

Methodologies In Sentiment Analysis

Palli Suryachandra, Dr. P. Venkata Subba Reddy

Abstract: Sentiment analysis uses data mining processes and techniques to extract and capture data for analysis in collection of documents, like blog posts, reviews, news articles and social media feeds like tweets and status updates. It has been gained order to distinguish the subjective opinion of a document or quite popularity in the recent years. Several techniques have been utilized frequently including machine learning approaches and vocabulary oriented semantic algorithms. This article presents an intellectual study of various techniques which are used in the sentiment analysis process.

Keywords: Sentiment Analysis, Machine Learning, Language, Techniques

1. INTRODUCTION

Sentiment analysis uses information retrieval and computational linguistics. Sentiment analysis has advantages in various forms such as in marketing or for business purposes. In marketing, it is used to notice about the favorable or negative points about their new product which helps to determine how successful the new product is. A specific view or notion can be depicted as ideas prompted, opinions, judgements or coloured by emotions or emotions (Boiy et al. 2007) [1]. In Computational Linguistics, the core is on feelings instead of sentiments, opinions or perceptions. The terms 'opinion's and 'sentiment's are frequently availed substitutable. In general, the information of a text is divided into two categories. 1. Based on text 2. Based on persuasion. Whereas actualities or facts are observational utterances about events, entities and their opinions, characteristics are particular utterances that depict opinions of people, events and their properties, feelings towards entities or appraisals (Liu 2010) [2]. A persuasion can be depicted by the following four terms: Sentiment, Claim, Holder and Topic . The Holder affirms a fact about a Topic, and frequently relates a persuasion, such as 'bad' or 'good', with the affirm. It depicts a persuasion as an implicit or explicit aspect in text of the holder's negative, positive or neutral notice into the requirement about the topic. Sentiment analysis suits with computational operation of persuasion, sentiment, opinion and individuality in text (Pang & Lee 2008) [3]. The document inception is likely in the pattern of unstructured data. Opinion mining or Sentiment analysis hit an important role in our regular selection construction progress. The above mentioned selections may has scope from acquisition a commodity such as mobile phone to scrutinize the movie to constructive contribution; all the choices will have a great favor on the regular life. Now-a-days folk before purchasing a service/commodity will accomplish a glimpse to scrutinize websites, online shopping websites and social media websites to get an assessment related about the required commodity or utility in merchandise. Opinion Analysis or Sentiment Analysis handles many technological objections such as feature extraction, opinion orientation classification and object identification.

Generally Opinion analysis is implemented by using unsupervised learning and supervised learning, those are generally Support Vector Machine, naive Bayes and Neural Networks. Among the given existing techniques which are in above, Support Vector Machine is recognized to be further applicable for Opinion Analysis. Opinion or Sentiment classification categorized into three major parts, namely:

- Feature level
- Sentence level
- Document level

In Sentence and Document level, the Opinion analysis extracts only a single opinion from the single opinion holder and make use of only a single object. But these considerations are not applicable for several context or conditions. Getting an opinion for whole document/blog is not an efficient as getting opinion by assuming appearance of each substance in the specific sentence. The article arranged in the following manner. Section 2 discussed about various approaches, Section 3 discussed about various techniques for sentiment analysis and section 4 discussed about conclusion.

2. Approaches

Based on the context and perception, there exist four major approaches [4] in the sentiment analysis, which are given below.

- Knowledge-based approach
- Relationship-driven approach
- Language model driven approach
- Discourse driven approach

2.1 Knowledge-based approach

The objective of this method is, developing lexicon words that represent negative class or positive class. Before the sentiment analysis task, the opinion values of the words in the lexicon are decided. There are many ways existing for creation of Lexicons. One way is start with core words and add some linguistic features to get more morphemes. Another way is start with core word and adds other words to it based on the occurrences in a document. For supporting opinion mining applications and sentiment classification , There exist a resource, SENTIWORDNET 3.0, which is a publicly available lexical resource explicitly devised [5].

2.2 Relationship-based approach

The objective of this method is, distinct relations among characteristics and components are evaluated for opinion classification work. Those relations may be either relationships between product features or relationships between different participants. For example one may compute it as a function of

- Palli Suryachandra, Research Scholar, CSE Dept., SVUCE, SVU, Tirupati. Email: surya.palle@gmail.com
- Prof. Dr. P. Venkata Subba Reddy, CSE Department, SVUCE, SVU, Tirupati. Email: vsrpoli@hotmail.com

the sentiments on different features or components of it, when one wants to know the sentiment of customers about a product brand.

2.3 Language models

The objective of this method is, n-gram language models are built. Presence or frequency of n-grams can be used. In text classification, frequency of n-grams gives better results. Normally, the frequency is converted to TFIDF to take term's importance for a document into account. However, Pang et al. [6] initiate that term existence provides enhanced results than the word term occurrences. They analyzed on of movie reviews and they determined, uni-gram existence is more applicable in opinion analysis. Dave et al. [7] initiate that tri-grams and bi-grams treated superior to uni-grams in opinion classification of product scrutiny.

2.4 Discourse structures and semantics

The objective of this method is, for classification, Discourse structures and semantics avail discourse relationship among text components. In several scrutinizes(reviews), the comprehensive opinion is generally considered at the end of the document [6]. In this discourse-driven approach the opinion of the full analysis is getting by deciding the sentiment among various discourse units and the discourse relationships that presented among themselves. Usually, In this context, the end paragraph of the scrutiny (review) will be given high importance in deciding the opinion of the aggregate review.

3. TECHNIQUES

Sentiment Classification techniques(shown in Fig. 1) are separated into two different techniques which is ML and Lexicon based Approaches. [8,9,10].

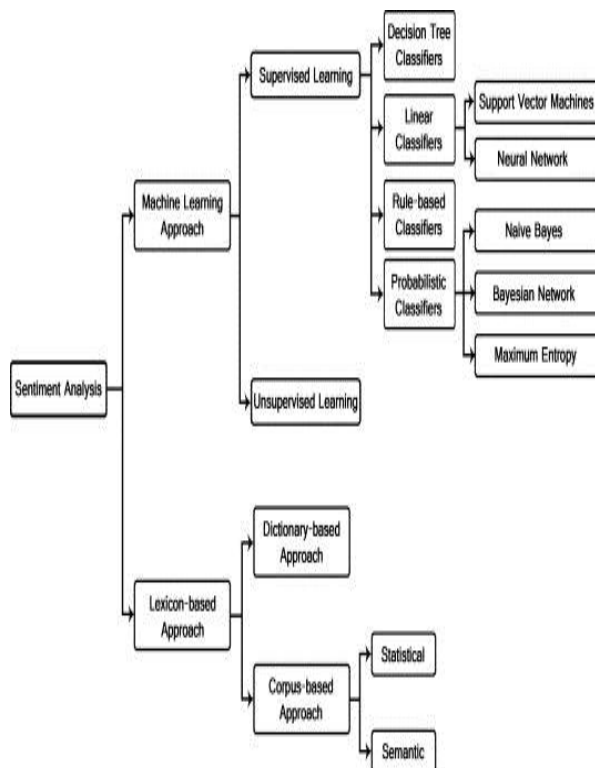


Fig 1 Sentiment Classification techniques

Machine Learning

Machine learning focuses on the development of computer programs that can access data and use it learn for themselves. The process of learning begins with observations or data, such as examples, direct experience, or instruction, in order to look for patterns in data and make better decisions in the future based on the examples that we provide. Machine learning methods builds based on the use of syntactic and/or linguistic features. Ssentiment analysis is usually a formal text classification issue that produce, use of linguistic and/or syntactic features. Text Classification Problem Definition: We have a set of training records $D = \{X_1, X_2, X_3, \dots, X_N\}$ where each record is labeled to a class. The hard classification problem is when one label is assigned to an instance.

3. A Supervised Learning:

Supervised techniques adapt the model to reproduce outputs known from a training set. In the beginning, the system receives input data as well as output data. Its task is to create appropriate rules that map the input to the output. The training process should continue until the level of performance is high enough. After training, the system should be able to assign an output objects which it has not seen during the training phase. In most cases, this process is really fast and accurate. These methods build upon the presence of labeled training texts. There exist more supervised classifiers like Regression and Classification. we showed some of the general classifiers in Sentiment Analysis in below subsections..

3. A.1 Probabilistic Classifiers

Probabilistic classifiers use blend of models for the grouping. [24] Mixture demonstrate expect that each class is a portion of the blend. Each blend fragment provides the likelihood of inspecting a specific term for that portion and it is a inventive model [25]. This type of distributions are also called generative classifiers. There are three probabilistic classifiers which are Naïve Byes, Maximum Entropy and Bayesian Network. [24,25]. Probabilistic classifiers use mixture models for classification. The mixture model assumes that each class is a component of the mixture. This type of distributions are also called generative classifiers for the sake of every mix unit is a constructive approach that gives the chance of sampling a specific word for that unit. There exist three popular probabilistic classifiers ,namely Naive Bayes Classifier , Bayesian Network and Maximum Entropy ,which are discussed in the following subsections.

3. A. 1. 1 Naive Bayes Classifier(NB)

Naïve Byes classifier is the simplest and most commonly used classifier. This classification model performs the posterior choice of a class based on the classification of the terms in the text. Bayes Theorem is used in this classification and to estimate the choice(probability) that a given characteristic set associates to a specific tag.. $P(\text{label}|\text{features}) = \frac{P(\text{label}) * P(\text{features} | \text{label})}{P(\text{features})}$ here $P(\text{label})$ is equal to prior probability of a label. $P(\text{features}|\text{label})$ is equal to prior probability that a given characteristic set and $P(\text{feature})$ is equal to prior probability that a given characteristic set is happened.

3. A. 1. 2 Bayesian Network(NB)

The inference of the Bayesian Network classifier is the autonomy of characteristics. Another inference is to suppose that all the characteristics are totally dependent. It derives, the NB(Bayesian Network) approach, is a DAG(Directed Acyclic Graph) whose vertices denote irregular(random) variables and edges denotes decision dependencies.

3. A. 1. 3 Maximum Entropy

The Maximum Entropy is a particular classifier usually provides service in NLP(Natural Language Processing), SLP(Speech Language Processing) and in IR(Information Retrieval) related issues. It does not consider, the characteristics are provisionally independent of each other. The Maximum Entropy is situated on the basis of Maximum Entropy. It states that all the models that fit into training data and choose the largest entropy. Comparing to Naive Bayes ,The Maximum Entropy model needs huge time to train, basically because of the optimization problem. and It needs to be determined, to predict the specifications of the model.

3. A. 2 Linear Classifiers

Given $X^i = \{X_1, X_2, X_3, \dots, X_N\}$ is the normalized document word frequency, vector $A^i = \{A_1, A_2, A_3, \dots, A_N\}$ is a vector of linear coefficients with the same dimensionality as the feature space, and b is a scalar; the output of the linear predictor is defined as, $p = A^i \cdot X^i + b$ which is the output of the linear classifier. The prediction p is a separating hyper plane between different classes.

3. A. 2. 1 Support Vector Machines Classifiers (SVM)

The objective of Support Vector Machines is to get linear classifiers in the search space and it can get best classifier in the distinctive classes. In the Fig. 2, there are two classes a, b and there are 3 hyper planes A, B and C. Hyperplane gives the good classification among the classes, due to the regular distance of any of the data points is the largest and indicates the highest boundary of classification.

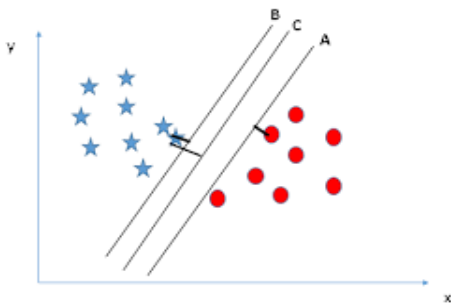


Fig-2: Using support vector machine on a classification problem.

For SVM classification, usually text data suited due to inadequate nature of text, in this case strange characteristics are extraneous, but they go for correlated with each other and naturally arranged into linear separable types. Support Vector Machine can develop a nonlinear conditional surface in the original characteristic space by set out the data illustrated non-linearly with a hyper plane [12].

3. A. 2. 2 Neural network classifier

A neural network consists of units (neurons), arranged in layers, which convert an input vector into some output. Each unit takes an input, applies a (often nonlinear) function to it and then passes the output on to the next layer. Generally the networks are defined to be feed-forward: a unit feeds its output to all the units on the next layer, but there is no feedback to the previous layer. Weightings are applied to the signals passing from one unit to another, and it is these weightings which are tuned in the training phase to adapt a neural network to the particular problem at hand. This is the learning phase. A neuron(see Fig.4) in an artificial neural network is

1. A set of input values (x_i) and associated weights (w_i).
2. A function (g) that sums the weights and maps the results to an output(y).

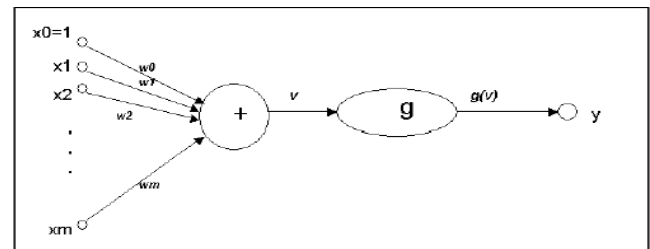


Fig-4: neuron architecture

Neurons are organized into layers: input, hidden and output. The input layer is composed not of full neurons, but rather consists simply of the record's values that are inputs to the next layer of neurons. The next layer is the hidden layer. Several hidden layers can exist in one neural network. The final layer is the output layer, where there is one node for each class. A single sweep forward through the network results in the assignment of a value to each output node, and the record is assigned to the class node with the highest value. In the training phase, the correct class for each record is known (termed supervised training), and the output nodes can be assigned correct values -- 1 for the node corresponding to the correct class, and 0 for the others. The training process normally uses some variant of the Delta Rule, which starts with the calculated difference between the actual outputs and the desired outputs. The connection weights are normally adjusted using the Delta Rule. This process proceeds for the previous layer(s) until the input layer is reached.

3. A. 3 Rule-based Technique:

In rule based technique, if a rule has "if-then" relation then it consists of an antecedent and its associated consequent. [20, 21, 22]

Antecedent -> Consequent

An antecedent describes a condition and can be either a token or a series of tokens that are concatenated by the " \wedge " operator. A token can be either "?" denoting a proper noun, a word or "#" denoting the result of the condition described by the antecedent.

3. A. 4 Decision Tree Classifiers:

A decision tree is a graphical model describing decisions and their possible outcomes. Decision trees consist of three types of nodes (see Fig.3):

1. Decision node: Often represented by squares showing decisions that can be made. Lines emanating from a square show all distinct options available at a node.
2. Chance node: Often represented by circles showing chance outcomes. Chance outcomes are events that can occur but are outside the ability of the decision maker to control.
3. Terminal node: Often represented by triangles or by lines having no further decision nodes or chance nodes. Terminal nodes depict the final outcomes of the decision making process.

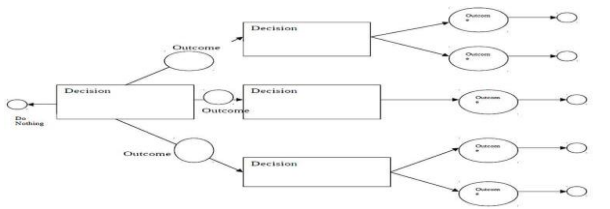


Fig 3: Decision trees are graphical models for describing sequential decision problems.

3. B Unsupervised Learning:

The main goal of the classification is to classify documents into number of categories. In large number of labeled training document, it is difficult to create labeled training document, but easy to collect the unlabeled documents. To solve this problem unsupervised learning methods are used. Koa and Seo present the research work in this field and they propose a method that divides document into sentences, and categorized each sentence using keyword lists of each category and sentence similarity measure [13].

3. C Lexicon Based Approach

In this unsupervised techniques the classification is held on features based .it is consider polarity of people's views and opinions and then analyze that it is more positive or it is negative , that's called sentiment lexicon. For compiling the sentiment word net used manually approach but it's tough and time taking process and do not perform any task alone. Its needed further two automated approaches that perform automatically sentiment .these two automated approaches are defined below Dictionary Based Approach: It first identify the polarity of the word after that find it in dictionary and its synonym and antonym and known as collection WorldNet [14], [15], [16] or word finder [17].Its follow the iterative process after finding new word then add it in stem words. Repeat this iterative process until not search new words.

3. C. 1 Corpus Based Approach :

Corpus based approach relatively perform good accuracy by solving the problem of finding specific opinion words with specific orientations ,this is the big advantage of Corpus based approach which don't have in dictionary based approach. IN the research[18] the author has used the Conditional Random Fields (CRFs) method for opinion extractions through the manually learning techniques.

3. C. 1. 1 Statistical Approach

Statistical Approach is used to search co-occurrence patterns and stem the opinion words in this technique. This is used for many related application to Sentiment analysis. For finding the polarity of big annotated word corpus [19].

3. C. 1. 2 Semantic Approach

Semantic Approach is rely on principle to analyze the similar calculation to show its semantic value to directly. This is used in different applications to make lexicon base model for the expressing of Noun, verbs, Adjective and adverbs are used in sentiment analysis as the author has used in [20].

3.C.2 Dictionary-Based Approach

Using a dictionary approach to compile sentiment words is an obvious approach because most dictionaries (e.g., WordNet (Miller, 1990)) list synonyms and antonyms for each word. Thus, a simple technique in this approach is to use a few seed sentiment words to bootstrap based on the synonym and antonym structure of a dictionary. Specifically, this method works as follows: A small set of sentiment words (seeds) with known positive or negative orientations is first collected manually, which is very easy. The algorithm then grows this set by searching in the WordNet or another online dictionary for their synonyms and antonyms.

3.D Hybrid Approach

The combination of Machine learning and lexicon based approach is called the Hybrid approach. Some research techniques have analyze that hybrid approach improves the performance in sentiment classification. In this paper author has used these both techniques together in concept level sentiment analysis and perform high accuracy [26].

Applications Areas Of Sentiment Analysis

Since the Opinion based or feedback based application are more fashionable, now a days, the natural language processing community shows much interest in Sentiment Analysis and Opinion Mining system. The explosion of internet has changed the people's life style, now they are more expressive on their views and opinions [1], and this tendency helped the researchers in getting user-generated content easily. The major applications of Opinion mining and sentiment analysis are the following:

- 1) Purchasing Product or Service
- 2) Quality Improvement in Product or service
- 3) Marketing research
- 4) Recommendation Systems
- 5) Detection of "flame"
- 6) Policy Making
- 7) Decision Making

Research Scope In Sentiment Analysis

The major research scope areas in sentiment analysis are:

- Spam Detection Sentiment Analysis;
- Sentiment Analysis on short Sentence like abbreviations;
- Improving sentiment word identification algorithm;
- Developing fully automatic analyzing tool;
- Effective Analysis of policy opinionated content;
- Successful handling of bi polar sentiments;
- Generation of highly content lexicon database.

4. CONCLUSIONS

Thus, Opinion Mining and Sentiment analysis has wide area of applications and it also facing many research challenges. Since the fast growth of internet and internet related applications, the Opinion Mining and Sentiment Analysis become a most interesting research area among natural language processing community. A more innovative and effective techniques required to be invented which should overcome the current challenges faced by Opinion Mining and Sentiment Analysis.

5. REFERENCES

- [1]. Boiy, E. Hens, P., Deschacht, K. & Moens, M.F., "Automatic Sentiment Analysis in Online Text," In Proceedings of the Conference on Electronic Publishing (ELPUB-2007), pp. 349-360, 2007.
- [2]. Liu, B., "Sentiment analysis and subjectivity," In Handbook of Natural Language Processing, Second Edition, N. Indurkha and F. J. Damerau, Eds. CRC Press, Taylor and Francis Group, Boca Raton, FL. ISBN 978-1420085921, 2010.
- [3]. Pang, B. & Lee, L., "Opinion Mining and Sentiment Analysis", Foundations and Trends in Information Retrieval, Vol. 2 (1-2), pp. 1 - 135, 2008.
- [4]. G. Gebremeskel, "Sentiment Analysis of Twitter Posts about News," Master's Thesis, University of Malta, May 2011.
- [5]. S. Baccianella, A. Esuli, and F. Sebastiani, "SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining," Proceedings of the Seventh conference on International Language Resources and Evaluation, 2010, pp. 2200-2204
- [6]. B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up?: sentiment classification using machine learning techniques," Proceedings of the ACL-02 conference on Empirical methods in natural language processing, vol.10, 2002, pp. 79- 86.
- [7]. K. Dave, S. Lawrence, and D. M. Pennock, "Mining the peanut gallery: Opinion extraction and semantic classification of product reviews," Proceedings of WWW, 2003, pp. 519- 528
- [8]. Xing Fang and Justin Zhan "sentiment analysis using product Review data" Department of computer science, North Carolina a&T State University Greensboro, NC, USA, 2015 Springer journal.
- [9]. Bogdan Batrinc, Philip C. Treleaven "Social media analytics: a survey of techniques, tools and platforms Department of computer science, University College London, Gower Street, London WC1E 6BT, UK Published on 26 July 2014.
- [10]. Joachims T. ,"Probabilistic analysis of the rocchio algorithm with TFIDF for text categorization", In: presented at the ICML conference; 1997.
- [11]. Aizerman M, Braveman E, Rozonoer L., "Theoretical foundations of the potential function method in pattern recognition learning",. Autom Rem Cont:821-37,1964.
- [12]. Koyoungioong, SeoJungyun, "Automatic text categorization by unsupervised learning",. In: Proceeding of COLING-00 the 18th international conference on computational linguistics;2000.
- [13]. Miller G, Beckwith R, Fellbaum C, Gross D, Miller K. WordNet: an on-line lexical database. Oxford Univ. Press; 1990.
- [14]. M. Kanakaraj, R. Mohana, and R. Guddeti, "NLP Based Sentiment Analysis on Twitter Data Using Ensemble Classifiers," in 3rd International Conference on Signal Processing, Communication and Networking (ICSCN), 2015.
- [15]. S. Gao, J. Hao, Y. Fu, "The Application and Comparison of Web Services for Sentiment Analysis in Tourism."
- [16]. Mohammad S, Dunne C, Dorr B. Generating high coverage semantic orientation lexicons from overly marked words and a thesaurus. In: Proceedings of the conference on Empirical Methods in Natural Language Processing (EMNLP'09); 2009.
- [17]. Lafferty J, McCallum A, Pereira F. Conditional random fields: probabilistic models for segmenting and labeling sequence data. In: Proceedings of International Conference on Machine Learning (ICML'01); 2001.
- [18]. Read J, Carroll J. Weakly supervised techniques for domain independent sentiment classification. In: Proceeding of the 1st international CIKM workshop on topic-sentiment analysis for mass opinion; 2009. p. 45- 52.
- [19]. Maks Isa, Vossen Piek. A lexicon model for deep sentiment analysis and opinion mining applications. Decis Support Syst 2012;53:680-8
- [20]. Chin-Shrng Yang, Hsiao-Ping Shih, Department of Information Management, Yuan Ze University, ChangLi, Taiwan, "A Rule-Based Approach For Effective Sentiment Analysis" PACIS 2012.
- [21]. Prern Chikersal, Soujanya Poria, and Erik Cambria, School of Computer Engineering Nanyang Technological University Singapore-639798, "SeNTU: Sentiment Analysis of Tweets By Combining a Rulebased Classifier with Supervised Learning" June 5 2015.
- [22]. Neethu M S Rajasree R Department of Computer Science and Engineering "Sentiment Analysis in Twitter using Machine Learning Techniques" IEEE 2013.
- [23]. Walaa Meddhat, Ahmed Hassan, Hoda Korashy "Sentiment analysis algorithms and applications: A survey, Ain Sham University, Faculty of Engineering, Computer & Systems Department, Egypt 19 April 2014.
- [24]. Bogdan Batrinc, Philip C. Treleaven "Social media analytics: a survey of techniques, tools and platforms Department of computer science, University College London, Gower Street, London WC1E 6BT, UK Published on 26 July 2014.
- [25]. A. Mudinas, D. Zhang, M. Levene, "Combining lexicon and learning based approaches for concept level sentiment analysis", Proceedings of the First International Workshop on Issues of Sentiment Discovery and Opinion Mining, ACM, New York, NY, USA, Article 5, pp. 1-8, 2012.