) = \min(1, \text{score}(\text{adj}) - sf * \text{score}(\text{adv}))$
- ❖ If adv is negative, then
 - If score(adj)>0
 - ◆ $f(\text{adv}, \text{adj}) = \max(-1, \text{score}(\text{adj}) + sf * \text{score}(\text{adv}))$
 - If score(adj)<0
 - $f(\text{adv}, \text{adj}) = \max(-1, \text{score}(\text{adj}) - sf * \text{score}(\text{adv}))$

2. For each extracted adv+verb combine do:

- ❖ If verb score=0, ignore it.
- ❖ If adv is affirmative, then
 - If score(verb)>0
 - ◆ $f(\text{adv}, \text{verb}) = \min(1, \text{score}(\text{verb}) + sf * \text{score}(\text{adv}))$
 - If score(verb)<0
 - ◆ $f(\text{adv}, \text{verb}) = \min(1, \text{score}(\text{verb}) - sf * \text{score}(\text{adv}))$
- ❖ If adv is negative, then
 - If score(verb)>0
 - ◆ $f(\text{adv}, \text{verb}) = \max(-1, \text{score}(\text{verb}) + sf * \text{score}(\text{adv}))$
 - If score(verb)<0
 - ◆ $f(\text{adv}, \text{verb}) = \max(-1, \text{score}(\text{verb}) - sf * \text{score}(\text{adv}))$

3. $f_{AAVC}(\text{sentence}) = f(\text{adv}, \text{adj}) + 0.3 * f(\text{adv}, \text{verb})$

Figure 3 Pseudo code: SWN(AAFC)

4.2 Rating Inference:

A review usually contains a mixture of positive and negative opinions towards different aspect of movies, and rating inference aims at determining the overall sentiment implied by the user.

We attempted to perform such task by aggregating the strengths of the opinion words in a review with respect to different sentiment classes, and then assigning an overall rating to the review to reflect the dominant sentiment class.

The movie features on which opinions are expressed may also be useful for determining weights of opinions, and this is facilitated by the feature generalization task. Opinions towards

a movie as a whole may be more useful for determining the SO of a review. This also allows easy integration with user-specified interest profiles if necessary (e.g. to address the new user problem in CF [19]). For example, if a certain user of a movie recommender system is particularly interested in a certain actor, then the acting of that actor in a movie may have stronger influence on the his overall sentiment towards the movie.

Chapter 5: SOLUTION APPROACH

In this chapter we quickly describe all the procedures and objectives of this work and we aim to succeed as a result. This work sorted into 3 stages. First phase is the web page crawling phase, in which we collect the data from movies' review websites. The 2nd stage is the dissecting phase, in which we parse the data, prepared and dissected to find valuable information. The 3rd stage is the visualization phase, in which information is visualized to clearly understand the results.

5.1 Review extraction process

Web blog are full of un-index and unstructured text that reflects the opinions of people. Many people make choices by taking the suggestions of other people into account. Thus, there is a need to crawl and process peoples' opinions, so that it can be used in decision making processes of potential Web review applications. In this study, we propose a blog mining system that will extract movie comments from Web blogs and that will show Web blog users what other people think about a particular movie. Fig. shows the overall process model.

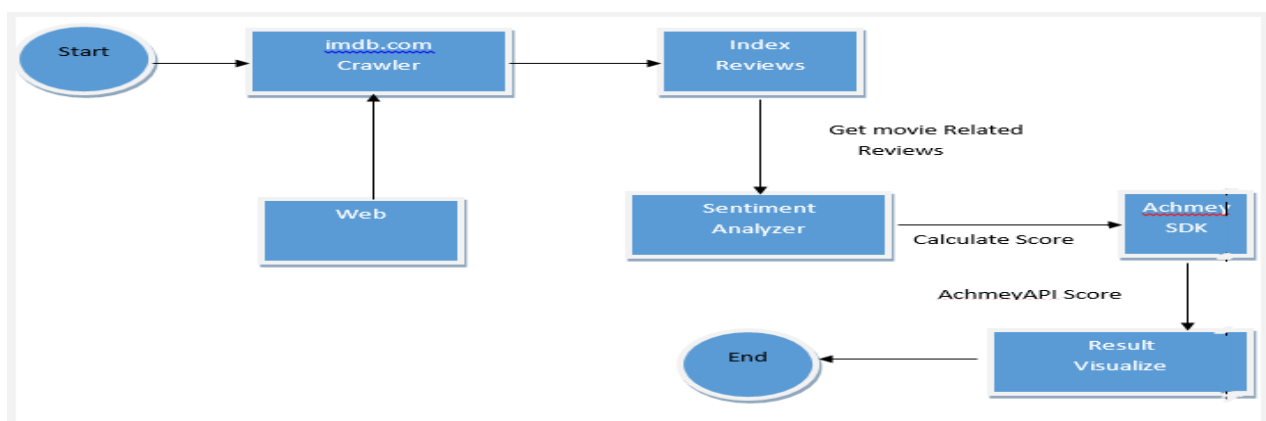


Figure 4 Overall process Model

This system architecture provides consist of several components like: Web crawler, sentiment analysis and web user interface.

5.2 Imdb.com Crawler

Web crawlers are the computer program that traverses the Websites in a systematic way with the purpose of collecting of data. A web crawler is use to download the Web pages for indexing and other purposes like structural analysis, page validation, visualization, update notification, for the spam purpose like collecting email addresses etc. the main objective of search engine is to provide more relevant results in faster time over rapidly expanding websites. There are 3 important sequential tasks a standard search engine does as shown[10]:

- a) Crawler
- b) Indexing
- c) Searching

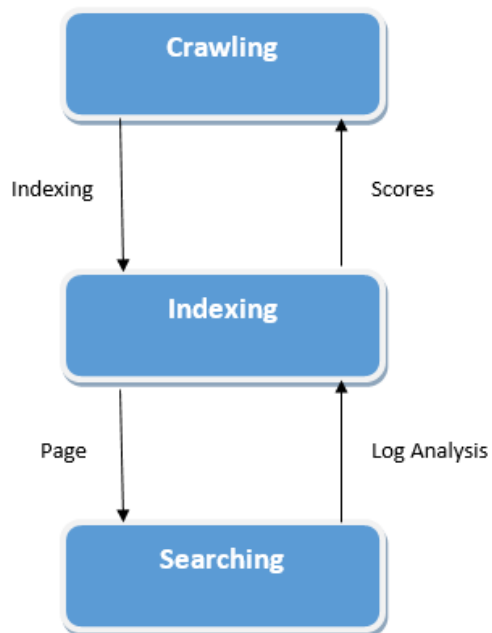


Figure 5 General task of search engine

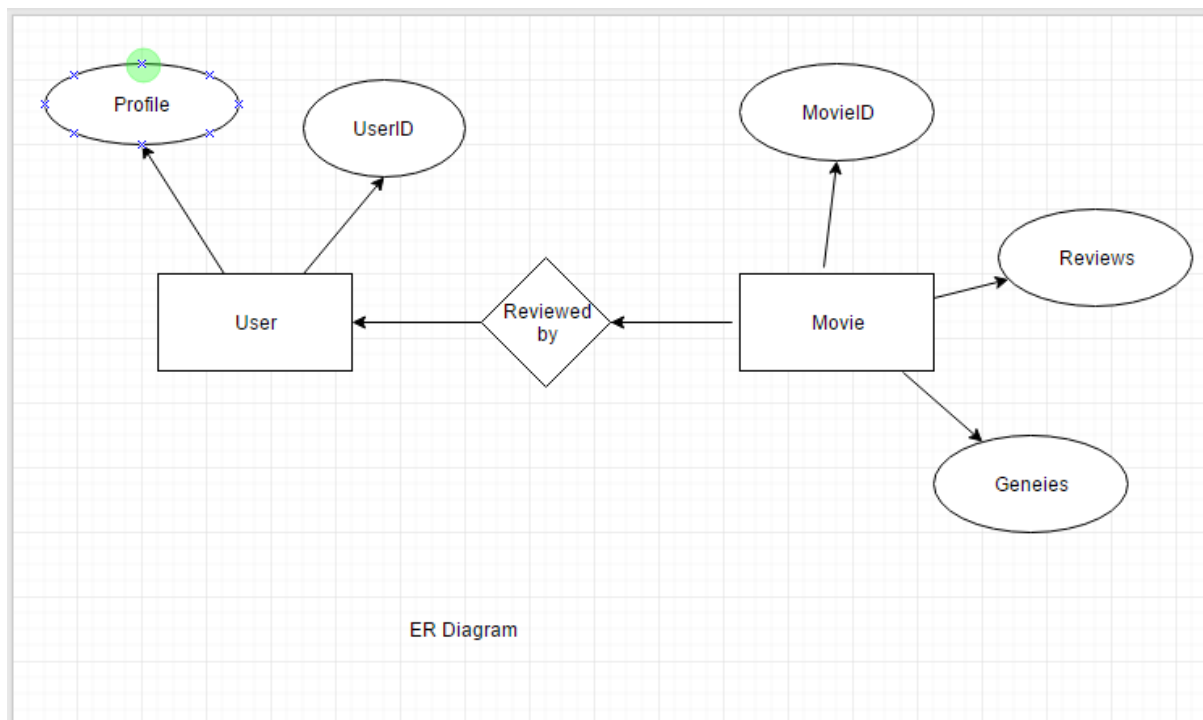


Figure 6 : Database structure (Tables ER)

5.3 Sentiment Analysis

'Sentiment Analysis is the task of identifying +ve and -ve opinions, emotions, and evaluations'. Sentiment Analysis has many different names. It's often referred to as Opinion mining, subjectivity analysis, and appraisal extraction with some connections in an affective computing. It is a technique for fetching opinions from unstructured human authored documents. In a simple word it is used to track the mood of public. It is an evolving field having roots in Natural Language Processing (NLP), Computational Linguistics and Textual Mining. There is a wide range of tools in the market that performs the automatic sentiment analysis of a given text. Many sentiment search engines exist in which users run typical queries on any of the topic of interest, and generate the text results. Usually results are coded and categorized into 2 or three polar categories. Some examples currently available are: Topsey, subjectivity analysis, BackTweets, Tweet Beep, Reachli, Twitterfall, Social Mention, Trackur, Sentiment.ly, Sentiment140, Twendz, Opinion Crawl, Amplified Analytics, Lithium, Open Amplify, SAS Sentiment Analysis Manager, IBM Social Sentiment Index, Twittratr, SAS

Sentiment Analysis Studio, Tweet Sentiments etc.

5.4 Web User Interface

The Web user interface is formed mainly under 2 categories. The 1st category is the selection. There are two types of options in selection. 1st is the selection of movies. In Which, the system lets the user to select a movie and then shows its sentiment score results corresponding to nine different keywords categories. 2nd is the selection of keywords categories. Here, the system lets the user specify only 1 category and shows the sentiment scores on different movies under the selected keyword category.

Chapter 6: DATASET AND PERFORMANCE MEASURES

The results are shown on two groups of experiments. We have performed our analysis on IMDB movie reviews to assist the task of finding the opinion strengths.

6.1 Dataset

We have collected movie reviews from IMDb.com for the movies in the MovieLens 100k dataset, courtesy of GroupLens Research [10]. The MovieLens dataset contains user ratings on 1692 movies. We removed movies that are duplicated or unidentifiable (movies without names), and crawled the IMDb to download user reviews for the remaining movies. We filtered out contributions from users who have provided fewer than 10 reviews and reviews without user-specified ratings, which will later be used for evaluating our proposed framework. The resulting dataset contains approximately 30k reviews on 1477 movies, provided by 1065 users. Each review contains a number of headers and a text body. The headers include movie ID, user ID, review date, summary, which is a one-line summary in natural language text written by the user, and a rating, which is a user-specified number ranging from 1 (awful) to 10 (excellent). The text body is the user's comments on the movie.

6.2 Analysis on the Use of Opinion Words

Determining opinion strengths would be a simple task if an explicit opinion word forever seems in reviews with an explicit rating, to illustrate, if the word “brilliant” forever seems in reviews rated as 10/10. This is, however, not going to be true. A review could contain each positive and negative opinions. This suggests a motion picture receiving a high rating may additionally have some unhealthy options, and the other way around.

We performed some preliminary experiments to investigate the employment of opinion words in user reviews. By doing thus, we have a tendency to hope to find fascinating usage patterns

of opinion words which will facilitate crucial opinion strengths. We have a tendency to initial performed the tasks represented in Chapter’s sub-section 4.1 and 4.2 on the dataset. we have a tendency to then indiscriminately sampled 3 coaching sets, namely T10, T5 and T1, every containing five hundred reviews whose user-specified ratings were 10/10, 5/10 and 1/10 severally. These ratings were chosen as they appear to be acceptable representative cases for Positive, Neutral and Negative sentiments. we have a tendency to use a program to extract opinion words, that are words labelled as adjectives [7], and reason their frequency counts in every of the coaching sets. Some frequent opinion words were any analyzed. The number of distinct opinion words appeared within the coaching sets is 4545, among that 839 (around 18.5%) appeared in 2 of the 3 coaching sets, and 738 (around 16.2%) appeared all told 3. we tend to more examined opinion words that appeared in additional than one coaching set. Table one lists, thanks to area constraint, the ten most frequent opinion words (top 10) of this type in every coaching set in dropping order of their frequency counts. Within the table, the amount in brackets following associate opinion word is its frequency within the explicit coaching set. Bold-face is employed to spotlight words having the very best frequency among the 3 coaching sets.

Table 3: Top 10 opinion words with relative frequencies.

Training set	Opinion words with relative frequencies
T1	bad (0.65), good (0.28), worst (0.89), much (0.49), more (0.46), other (0.28), first (0.28), better (0.29), many (0.24), great (0.14)
T5	good (0.39), more (0.54), much (0.51), bad (0.35), better (0.41), other (0.32), few (0.73), great (0.21), first (0.34), best (0.19)
T10	best (0.68), great (0.66), good (0.33), many (0.47), first (0.38), classic (0.71), better (0.30), favorite (0.75), perfect (0.75), greatest (0.85)

Our observations are summarized as follows. Firstly, the relative frequencies of positive opinion words are typically, however not forever, the very best in T10 and therefore the lowest in T1, and the other way around for negative opinion words. Table two lists as examples the relative frequencies of the foremost frequent opinion word (top 1) in every coaching set. Boldface is employed to spotlight the very best frequency of every opinion word. Such observation suggests that relative frequencies of opinion words could facilitate determinative their thus and strengths. As an example, the word “best” appeared in T10 sixty eight of the time. it should thus be thought of a positive opinion word with the strength zero.68. Secondly, nearly thirty fifth of all opinion words, as well as those having clear and robust understood thus (e.g. “best”), appeared in additional than one coaching set. we have a tendency to model this reality by adopting the fuzzy set thought [24], which suggests that associate attribute is a member of some fuzzy sets to sure membership degrees within the vary [0,1], determined by some membership functions. Within the context of our work, the “membership degree” of a word with reference to a sentiment category is decided by the frequency of the word within the corresponding coaching set. Let’s say, the word “best” has thus Positive, Neutral and Negative with strengths 0.68, 0.19 and 0.13 severally. the utilization of fuzzy sets to model user ratings in CF has recently been planned in [11], however our work deals with a unique drawback as we have a tendency to adopt the fuzzy set thought to model thus and opinion strengths.

Thirdly, the SO of opinion words determined by the relativefrequency-based technique might not believe their usually understood thus. Associate in Nursing example is that the word “frightening” that appears to be a negative sentiment. Its ratio in T1, however, is only 0.29. supported this observation, we tend to additional conclude that synonyms don't essentially have constant thus. maybe, “terrible” may be a word of “frightening” in WordNet, however its ratio in T1 is zero.75. Recall that word-similarity-based strategies create use of a group of seed adjectives and also the similarities between word meanings to work out thus of opinion words

[7, 9]. Our analysis, however, indicates that similar meanings might not imply similar sentiments. this means that our relative-frequency-based technique overcomes a serious limitation of the word-similarity-based strategies, as a result of it permits similar words to possess totally different thus.

Chapter 7: CONCLUSION AND FUTURE WORK

We propose a hybrid approach of recommendation using the reviews which are written in Natural language and then providing the input to CF. Using Alcamey API the unstructured, natural language is transformed into a numerical value which can be easily feed to existing CF algorithms.

This work conjointly outlines preliminary results of associate analysis of the projected framework. Any development of the framework remains current. A lot of elaborated descriptions regarding the framework and comprehensive results are rumored in a very follow-up article. As noted, our rating illation approach transforms matter reviews into ratings to alter simple integration of sentiment analysis and CF. we have a tendency to even so acknowledge the chance to perform text-based CF directly from a set of user reviews. An attainable answer is to model text-based CF as associate info retrieval (IR) drawback, having reviews written by a target user because the “query” and people written by different similar users because the “relevant documents”, from that recommendations for the target user may be generated. This remains as a noteworthy analysis direction for future work.

References

- [1] E. Brill, "A simple rule-based part-of-speech tagger", in *Proceedings of the 3rd Conference on Applied Natural Language Processing*, pp. 152–155, 1992.
- [2] S. Das and M Chen, 'Yahoo! for Amazon: Extracting market sentiment from stock message boards.', in *Proceedings of the Asia Pacific Finance Association Annual Conference*, 2001.
- [3] K. Dave, S. Lawrence, and D. M. Pennock, 'Mining the peanut gallery: Opinion extraction and semantic classification of product reviews', in *Proceedings of the 12th International World Wide Web Conference*, pp. 519–528, 2003.
- [4] A. Esuli and F. Sebastiani, 'Determining the semantic orientation of terms through gloss classification.', in *Proceedings of the ACM International Conference on Information and Knowledge Management*, pp. 617–624, 2005.
- [5] D. Goldberg, D. Nichols, B. Oki, and D. Terry, 'Using collaborative filtering to weave an information tapestry', *Communications of the ACM*, 35(12), 61–70, 1992.
- [6] "GroupLens", *GroupLens*, 2016. [Online]. Available: <http://www.grouplens.org/>, 2016.
- [7] M. Hu and B. Liu, 'Mining and summarizing customer reviews', in *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 168–177, 2004.
- [8] M. Hu and B. Liu, 'Mining opinion features in customer reviews', in *Proceedings of the 19th National Conference on Artificial Intelligence*, pp. 755–760, 2004.
- [9] J. Kamps, M. Marx, R.J. Mokken, and M. de Rijke, 'Using Wordnet to measure semantic orientations of adjectives', in *Proceedings of the 4th International Conference on Language Resources and Evaluation*, pp. 1115–1118, 2004.
- [10] J. A. Konstan, B. N. Miller, D. Maltz, J. L. Herlocker, L R. Gordon, and J. Riedl, 'Grouplens: Applying collaborative filtering to usenet news', *Communications of the ACM*, 40(3), 77–87, 1997.
- [11] C. W. K. Leung, S. C. F. Chan, and F. L. Chung, 'A collaborative filtering framework based on fuzzy association rules and multiple-level similarity', *Knowledge and Information Systems*, (forthcoming).
- [12] G. Linden, B. Smith, and J. York, 'Amazon.com recommendations: Item-to-item collaborative filtering', *IEEE Internet Computing*, 7(1), 76–80, 2003.
- [13] H. Liu. MontyLingua: An end-to-end natural language processor with common sense, 2004. <http://web.media.mit.edu/hugo/montylingua>.
- [14] G. Miller, R. Beckwith, C. Fellbaum, D. Gross, and K. Miller, 'Introduction to Wordnet: An online lexical database', *International Journal of Lexicography (Special Issue)*, 3(4), 235–312, 1990.
- [15] R. Mukras, A comparison of machine learning techniques applied to sentiment classification, Master's thesis, University of Sussex, Brighton, UK., 2004.

- [16] T. Nasukawa and J. Yi, 'Sentiment analysis: Capturing favorability using natural language processing', in Proceedings of the *2nd International Conference on Knowledge Capture*, pp. 70–77, 2003.
- [17] B. Pang and L. Lee, 'Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales', in Proceedings of the *43rd Annual Meeting of the Association for Computational Linguistics*, pp. 115–124, 2005.
- [18] B. Pang, L. Lee, and S. Vaithyanathan, 'Thumbs up: Sentiment classification using machine learning techniques', in Proceedings of the *Conference on Empirical Methods in Natural Language Processing*, pp. 79–86, 2002.
- [19] A.M. Rashid, I. Albert, D. Cosley, S.K. Lam, S. McNee, J.A. Konstan, and J. Riedl, 'Getting to know you: Learning new user preferences in recommender systems', in Proceedings of the *2002 International Conference on Intelligent User Interfaces*, pp. 127–134, 2002.
- [20] F. Ricci, 'Travel recommender systems', *IEEE Intelligent Systems*, 17(6), 55–57, 2002.
- [21] L. Terveen, W. Hill, B. Amento, D. McDonald, and J. Creter, 'Phoaks: A system for sharing recommendations', *Communications of the ACM*, 40(3), 59–62, 1997.
- [22] Traveljournals.net. <http://www.traveljournals.net>.
- [23] P.D. Turney, 'Thumbs up or thumbs down? Semantic orientation applied to unsupervised classification of reviews', in Proceedings of the *40th Annual Meeting of the Association for Computational Linguistics*, pp. 417–424, 2002.
- [24] L. A. Zadeh, 'Knowledge representation in fuzzy logic', *IEEE Transactions on Knowledge and Data Engineering*, 1(1), 89–100, 1989.
- [25] Abbasi, A., Chen, H., and Salem, A. Sentiment analysis in multiple languages: Feature selection for opinion classification in Web forums. *ACM Transactions on Information Systems*, 26, 3 1-34. ,2008
- [26] Andreevskaya A., Bergler S. (2006). Mining WordNet for Fuzzy Sentiment: Sentiment Tag Extraction from WordNet Glosses. In Proceedings of the *11th Conference of the European Chapter of the Association for Computational Linguistics – EACL 2006*.
- [27] Chapman W, Bridewell W, Hanbury P, Cooper G, Buchanan B. (2001). Evaluation of Negation Phrases in Narrative Clinical Report. Proceedings of *AMIA Symposium*, 105-109,2001
- [28] Dave K, Lawrence S, Pennock D. (2003). Mining the Peanut Gallery: Opinion Extraction and Semantic Classification in Product Reviews. Proceedings of the *12th International conference on the World Wide Web - ACM WWW2003*, Budapest, Hungary,2003
- [29] Esuli, A. and Sebastiani, F. (2005). Determining the semantic orientation of terms through gloss classification. Proceedings of the *14th ACM international Conference on information and Knowledge Management* (Bremen, Germany, October 31 - November 05, 2005). CIKM '05. ACM, New York, NY, 617-624.

[30] [Online]. Available: <http://www.sentiwordnet.isti.cnr.it>, Accessed: 27- Jul- 2016.

[31] [Online]. Available: <http://www.imdb.com>. Accessed: 27- Jul- 2016.