

# CORPUS-BASED APPROACHES FOR SENTIMENT ANALYSIS: A REVIEW

## ABSTRACT

The investigation studied the state of the art of corpus-based approaches for sentiment analysis. Thus, detailing its methodologies, evaluation metrics, limitations, and future directions. The importance of sentiment analysis in fields such as marketing, customer feedback analysis, social media monitoring, financial analysis, and political science is emphasized. The methodology for corpus-based approaches in sentiment analysis includes the following key steps: data collection, preprocessing, feature extraction, and sentiment classification. The lexicon-based approaches include the corpus-based or bag of words (BOW) and dictionary (also called opinion lexicon). Evaluation of the corpus-based sentiment analysis approach is addressed through performance metrics such as accuracy, precision, recall, F1-score, and comparative analysis with other approaches including hybrid and rule-based systems. Limitations of corpus-based sentiment analysis, such as data sparsity and domain adaptation, are acknowledged, alongside potential enhancements and research directions including ensemble learning, deep learning architectures, and multimodal data integration. The conclusion emphasizes the versatility and scalability of corpus-based sentiment analysis, while ongoing research efforts aim to address its limitations and further enhance its applicability in diverse domains.

**KEYWORDS:** Corpus-based approach, sentiment analysis, NLP, feature extraction

## 1. Introduction

The definition of sentiment analysis is clarified as the automatic extraction and classification of sentiments from textual data, categorizing them as positive, negative, or neutral using natural language processing (NLP) and machine learning techniques.

Sentiment analysis, also known as opinion mining, is a computational technique under Natural Language Processing (NLP) used to determine the sentiment, emotion or opinion expressed in natural language text (Darwich, Mohammad & Mohd Noah, Shahrul Azman & Omar, Nazlia & Osman, Nurul, 2019). It plays a crucial role in various fields such as marketing, customer feedback analysis, social media monitoring, and political analysis (Yasen, Mais & Tedmori, Sara., 2019).

Summarily, it involves the automatic extraction and classification of sentiments from textual data, categorizing them as positive, negative, or neutral (Cyril, C. P. D., Beulah, J. R., Subramani, N., Mohan, P., Harshavardhan, A., & Sivabalaselvamani, D., 2021). In essence, it utilizes natural language processing (NLP) and machine learning techniques to analyses the subjective information present in the text.

Sentiment classification is composed of machine learning approaches, lexicon-based approaches, and hybrid-based approaches (Kandukuri & Gopal, 2019). The lexicon-based approaches

include the corpus-based or bag of words (BOW) and dictionary also called opinion lexicon (Umar, Ahmad, & Zainal, 2020).

The significance of sentiment analysis in diverse fields will be underscored, elucidating its role in extracting valuable insights from textual data, understanding public sentiment, and facilitating informed decision-making. This paper provides an overview of the corpus-based approach for sentiment analysis, focusing on its methodologies, evaluation, limitations, and future directions.

### 1.1. Role of Sentiment Analysis in Various Fields

Sentiment analysis, also known as opinion mining, is a valuable tool in data analysis across various domains. Its application extends to a wide range of fields, including marketing, customer service, social media monitoring, financial analysis, and political science (Ms. Binju Saju; Ms. Siji Jose; Mr. Amal Antony, 2020). Here are some specific applications of sentiment analysis in data analysis as highlighted by (Priya Jadhav; Aditya Saha, 2023):

- i **Customer Feedback Analysis:** Sentiment analysis is widely used to analyze customer feedback, such as product reviews, surveys, and social media comments. By automatically categorizing sentiments expressed in customer feedback, businesses can identify areas for improvement, monitor brand perception, and tailor their products or services to meet customer needs effectively.
- ii **Market Research:** In market research, sentiment analysis helps companies gauge public opinion about their products, services, or brand reputation. By analyzing sentiments expressed in online discussions, forums, and social media platforms, businesses can identify emerging trends, assess market sentiment, and make data-driven decisions regarding product development, marketing strategies, and customer engagement.
- iii **Brand Monitoring and Reputation Management:** Sentiment analysis enables businesses to monitor their brand's online reputation by analyzing sentiments expressed in online mentions, news articles, and social media conversations. By tracking positive and negative sentiments associated with their brand, businesses can address customer concerns promptly, mitigate potential PR crises, and maintain a positive brand image.
- iv **Customer Service and Support:** Sentiment analysis is utilized in customer service to analyze customer inquiries, complaints, and feedback received through various channels, such as emails, live chats, and social media messages. By automatically categorizing the sentiment of customer interactions, businesses can prioritize and address issues more efficiently, improve response times, and enhance overall customer satisfaction.
- v **Financial Analysis:** In finance, sentiment analysis is employed to analyze sentiments expressed in financial news, social media discussions, and analyst reports to gauge market sentiment and predict stock price movements. By monitoring investor sentiment and market sentiment indicators, financial analysts and traders can make more informed investment decisions and identify potential market trends or anomalies.

- vi **Political Analysis:** Sentiment analysis is used in political analysis to analyze public sentiment towards political candidates, parties, policies, and current events. By analyzing sentiments expressed in social media discussions, news articles, and public opinion polls, political analysts can assess voter sentiment, track electoral trends, and predict election outcomes.
- vii **Product and Service Evaluation:** Sentiment analysis is utilized to evaluate the performance of products and services by analyzing sentiments expressed in customer reviews, ratings, and feedback. By aggregating and analyzing sentiment data, businesses can identify the strengths and weaknesses of their offerings, benchmark against competitors, and prioritize areas for improvement.
- viii **Event Monitoring and Crisis Management:** Sentiment analysis is employed to monitor public sentiment during events, crises, or public relations campaigns. By analyzing sentiments expressed in real-time social media conversations, news articles, and online discussions, organizations can assess public perception, identify potential issues or crises, and take proactive measures to address them.

## 2. Review of Some Related Works

### 2.1 Overview of Corpus-Based Approach for Sentiment Analysis

The corpus-based approach relies on large collections of annotated text data, known as corpora, to train machine learning models for sentiment analysis (Darwich et al., 2019). It involves several key steps, including data collection, preprocessing, feature extraction, and sentiment classification.

### 2.2. Data Collection and Preprocessing

Data collection involves gathering text data from various sources such as social media, product reviews, and news articles. Preprocessing steps include tokenization, lowercasing, removal of stop words, punctuation, and special characters, as well as stemming or lemmatization to normalize the text (Ms. Binju Saju, Ms. Siji Jose, and Mr. Amal Antony, 2020).

### 2.3. Feature Extraction and Selection

Feature extraction involves transforming the preprocessed text data into numerical representations, such as bag-of-words, TF-IDF (Term Frequency-Inverse Document Frequency), or word embeddings (Binju et al., 2020). Feature selection techniques may be employed to reduce dimensionality and improve model performance (Yasen, Mais & Tedmori, Sara., 2019).

### 2.4. Sentiment Classification Techniques

Several machine learning algorithms, including but not limited to, Naive Bayes, Support Vector Machines (SVM) (Cyril et al., 2021), Logistic Regression, and Neural Networks (Haisal Dauda Abubakar; Sharin Hazlin Huspi; Mahmood Umar, 2022), are commonly used for sentiment classification. These algorithms learn to classify text data into positive, negative, or neutral sentiment categories based on the extracted features.

### 3. Evaluation of Corpus-Based Approach

#### 3.1. Performance Metrics for Sentiment Analysis

Performance metrics such as accuracy, precision, recall, F1-measure, confusion matrix, **k-fold**, and cross-validation are used to evaluate the effectiveness of sentiment analysis models (RAGHUNATHAN & KANDASAMY, 2023). These metrics measure the model's ability to correctly classify sentiments and handle class imbalances.

#### 3.2. Comparison with Other Approaches

The corpus-based approach is compared with machine learning approaches, dictionary-based approaches, rule-based approaches, and hybrid methods to assess their relative performance in sentiment analysis tasks. Corpus-based algorithms often outperform simpler approaches, especially in scenarios with complex or domain-specific language (Darwich, Mohammad & Mohd Noah, Shahrul Azman & Omar, Nazlia & Osman, Nurul, 2019).

In sentiment analysis, various approaches are employed to classify the sentiment expressed in textual data. The comparison between the lexicon-based approaches (corpus-based approach in particular), machine learning approaches, rule-based systems, and hybrid methods allows us to assess their relative performance in sentiment analysis tasks, **are shown in Table:1.**

**Table 1: Observation the comparison between Lexicon-based approaches and Others**

Approach	Methodology	Strengths	Weaknesses
Machine Learning Approach	<ul style="list-style-type: none"><li>Relies on large collections of annotated text data (corpora) to train ML models for sentiment analysis (Darwich et al., 2019).</li></ul>	<ul style="list-style-type: none"><li>Can handle complex language patterns and adapt to different domains (Darwich et al., 2019).</li><li>Effective in capturing context and nuances in sentiment expression.</li></ul>	<ul style="list-style-type: none"><li>Requires large amounts of data for training.</li><li>May suffer from data sparsity issues in specific domains (Haisal Dauda Abubakar; Sharin Hazlin Huspi; Mahmood Umar, 2022).</li><li>Performance heavily depends on the quality of the training corpus.</li></ul>
Lexicon-Based Approach	<ul style="list-style-type: none"><li>Utilizes sentiment lexicons or dictionaries containing predefined lists of words associated with positive or negative sentiments (Umar, Aliyu, &amp; Modi, 2022).</li><li>It's classified as corpus-based and dictionary-based approaches by (RAGHUNATHAN &amp; KANDASAMY, 2023).</li><li>Corpus-based approaches employ two</li></ul>	<ul style="list-style-type: none"><li>Simple and interpretable.</li><li>Can handle out-of-vocabulary words.</li><li>Less computationally intensive compared to machine learning approach.</li></ul>	<ul style="list-style-type: none"><li>Limited coverage of sentiment lexicons may lead to inaccuracies, especially with domain (Darwich et al., 2019).</li><li>Specific or ambiguous terms.</li><li>May struggle with context-dependent sentiment expressions.</li></ul>

	techniques: Semantic and Statistical (RAGHUNATHAN & KANDASAMY, 2023)		
Rule-Based Systems	<ul style="list-style-type: none"> <li>• Defines explicit rules or heuristics based on linguistic patterns, syntactic structures, or semantic rules.</li> <li>• It does sentiment analysis on the basis a set of human-created rules to identify subject, polarity, or the opinion (Binju et al., 2020).</li> </ul>	<ul style="list-style-type: none"> <li>• Allows for fine-grained control over sentiment analysis process.</li> <li>• Can incorporate domain-specific knowledge and linguistic rules.</li> </ul>	<ul style="list-style-type: none"> <li>• Limited scalability and generalization.</li> <li>• Requires manual effort to design and maintain rules.</li> <li>• May struggle with capturing nuances and variations in sentiment expression.</li> <li>• not very efficient as it will not analyse how words are combined in a sequence (Binju et al., 2020).</li> </ul>
Hybrid Methods	<ul style="list-style-type: none"> <li>• Combines multiple approaches (machine learning, lexicon-based, rule-based) to leverage their strengths and mitigate their weaknesses (Binju et al., 2020).</li> </ul>	<ul style="list-style-type: none"> <li>• Offers improved robustness and performance by combining complementary techniques.</li> <li>• Can adapt to diverse datasets and domains.</li> </ul>	<ul style="list-style-type: none"> <li>• Increased complexity in implementation and tuning.</li> <li>• Requires careful integration and coordination of different components.</li> </ul>

Each approach in Table 1 has its own strengths and weaknesses in sentiment analysis tasks. Corpus-based approach excels in capturing context and nuances but requires substantial labeled data. Lexicon-based approaches are simple and interpretable but may lack coverage and struggle with context. Rule-based systems offer fine-grained control but are limited in scalability and generalization. Hybrid methods aim to leverage the strengths of multiple approaches but require careful integration and tuning. The choice of approach depends on factors such as the availability of labeled data, the complexity of the sentiment analysis task, and the desired balance between accuracy and interpretability.

### 3.3. Case Studies and Experimental Results

In this section, case studies and experimental results on corpus-based approaches in real-world sentiment analysis tasks will be demonstrated and summarized in tabular form. The tables showcase the performance of lexicon-based approaches especially the corpus-based approaches across different domains and datasets, highlighting their versatility and scalability.

A review of studies (in Table 2.) shows that R and Python programming are the commonly used tools for sentiment analysis. The review also shows that scholars design their own dictionaries for sentiment classification. The following hint for an easy understanding of the abbreviations are shown in Table 2.

Hint: (Text Sentiment Score, or TSS) and Text Sentiment Intensity (TSI)

**Table 2: Observation the experimental results on Lexicon-based approaches for sentiment analysis**

Studies	Level	Feature Extraction	Source of Data	Algorithm	Tool	Evaluation	
						Metrics	Results
(Kandukuri & Gopal, 2019)	Aspect level	Tf-IDF	Email	Lexicon-based	R-Tidy text package	Positive	>100%
						Negative	-50%
(Vargas-Sierra & Orts, 2023)	Sentence level	Polarity and intensity	Newspapers; English (Economist) and Spanish(Expansion)	Lexicon-based	Lingmotif 2 software	<b>IN COVID</b>	
						TSI-English	76%
						TSI-Spanish	81%
						TSS-English	40%
						TSS-Spanish	49%
(Alves & Bekavac, 2023)	Sentence Level	lang2vec a	News and Wikipedia	corpus-based quantitative methods: lang2vec clustering and Marsagram linear properties clustering	NLP tools:lang2vec andMarsagram tool	Language Cluster comparison:	N/A
(Munnes, Harsch, Knobloch, & Vogel, 2022)	Document Level	Word count,	Dictionaries: SentiWS, Rauh's German Political Sentiment Dictionary, GerVADER, GloVeDictionary	Dictionary-based approaches:SentiWS, Rauh's German Political Sentiment Dictionary, GerVADER, GloVe algorithm,wordscores and wordfish	R package	<b>Maximal Correlations</b>	
						SentiWS	0.29
						Rauh	0.37
						GerVADER	0.32
(Heidarypur, Pahlavannezhad, & Kahani, 2023)	Aspect level	TFIDF	News data from Website	Sentiment Dictionary approaches;Keyword baseline (M1) and PMI (M2)	Python Classifiers	<b>Accuracy</b>	
						Keyword baseline (M1)	37.01%
						PMI (M2)	42.01%

#### 4. Limitations and Future Directions

#### 4.1. Challenges in Corpus-Based Sentiment Analysis

According to Siyu Lei and Chu-Ren Huang, (2022) challenges of sentiment analysis include the difficulty of dealing with non-standard language and semantic ambiguity. Technical problems are also posing a significant challenge to sentiment analysis, such as the lack of training data that can be used across domains and, more in general, the lack of datasets for specialized domains.

Challenges in corpus-based sentiment analysis include data sparsity, domain adaptation, context, and sensitivity which need a large annotated corpora (RAGHUNATHAN & KANDASAMY, 2023). Addressing these challenges requires advanced techniques in machine learning, NLP, and domain-specific knowledge incorporation.

#### 4.2. Recommendations and Future Works

Future research directions include exploring ensemble learning techniques, deep learning architectures, incorporating multimodal data, merging lexical, and machine learning model for improved sentiment analysis (Islam, Sheakh, Sadik, Mst, & Tahosin, 2024). Additionally, efforts towards developing domain-specific sentiment lexicons and annotated corpora can further enhance the performance of corpus-based algorithms (Alves & Bekavac, 2023).

#### Conclusion

Corpus-based approach for sentiment analysis offers a powerful and versatile approach to extracting sentiment from textual data. By leveraging large annotated corpora and machine learning techniques, these algorithms can provide valuable insights into public opinion, customer sentiment, and market trends. Despite their limitations, there are ongoing research efforts that enhance the performance and applicability of corpus-based approaches for sentiment analysis in diverse domains.

#### Reference

- ALNAWAS, A., & ARICI, N. (2018). The Corpus Based Approach to Sentiment Analysis in Modern Standard Arabic and Arabic Dialects: A Literature Review. *Journal of Polytechnic*, 461-470.
- Alves, D., & Bekavac, B. (2023). Analysis of Corpus-based Word-Order Typological Methods. *Proceedings of the Sixth Workshop on Universal Dependencies (UDW, GURT/SyntaxFest 2023)* (pp. 36-46). .. Association for Computational Linguistics.
- Cyril, C. P. D., Beulah, J. R., Subramani, N., Mohan, P., Harshavardhan, A., & Sivabalaselvamani, D. (2021). An automated learning model for sentiment analysis and data classification of Twitter data using balanced CA-SVM. *Concurrent Engineering*, 29(4), 1-10.
- Darwich, Mohammad & Mohd Noah, Shahrul Azman & Omar, Nazlia & Osman, Nurul. (2019). Corpus-Based Techniques for Sentiment Lexicon Generation: A Review. *Journal of Digital Information Management*, 17, 296-305.

- Haisal Dauda Abubakar; Sharin Hazlin Huspi; Mahmood Umar. (2022). A Scheme of Pairwise Feature Combinations to Improve Sentiment Classification Using Book Review Dataset. *International Journal of Innovative Computing*, 25-33.
- Heidarypur, M., Pahlavannezhad, M. R., & Kahani, M. (2023). The role of corpus linguistics in sentiment analysis of Persian texts, case study: a Farsi news agency website. *Journal of linguistics, philology and translation*, 106-121.
- Islam, T., Sheakh, M. A., Sadik, M. R., Mst, & Tahosin, S. (2024). Lexicon and Deep Learning-Based Approaches in Sentiment Analysis on Short Texts. *Journal of Computer and Communications*, 11-34.
- Kandukuri, M., & Gopal, V. H. (2019). Sentiment Analysis on Email Database Corpus-based Approach. *Singapore Journal of Scientific Research*, 9 (2), 45-51.
- Ms. Binju Saju; Ms. Siji Jose; Mr. Amal Antony. (2020). Comprehensive Study on Sentiment Analysis: Types, Approaches, Recent.
- Munnes, S., Harsch, C., Knobloch, M., & Vogel, J. S. (2022). Examining Sentiment in Complex texts: a Comparison of different Computational Approaches. *Frontiers of Big data*, 1-16.
- Priya Jadhav; Aditya Saha. (2023). Sentiment Analysis of Mobile App Reviews Using Robotic Process automation. *2023 7th International Conference on Electronics, Materials Engineering & Nano-Technology (IEMENTech)* (pp. 1-6). .: IEEE Xplore.
- RAGHUNATHAN, N., & KANDASAMY, S. (2023). Challenges and Issues in Sentiment Analysis: A Comprehensive Survey. *IEEE ACCESS*, 69626-69642.
- Siyu Lei<sup>1</sup>; Chu-Ren Huang. (2022). Conducting Sentiment Analysis: Lei L. & Liu D. Elements in Corpus Linguistics, CUP. *Springer*, 29(1), .
- Umar, M., Ahmad, N. B., & Zainal, A. (2020). Sentiment Analysis of Student's Opinion on Programming Assessment: Evaluation of Naïve Bayes over Support Vector Machine. *International Journal of Innovative Computing*, 51-58.
- Umar, M., Aliyu, M., & Modi, a. S. (2022). Sentiment Analysis in the Era of Web 2.0: Applications, Implementation Tools and Approaches for the Novice Researcher. *Caliphate Journal of Science & Technology (CaJoST)*, 1-9.
- Vargas-Sierra, C., & Orts, M. Á. (2023). Sentiment and emotion in financial journalism: a corpus-based, cross-linguistic analysis of the effects of COVID 19. *HUMANITIES AND SOCIAL SCIENCES COMMUNICATIONS*, 1-17.
- Yasen, Mais & Tedmori, Sara. (2019). Movies Reviews Sentiment Analysis and Classification. *2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT)*. Amman, Jordan: IEEE Xplore.