



# Comparison of Data Analytic Techniques for a Spatial Opinion Mining in Literary Works: A Review Paper

Sea Yun Ying<sup>1</sup>, Pantea Keikhosrokiani<sup>1</sup>(✉), and Moussa Pourya Asl<sup>2</sup>

<sup>1</sup> School of Computer Sciences, Universiti Sains Malaysia, 11800 Minden, Penang, Malaysia  
pantea@usm.my

<sup>2</sup> School of Humanities, Universiti Sains Malaysia, 11800 Minden, Penang, Malaysia

**Abstract.** Opinion mining is the use of analytic methods to extract subjective information. A study was conducted to apply spatial opinion mining in literary works to examine the writers' opinions about how matters of space and place are experienced. For this reason, this paper conducts a review study to identify and compare different analytical techniques for opinion mining in fictional writings. This review study focused on sentiment analysis and topic modeling as two main techniques for spatial opinion mining in literary works. The comparison results are reported and the limitations of different techniques are mentioned. The results of this study can assist researchers in the field of opinion and text mining.

**Keywords:** Big data analytics · Opinion mining · Text mining · Sentiment analysis · Topic modeling · Literary works

## 1 Introduction

Opinion Mining (OM) is a technique that is used to detect and extract the prevalent opinion about entities. It utilizes text mining to detect the sentiment orientation of a text which could be positive, negative or neutral. It can be described as a fiend of knowledge discovery and data mining (KDD) that applies Natural language processing (NLP) and statistical machine learning (ML) techniques to categorize opinionated text from factual text. Therefore, the OM task involves opinion identification, opinion classification, target identification, source identification and opinion summarization as stated by [1].

Text analytics or text mining, is the methodology and process that allows machines to derive quality information and insights from textual data. This process involves using NLP, information retrieval and (ML) techniques to parse unstructured text data into more structured forms and extracting patterns and information from this kind of data that might bring benefits to the end user [2]. Human analysis of textual information is subject to prejudice and bias because people tend to give opinions that are consistent with their preferences. It is commonly believed that individuals can make decisions based on a rational analysis of available alternatives. However, it cannot be ignored that emotions exert a profound impact on the decisions that humans make in reality. Emotion is the

main ingredient that cannot be neglected in any decision making processes. Furthermore, humans have difficulty in producing consistent results when the amount of information to be processed is huge. Thus, automated text mining and summarization systems are required to overcome subjective biases and human limitations with an objective sentiment analysis system [3].

For this reason, a study was conducted to propose a computerized method to examine literary writers' opinions in how matters of space and spatiality are addressed in their fictional works [4, 5]. Therefore, a review study was required to compare different data analytic techniques used to find how spatial experiences are portrayed in certain literary texts by diasporic writers [6, 7]. Therefore, this paper focused on the review and comparison of data analytic techniques for opinion mining in literary writings.

## 2 Methods for a Spatial Opinion Mining in Literary Works

### 2.1 Natural Language Processing and Text Analytics

Emergence of big data technologies and artificial intelligence transformed the way researches from various disciplines are conducted [8–10]. Natural Language Processing (NLP) can be defined as a subfield of data science and Artificial Intelligence (AI) with roots in computational linguistics. The field is closely related to AI and was useful in the 1950's automatic translation efforts [11]. It is mainly concerned with designing and constructing applications and systems that enable interaction between machines and natural languages evolved for utilization by humans [2]. NLP techniques enable computers to understand and manipulate natural language texts or speeches to perform useful outputs [12].

Text analytics, also referred to as text mining, is the methodology and process that allows machines to derive quality information and insights from textual data. Text analytics can be defined as an application which applies text mining techniques to sort out data sets. This process involves using NLP, information retrieval and machine learning (ML) techniques to parse unstructured text data into more structured forms and extracting patterns and information from these kinds of data that might bring benefits to the end user [2]. Text mining has become popular due to the development of big data platforms and deep learning algorithms that are able to analyze massive sets of unstructured data.

There are two sub-domains for mining knowledge from user-generated discourses (subjectivity analysis): opinion mining and sentiment analysis [1]. Some of the researchers agree to use these domains interchangeably [13], while some of them have considered sentiment analysis to be a subfield of OM [14].

### 2.2 Opinion Mining

Opinion Mining (OM) can be defined as the science of using text mining to detect the sentiment orientation of a text which could be positive, negative or neutral. The term OM was first mentioned in 2003 by Dave et al. who described it as analysis of reviews about entity and presented it as a model for document polarity classification for recommended or not recommended. OM is a process used to extract information

or opinion about entities. It can be defined as a field of knowledge discovery and data mining (KDD) that applies NLP and statistical ML techniques to categorize opinionated text from factual text. Therefore, the OM task involves opinion identification, opinion classification, target identification, source identification and opinion summarization [1]. The main concern that Khan et al. mentioned in their paper is how to automatically detect opinion components from unstructured text data and summarize the opinion about an entity from a large volume of unstructured text data.

### 2.3 Topic Modelling

Topic Modelling is one of the unsupervised techniques used to perform text clustering in large document collections. It is a statistical model that helps to search a group of keywords or topics for a text. It assumes that each document consists of a group of topics or keywords. Each topic or keyword in the text consists of a collection of words. It is a form of opinion mining that is able to obtain recurring patterns of words in a textual document. Latent Dirichlet Allocation (LDA) and Latent Semantic Analysis (LSA) are two of the widely used topic modelling techniques nowadays. Topic Modelling is one type of probabilistic generative model that is usually used for discovering the hidden semantic structures in a text body. It is also used for annotating documents based on its topics and these annotations can be used to organize, search and summarize the whole texts. In short, topic modelling can handle the tasks of tag recommendation, text categorization, keyword extraction, and similarity search. It can be described as an approach for searching a set of words or topics from a large document collection that best describes the information in the collection. All topic models have the basic assumption in which each document contains a mixture of topics/keywords and each topic/keyword contains a collection of words [15].

### 2.4 Sentiment Analysis (SA)

Sentiment Analysis (SA) is one of the recent techniques used to extract and analyze emotional and sentiment statements in a text [16–18]. It can be referred to as emotional polarity computation as it is used to detect the sentiments and categorize them based on their polarity. The used polarities can be positive, neutral or negative. In this case, traditional machine learning techniques on n-grams, parts of speech, and other bag of words features able to be applied when the data is labelled. Knowledge-based method that was introduced by [19] is another method in using labelled data. Both of these methods rely on crowdsourcing. Sentiment analysis (SA) is related to the extraction and analysis of emotional and sentimental statements in a text. The aim of SA is to detect opinions, identify the sentiments they express, and then categorize them based on their polarity. It uses polarities such as positive or negative or a scale of ratings (e.g., 1–5) and relies on features about emotion, affect, review and subjectivity.

SA is a classification process which is divided into three main levels: document-level [20], sentence-level [21], and aspect-level [22]. The purpose of document-level SA is to categorize an opinion document as expressing a positive or negative opinion. In this case, the whole document is considered as a basic information unit (topic). The aim of sentence-level SA is to classify sentiment expressed in every sentence. It is

necessary to determine whether the sentence is subjective or objective before determining whether the sentence expresses positive or negative opinions. However, document level and sentiment level SA do not find out what exactly people preferred or did not prefer. Aspect-level SA performs finer-grained analysis. The goal of performing an aspect-level SA is to determine the sentiment with respect to the specific aspects of entities. Aspect-level directly looks at the opinion rather than looking at language constructs which include documents, paragraphs, sentences, clauses or phases. Sentiment analysis is generally carried out by three common approaches: Lexicon-based approach, Learn-based approach or Machine Learning approach, and Hybrid approach.

### 3 Comparison and Discussion Between Sentiment Analysis and Topic Modelling Approaches

#### 3.1 Comparison and Discussion on Different Sentiment Analysis Approaches

Performing SA with different approaches produces different results. Each approach has its own advantages and disadvantages. The comparison and discussion of the two main approaches used in SA is given in Table 1.

**Table 1.** Comparison of two approaches

Criteria	Lexicon-based	Learn-based
<i>Classification</i>	Unsupervised learning	Supervised, Semi-supervised and Unsupervised learning
<i>Advantages</i>	<ul style="list-style-type: none"> <li>• Domain independent</li> <li>• Does not need labelled data and the procedure of learning</li> <li>• Fast time to get the results</li> </ul>	<ul style="list-style-type: none"> <li>• Dictionary is not necessary</li> <li>• High accuracy of classification</li> <li>• High precision and adaptability</li> <li>• Do not need maintenance</li> </ul>
<i>Disadvantages</i>	<ul style="list-style-type: none"> <li>• Need maintenance for corpus/corpora</li> <li>• Requires strong linguistic resources which is not always available</li> <li>• Needs dictionaries that covers plenty of opinion words</li> <li>• Low accuracy if compared with learn-based approach</li> </ul>	<ul style="list-style-type: none"> <li>• Dependent to the domain so classifier trained on the texts in one domain does not work with other domains</li> <li>• Needs labelled data and procedure of learning</li> <li>• Slow time to get the results</li> </ul>

**Lexicon-based Approach:** Lexicon-based approach can be referred to as a rule-based approach because the dictionaries are used following certain rules. It depends on finding the opinion lexicon which is used to analyze a text. It relies on a sentiment lexicon, a collection of known and precompiled sentiment terms. Lexicon-based approach is one of the methods to do SA in document level, and it can be considered as an unsupervised approach as mentioned by [13].

Lexicon-based can be grouped into dictionary-based approach, corpus-based approach and manual approach. The dictionary-based method depends on searching opinion seed words before looking for the dictionary of their synonyms and antonyms. Dictionary based method can use existing dictionaries such as SentiWordNet.

The corpus-based method starts with a seed list of opinion words before looking for other opinion words in a large corpus to help in searching opinion words with context specific orientations. This approach needs expensive manual annotation effort because it involves large corpus as mentioned by [23]. Corpus-based method is not as effective as applying a dictionary-based approach alone because it is difficult to prepare a large corpus to cover all English words [24]. Manual approach is a very time-consuming method and it is usually combined with dictionary-based and corpus-based approaches to prevent the mistakes that result from dictionary-based and corpus-based methods.

**Recent Studies of Lexicon-based Approach in SA:** Lexicon-based approach is an unsupervised learning in that it does not need prior training for mining data. Most researches create their own lexicon in order to improve the performance. Several recent studies of lexicon-based approach in SA are summarized in Table 2.

**Learn-based Approach:** Learn-based or ML Approach uses the famous ML algorithms and linguistic features. It is a classification algorithm which trains labelled document over corpus so that the features can be recognized for classifying the sentiment (Gianakopoulos et al., 2012). This approach can be supervised, semi-supervised or unsupervised. Supervised methods need large number of labelled training data which makes them expensive while unsupervised methods do not need labelled data and are therefore easy to apply for unlabeled data. There are several recent studies in SA using learn-based approach that are summarized in Table 3.

Table 4 shows the advantages and disadvantages of different learn-based methods. It is a guideline to decide which method is and can be used. The machine learning methods given in the table below is the most common used method.

**Hybrid Approach:** The hybrid approach combines both lexicon-based and learn-based approaches. It is very common with sentiment lexicons and plays a crucial role in the majority of methods. It applies the lexicon-based approach for sentiment score and then these scored documents represent the training data for the learn-based part. According to [29], hybrid approach is widely adopted due to its high accuracy and stability inherited from lexicon based approach and ML approach respectively. It uses a lexicon or learning symbiosis to attain the best of both worlds-stability and readability from a carefully designed lexicon, while the high accuracy from a powerful supervised algorithm.

### 3.2 Comparison of Topic Modelling: Latent Dirichlet Allocation (LDA) and Latent Semantic Analysis (LSA)

Latent Dirichlet Allocation (LDA) and Latent Semantic Analysis (LSA) are two common topic modelling techniques that are used nowadays. Both of them are able to find hidden topics from given documents without labelled training data. Discovering hidden topics provide benefits for different purposes such as clustering documents and

**Table 2.** Comparative Study of lexicon-based approach in SA

Author	Dataset	Techniques	Tool	Accuracy
[25]	Newspaper articles (the set of 1292 quotes)	WordNet- lexicon based	WordNet Affect, SentiWordNet, MicroWNOp, JRC Tonality	82% improve the baseline 21%
[26]	40 216 tweets from Stanford Twitter Sentiment, and 3 269 tweets from ObamaMcCain Debate	General Inquirer LBM, MPQA, K-means, ONMTF, Moodlens, CFMS and ESSA	A novel framework make use of ESSA	General inquirer with improvement 21.4% and 17.87% for both datasets
[27]	Online customers reviews (Spam & fake reviews)	Combine lexicon and use shallow dependency parser	SentiWordNet and MPQA lexicon	85.7% for sentiment method but word counting approach 76.7%
[28]	Data set of 1 600 Facebook messages	n-gram (uni and bi-grams)	Lexicon of Hu and Liu (HL) and MPQA lexicon	70%
[29]	Three datasets (training set, test set and the verified set) 1000, 5000, and 10 000	Enhancement BOW model	New lexicon	83.5%
[30]	335 022 restaurant reviews	Multilevel model	AFINN sentiment lexicons	Find out relation between the words
[31]	3 000 tweets (movie tweets)	Lexicon based technique	TextBlob, SentiWordNet and WSD	TextBlob results were relatively better
[32]	6 250 tweets (political views)	Lexicon based technique and ML approach to check the accuracy	TextBlob, SentiWordNet and W-WSD	62.67% (TextBlob); 53.33% (SentiWordNet); 62.33% (WSD)
[33]	11 861 sentence-level snippets (Movie reviews)	Lexicon based technique	VADER, Textblob and NLTK	77% (VADER); 74% (Textblob); 62% (NLTK)
[34]	1 828 patients' opinions in healthcare	Lexicon based technique	VADER and TextBlob	71.9% (VADER); 73% (TextBlob)

**Table 3.** Comparative study of learn-based approach in SA

Author	Dataset	Techniques	Tool	Accuracy
[35]	70 103 hotel reviews	Naïve Bayes	Natural language toolkit (NLTK)	Find out relationship between sentiment variables and driving number of reviews
[36]	185 English song lyrics (textual)	Naïve Bayes, KNN, SVM	RapidMiner	SVM based classifiers indicate promising results
[37]	1 200 electronic product review	Naïve Bayes, SVM, Maximum Entropy and ensemble classifier	Matlab simulator	SVM and MEM have equal accuracy of 90%; NB has 89.5% but NB has better precision
[38]	56 483 restaurant reviews	Naïve Bayes algorithm	Ictcla50	74%
[39]	24 000 sentences of multi domain sentiment data (12 domain) from 2 000 reviews	Naïve Bayes, Binarized Multinomial Naïve Bayes (BMNB), Multinomial Naïve Bayes, SVM and J48	NA	BMNB has the best performance in six out of twelve domain, followed by SVM in four out of twelve data domain, the best feature selection is Information Gain
[40]	1 346 545 business reviews	SVM, Naïve Bayes, Logistic Regression, SGD	Natural language toolkit (NLTK)	Linear SVC and SGD have an accuracy of 94.4%; NB and Logistic Regression tend to have slightly worst results
[41]	Online Hotel Reviews: 19 650 data in Chinese and 2 008 in English	Multinomial Naïve Bayes	Natural language toolkit (NLTK)	Find out general consistent interrelationship between structured and unstructured UGC

*(continued)*

**Table 3.** (continued)

Author	Dataset	Techniques	Tool	Accuracy
[42]	Hotel reviews of 3 hotels located in Astana	Language processing SVM, MEM	NLTK	Extended Model able to express more accurate sentiment polarity
[43]	13 541 tweets from E-Twitter and 479 tweets from Twitter-sanders	SVM, Decision Tree and Naïve Bayes	WEKA	SVM provides more accurate results than DT and NB
[44]	4 datasets: 1 600 000 tweets; 888 tweets; 10 729 tweets; 99 989 tweets	Naïve Bayes, Random Forest, SVM, Logistic Regression, Majority Voting, and a proposed ensemble	NA	Proposed ensemble classifier performs better than the other classifier

**Table 4.** Comparison of different machine learning methods

Methods	Advantages	Disadvantages
<i>SVM</i>	<ul style="list-style-type: none"> <li>• High dimensional input space</li> <li>• Involve few inappropriate features</li> <li>• Document vectors are scarce</li> </ul>	<ul style="list-style-type: none"> <li>• A huge amount of training data set is needed</li> <li>• Data collection process is tedious</li> </ul>
<i>NB</i>	<ul style="list-style-type: none"> <li>• Simple and intuitive approach</li> <li>• It is an efficiency method with reasonable accuracy</li> </ul>	<ul style="list-style-type: none"> <li>• Mainly used for small training set</li> <li>• It always assume conditional independence among the linguistic features</li> </ul>
<i>MEM</i>	<ul style="list-style-type: none"> <li>• It does not assume the independent features like NB</li> <li>• Able to handle large amount of data</li> </ul>	<ul style="list-style-type: none"> <li>• Simplicity is hard</li> </ul>
<i>KNN</i>	<ul style="list-style-type: none"> <li>• It is computationally efficient</li> <li>• Classification of an instance will be similar to those closer to it in the vector space</li> </ul>	<ul style="list-style-type: none"> <li>• Big storage is needed</li> <li>• Computationally intensive recall</li> </ul>

organizing online available texts for information retrieval and are also able to provide accurate recommendations. Latent Semantic Analysis (LSA) and Latent Dirichlet Allocation (LDA) are two common text data computer algorithms that have gained much attention individually in the text analysis for topic extraction studies but are not famous for neither document classification nor comparison studies [45]. Based on Anaya's study, the accuracy rates for both LDA and LSA at the high level of abstraction are 84% while



the accuracy of LDA and LSA are 64% and 67% respectively at the lower level of abstraction. [46] and [47] summarize LDA and LSA as given in Table 5.

**Table 5.** Comparison of LDA and LSA

Criteria	LDA	LSA
<i>Characteristics</i>	<ul style="list-style-type: none"> <li>• A generative probabilistic model</li> <li>• Need manually remove stop words to increase the performance of model</li> <li>• Words are categorized into topics and can exist in more than one topic</li> </ul>	<ul style="list-style-type: none"> <li>• Examines words existed in a document with same or similar meaning</li> <li>• Low dimension representation of documents and words by applying SVD</li> <li>• Create latent semantic space</li> </ul>
<i>Advantages</i>	<ul style="list-style-type: none"> <li>• Works well with a huge corpus</li> <li>• Avoid from issue overfitting</li> <li>• Able embedded in other complicated models</li> <li>• Noise reduction is possible</li> <li>• Easier to implement than LSA</li> <li>• Best learns descriptive topics</li> </ul>	<ul style="list-style-type: none"> <li>• Able to cluster the words and documents in the space</li> <li>• Unable to capture the multiple meanings of words</li> <li>• Error reduction available by using dimension reduction</li> <li>• Best at creating a compact semantic representation of documents and words in a corpus</li> </ul>
<i>Disadvantages</i>	<ul style="list-style-type: none"> <li>• The number of topics should be set in advance</li> <li>• Uncorrelated topics</li> </ul>	<ul style="list-style-type: none"> <li>• Not presenting well defined probabilities</li> <li>• Lack of interpretable embedding</li> <li>• Less efficient presentation</li> <li>• Offer lower accuracy then LDA</li> </ul>

The primary points that need to be taken into consideration when using a topic model technique is the degree to which the learned topics match human judgments and are able to help humans differentiate between ideas, as suggested by [46]. The evaluation of the topic model has been ad hoc and application-specific. The existing evaluations range from fully-automated intrinsic evaluations to manually crafted extrinsic evaluations. Extrinsic evaluations are normally hand constructed and are often expensive to perform for domain-specific topics while intrinsic evaluations are able to evaluate the amount of information encoded by the topics easily.

Perplexity is one of the common examples of intrinsic evaluations as suggested by [48] to evaluate the performance of the topic model technique. [49] found that perplexity may not yield human interpretable topics. As a result, researchers have introduced topic coherence measures – a qualitative approach to automatically discover the coherence of a topic [34, 50–52]. A number of measures have been combined into a framework in order to evaluate the coherence between topics inferred by a model. The degree of semantic similarity between topic-related words in the topic is measured by using topic coherence measures [53]. The higher the topic coherence score, the more the semantically meaningful topic is generated [46].

## 4 Conclusion

This study focused on reviewing different sentiment analysis and topic modeling techniques which are suitable for a spatial opinion mining in literary works. Based on the results of this review study, LDA is an unsupervised learning approach that is applied in this study due to lack of labelled training data. Unsupervised learning is normally used for finding hidden patterns of data to improve the performance of the model. In other words, it might not be used alone but combined with supervised learning approach in order to achieve a higher quality of model. Unsupervised learning approach can be conducted when (1) there are no labels on training data; (2) the data cannot be labelled manually or it is expensive to do so; and (3) most of the supervised learning algorithms fail to fit well with the underlying distribution of the data renders. In this study, the first and the second criteria are matched so LDA is used. However, it will never be the first choice if a big and good quality labelled training data is provided. Lexicon-based approach is performed for the sentiment analysis task in this project. It is a method that is designed for all domains as it is unsupervised and domain independent. In short, it can achieve a more robust performance across domains than learn-based approach. However, the limitation of the lexicon-based approach is that the maintenance of sentiment lexicons for different domains is needed. Lexicon-based performs faster than learn-based method but the accuracy of lexicon-based method is always lower than the learn-based method. In order to solve these kinds of problems, a big and good quality of labelled data should be provided.

For future studies, the machine learning approach is suggested to be applied. The main reason is the limitation of unsupervised learning and that the lexicon-based approach of sentiment analysis can be ignored if ML approach is applied. However, the labelled training data set should be provided. The lack of data and lack of good data will generate a poor model. Most of the available machine learning algorithms require large amounts of data before they start to build a model. The good data ensures the model to capture good features from the training data set and then the algorithm will perform well.

**Acknowledgment.** The authors are thankful to School of Computer Sciences and School of Humanities, Universiti Sains Malaysia for unlimited supports to finish this project. In addition, the authors are grateful to Division of Research & Innovation, USM for financial support from Short Term Grant (304/PHUMANITI/6315300) granted to Dr Moussa Pourya Asl.

## References

1. Khan, K., et al.: Mining opinion components from unstructured reviews: a review. *J. King Saud Univ. – Comput. Inf. Sci.* **26**(3), 258–275 (2014)
2. Sarkar, D.: *Text Analytics With Python: A Practical Real-World Approach to Gaining Actionable Insights from Your Data*. Apress, New York (2016)
3. Lum, K.: Limitations of mitigating judicial bias with machine learning. *Nat. Hum. Behav.* **1**(7), 0141 (2017)
4. Asl, M.P.: The politics of space: vietnam as a communist heterotopia in Viet Thanh Nguyen's the refugees. *Lang. Linguist. Lit.* **26**(1), 156–170 (2020)

5. Asl, M.P.: Micro-Physics of discipline: Spaces of the self in middle Eastern women life writings. *Int. J. Arabic-English Studies* **20**(2), 223 (2020)
6. Asl, M.P.: Leisure as a space of political practice in Middle East women life writings. *GEMA Online®. J. Lang. Stud.* **19**(3), 43–56 (2019)
7. Asl, M.P.: Practices of counter-conduct as a mode of resistance in Middle East women's life writings. *Lang. Linguist. Lit.* **24**(2), 195–205 (2018)
8. Keikhosrokiani, P.: Chapter 1 - Introduction to Mobile Medical Information System (mMIS) Development, in *Perspectives in the Development of Mobile Medical Information Systems*, P. Keikhosrokiani, Editor. 2020, Academic Press pp. 1–22 (2020)
9. Keikhosrokiani, P., *Perspectives in the Development of Mobile Medical Information Systems: Life Cycle, Management, Methodological Approach and Application*, Academic Press, Cambridge (2019)
10. Abdelrahman, O., Keikhosrokiani, P.: Assembly line anomaly detection and root cause analysis using machine learning. *IEEE Access* **8**, 189661–189672 (2020)
11. Hilborg, P.H., Nygaard, E.B.: Viability of sentiment analysis in business. 2015, The Copenhagen Business School. <http://studenttheses.cbs.dk>
12. Chowdhary, K.R.: Natural language processing. In: Chowdhary, K.R. (ed.) *Fundamentals of Artificial Intelligence*, pp. 603–649. Springer India, New Delhi (2020)
13. Liu, B.: Sentiment analysis and opinion mining. *Synth. Lect. Hum. Lang. Technol.* **5**(1), 1–167 (2012)
14. Tang, H., Tan, S., Cheng, X.: A survey on sentiment detection of reviews. *Expert Syst. Appl.* **36**(7), 10760–10773 (2009)
15. Kumar, S.A., et al.: Computational intelligence for data analytics. In: *Recent Advances in Computational Intelligence*, Springer. pp. 27–43 (2019)
16. Bakshi, R.K., et al.: Opinion mining and sentiment analysis. In: *2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom)*. IEEE (2016)
17. Ravi, K., Ravi, V.: A survey on opinion mining and sentiment analysis: tasks, approaches and applications. *Knowl.-Based Syst.* **89**, 14–46 (2015)
18. Li, N., Wu, D.D.: Using text mining and sentiment analysis for online forums hotspot detection and forecast. *Decis. Supp. Syst.* **48**(2), 354–368 (2010)
19. Andreevskaia, A., Bergler, S.: CLaC and CLaC-NB: Knowledge-based and corpus-based approaches to sentiment tagging. In: *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)* (2007)
20. Yessenalina, A., Yue, Y., Cardie, C.: Multi-level structured models for document-level sentiment classification. In: *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*. (2010)
21. Farra, N., et al.: Sentence-level and document-level sentiment mining for Arabic texts. In: *2010 IEEE International Conference on Data Mining Workshops* (2010)
22. Zhou, H., Song, F.: Aspect-level sentiment analysis based on a generalized probabilistic topic and syntax model (2015)
23. He, Y., Zhou, D.: Self-training from labeled features for sentiment analysis. *Inf. Process. Manag.* **47**(4), 606–616 (2011)
24. Medhat, W., Hassan, A., Korashy, H.: Sentiment analysis algorithms and applications: a survey. *Ain Shams Eng. J.* **5**(4), 1093–1113 (2014)
25. Balahur, A., et al.: Sentiment analysis in the news. *arXiv preprint arXiv:1309.6202* (2013)
26. Hu, X., et al.: Unsupervised sentiment analysis with emotional signals. In: *Proceedings of the 22nd International Conference on World Wide Web* (2013)
27. Peng, Q., Zhong, M.: Detecting spam review through sentiment analysis. *JSW* **9**(8), 2065–2072 (2014)

28. Flekova, L., Preotiuc-Pietro, D., Ruppert, E.: Analysing domain suitability of a sentiment lexicon by identifying distributionally bipolar words. In: *Proceedings of the 6th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis* (2015)
29. El Alaoui, I., et al.: A novel adaptable approach for sentiment analysis on big social data. *J. Big Data* **5**(1), 12 (2018)
30. Gan, Q., et al.: A text mining and multidimensional sentiment analysis of online restaurant reviews. *J. Qual. Assur. Hosp. Tourism* **18**(4), 465–492 (2017)
31. Gupta, M., Sharma, P.: Sentimental Analysis of Movies Tweets with Different Analyzer
32. Hasan, A., et al.: Machine learning-based sentiment analysis for twitter accounts. *Math. Comput. Appl.* **23**(1), 11 (2018)
33. Bonta, V., Janardhan, N., Kumares, N.: A Comprehensive study on lexicon based approaches for sentiment analysis. *Asian J. Comput. Sci. Technol.* **8**(S2), pp. 1–6 (2019)
34. RamyaSri, V., et al.: Sentiment analysis of patients' opinions in healthcare using lexicon-based method
35. Duan, W., et al.: Mining online user-generated content: using sentiment analysis technique to study hotel service quality. In: *2013 46th Hawaii International Conference on System Sciences* (2013)
36. Kumar, V., Minz, S.: Mood classification of lyrics using SentiWordNet. In: *2013 International Conference on Computer Communication and Informatics* (2013)
37. Neethu, M.S., Rajasree, R.: Sentiment analysis in twitter using machine learning techniques. In: *2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT)* (2013)
38. Chen, R.Y., Guo, J.Y., Deng, X.L.: Detecting fake reviews of hype about restaurants by sentiment analysis. In: *Web-Age Information Management*. Cham: Springer International Publishing (2014)
39. Saad, F.: Baseline evaluation: an empirical study of the performance of machine learning algorithms in short snippet sentiment analysis. In: *Proceedings of the 14th International Conference on Knowledge Technologies and Data-driven Business* (2014)
40. Salinca, A.: Business reviews classification using sentiment analysis. In: *2015 17th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC)* (2015)
41. Zhang, X., et al.: Sentimental interplay between structured and unstructured user-generated contents: An empirical study on online hotel reviews. *Online Inf. Rev.* **40**(1), 119–145 (2016)
42. Yergesh, B., Bekmanova, G., Sharipbay, A.: Sentiment analysis on the hotel reviews in the Kazakh language. In: *2017 International Conference on Computer Science and Engineering (UBMK)* (2017)
43. Mathur, R.: Analyzing sentiment of twitter data using machine learning algorithm. *GADL J. Invent. Comput. Sci. Commun. Technol. (JICSCT)* **4**(2), 1–7 (2018)
44. Saleena, A.N.: An ensemble classification system for twitter sentiment analysis. *Procedia Comput. Sci.* **132**, 937–946 (2018)
45. Anaya, L.H.: Comparing Latent Dirichlet Allocation and Latent Semantic Analysis as Classifiers: ERIC (2011)
46. Stevens, K., et al.: Exploring topic coherence over many models and many topics. In: *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning* (2012)
47. George, M., Soundarabai, P.B., Krishnamurthi, K.: Impact of topic modelling methods and text classification techniques in text mining: a survey. *Int. J. Adv. Electron. Comput. Sci.* **4**(3) (2017)
48. Wallach, H.M., et al.: Evaluation methods for topic models. In: *Proceedings of the 26th Annual International Conference on Machine Learning* (2009)

49. Chang, J., et al.: Reading tea leaves: how humans interpret topic models. In: *Advances in Neural Information Processing Systems*. (2009)
50. Aletras, N., Stevenson, M.: Evaluating topic coherence using distributional semantics. In: *Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013)–Long Papers* (2013)
51. Lau, J.H., Newman, D., Baldwin, T.: Machine reading tea leaves: automatically evaluating topic coherence and topic model quality. In: *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics* (2014)
52. Röder, M., Both, A., Hinneburg, A.: Exploring the space of topic coherence measures. In: *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining* (2015)
53. Korenčić, D., Ristov, S., Šnajder, J.: Document-based topic coherence measures for news media text. *Expert Syst. Appl.* **114**, 357–373 (2018)