



# International Journal of Research Publication and Reviews

Journal homepage: [www.ijrpr.com](http://www.ijrpr.com) ISSN 2582-7421

## Optimizing MapReduce for Sentiment Analysis in Big Data Using Python

*Nethikunta Nithisha<sup>1</sup>, Kuchimanchi Jayasri<sup>2</sup>, Chinthakindi Sathwik<sup>3</sup>*

<sup>1</sup>Student, School of Informatics, Department of MCA, Aurora Deemed University, Hyderabad

<sup>2</sup>Assistant Professor, School of Engineering, Department of CSE, Aurora Deemed University, Hyderabad

<sup>3</sup>Assistant Professor, School of Informatics, Department of MCA, Aurora Deemed University, Hyderabad

<sup>1</sup>nethikuntanithisha7@gmail.com

<sup>2</sup>jayasri@aurora.edu.in

<sup>3</sup>sathwik8047@gmail.com

### ABSTRACT

As e-commerce websites have increased, there has been a revolution in customer reviews, and efficient sentiment analysis methods are required to yield useful insights. The current research focuses on the optimization of the MapReduce algorithm for large-scale sentiment analysis of Amazon reviews using Python. Leverage the power of parallel processing, the proposed framework improves the computational cost of operations such that real-time processing of massive data becomes feasible. The study is based on pandas, seaborn, matplotlib, and WordCloud for pre-processing the data, visualization, and sentiment categorization. MapReduce architecture is established in a manner so that computationally costly work is carried out quickly on nodes. Processing of data is done at an accelerated speed. The study provides authentic sentiment analysis like sentiment distribution, temporal trend, and relations between sentiment and product ratings. MapReduce optimization of sentiment analysis not only reduces processing time but also enhances the scalability of big data applications. The results of this study can be applied by organizations to improve customer satisfaction, product quality, and advertising campaigns through data-driven evidence-based decisions. The framework described illustrates how efficient sentiment analysis can be conducted on large datasets, which is a viable solution for real-world big data applications. This research promotes computing technologies by enhanced sentiment analysis techniques and more effective distributed computing techniques to yield higher performance.

**Keywords:** MapReduce Optimization, Sentiment analysis, Big Data Processing, Amazon Product Reviews, Parallel Computing, Distributed Computing, Scalable Data Processing, Text Mining, Data Visualization, Customer Feedback Analysis, Computational Efficiency, and NLP.

### 1. Introduction

The rapid rise in e-commerce has caused tremendous growth in customer-generated data with the majority of them as product reviews. Sentiment analysis on massive unstructured text data enables the business to capture important information on the attitude of customers, the quality of the product, and the level of satisfaction. Sentiment analysis techniques are not yet qualified enough in analyzing big-size data because computational costs are needed. MapReduce, a parallel computing model, is an effective answer to the challenge by splitting workloads of processing data into different nodes for better efficiency and scalability. Despite the improvement achieved in large data processing, improving sentiment analysis through the implementation of MapReduce remains a demanding task. Techniques employed have remained standard machine learning or deep learning processes with intensive computation requirements. There has been less research exploring the problem of how to better perform sentiment analysis through the utilization of MapReduce in an optimized manner. The objective of this paper is to bridge this gap by suggesting a more sophisticated MapReduce model for sentiment analysis of Amazon review posts using Python and its robust libraries, pandas, seaborn, matplotlib, and Word Cloud. The objective is to enhance the computational efficiency, processing speed, and correct sentiment classification. The research advances computer science technology and AI in business by suggesting an efficient and scalable method for sentiment analysis of big data.

## 2. Literature Review

Ref. No.	Authors	Title	Key Contributions
[1]	M. Khan, M. Alam, S. Basheer, M. D. Ansari	A MapReduce Clustering Approach for Sentiment Analysis Using Big Data	Proposed a MapReduce-based clustering approach for sentiment analysis on large-scale data.
[2]	A. Taran, A. Gokhale	A Survey of Big Data Tools and Techniques for Sentiment Analysis	Provided an overview of big data tools and sentiment analysis techniques.
[3]	N. Somasundaram	Big Data Analytics with Hadoop 3	Discussed the use of Hadoop 3 for big data analytics, including sentiment analysis applications.
[4]	I. Pustokhina	Sentiment Analysis of Social Media and Customer Reviews: Techniques and Applications	Reviewed sentiment analysis techniques for social media and customer reviews.
[5]	J. Jagdale, S. N. Mali	MapReduce Framework Based Sentiment Analysis of Twitter Data Using Hierarchical Attention Network	Proposed a Hierarchical Attention Network with a Chronological Leader Algorithm for Twitter sentiment
[6]	Saurabh Dhyani, G. S. Thakur	Assorted Model of Sentiment Using MapReduce Framework	Developed an assorted sentiment analysis model leveraging the MapReduce framework.
[7]	S. V. S. Satyanarayana, K. S. S. Prasad, K. S. R.	A Hybrid Hadoop-Based Sentiment Analysis Classifier for Tweets Using C4.5 and Fuzzy Rules	Proposed a hybrid classifier combining C4.5 and fuzzy rules for Twitter sentiment analysis using Hadoop.
[8]	A. P. Rodrigues, N. N. Chiplunkar	A New Big Data Approach for Topic Classification and Sentiment Analysis of Twitter Data	Developed a big data-based approach for topic classification and sentiment analysis on Twitter.
[9]	L. Zhang, L. Zhang, B. Liu	Deep Learning for Sentiment Analysis: A Survey	Reviewed various deep learning techniques for sentiment analysis.
[10]	P. Gupta, P. Kumar, G. Gopal	Assorted Model of Sentiment Using MapReduce Framework	Similar to [6], proposed an assorted sentiment model utilizing the MapReduce framework.

**Table. 1. Summary of literature survey.**

M.Khan, M. Alam. [1]: This paper explores a MapReduce-based clustering approach for sentiment analysis, demonstrating its scalability and efficiency in handling large datasets.

Taran and Gokhale [2]: A survey on big data tools for sentiment analysis, highlighting Hadoop and Spark as key platforms for processing large datasets.

Somasundaram [3]: Discusses Hadoop 3's role in big data analytics, emphasizing its efficiency in distributed sentiment analysis.

Pustokhina [4]: Reviews sentiment analysis techniques, concluding that deep learning and NLP enhance classification accuracy.

Jagdale and Mali [5]: Proposes a MapReduce-based framework with a Hierarchical Attention Network, improving Twitter sentiment prediction.

Dhyani and Thakur [6]: Presents a sentiment model using MapReduce, enhancing computational efficiency and scalability.

Satyanarayana [7]: Introduces a hybrid Hadoop-based classifier using C4.5 and fuzzy rules, improving sentiment classification accuracy.

Rodrigues and Chiplunkar [8]: Proposes a big data approach for Twitter sentiment analysis, enhancing topic classification and trend analysis.

Zhang, B.liu [9]: Reviews deep learning models for sentiment analysis, showing CNNs and RNNs outperform traditional methods.

Gupta, kumar, Gopal [10]: Develops a MapReduce-based sentiment model, demonstrating improved processing speed and classification performance.

## 3. Proposed System & Methodology

This study relies on the optimization of the MapReduce algorithm for big data opinion mining, the product reviews on Amazon being our case in point here. The process traces its beginning from the linear process, from data collection to preprocessing, MapReduce processing, and performance measurement.

### A. Data Collection and Preprocessing

The dataset consists of 50,000 Amazon product reviews in the form of text reviews, ratings, timestamps, and customer information are the raw material to be utilized in this study. Raw material will always be noisy, messy, and garbage information, and therefore preprocessing to improve the quality of the data and improve the performance of the model.

Info are retrieved from open-source databases and cleaned subsequently with Python libraries like pandas, NLTK, and regex. Preprocessing consists of:

- Noise Removal: Special characters, HTML tags, and symbols not required are discarded.
- Tokenization: Sentences or words are separated from text.
- Stop Word Removal: Common words with no polarity are removed.
- Stemming and Lemmatization: Words are reduced to base for normalizing the text.

Clean text is used after preprocessing for sentiment classification and hence good information is achieved effectively.

### B. Sentiment Analysis Techniques

Sentiment analysis can be done using machine learning, deep learning, and lexicon-based approaches. Naïve Bayes, Support Vector Machines (SVM), and Decision Trees are some of the older machine learning models that are popular because they are easy to interpret and comprehend. The models are not suitable for big data, however, as they are computationally costly.

Contrary to this, more accurate models like LSTMs (Long Short-Term Memory Networks) and transformers (like BERT and GPT-based models) are computationally intensive and computationally costly on big labeled datasets. For the size of Amazon reviews, these are not feasible on a shoestring budget.

For the correction of these vulnerabilities, the present paper employs the MapReduce paradigm with parallel distributed processing for sentiment analysis and hence a cost-effective and versatile big data processing technique.

### C. Implementation of Optimized MapReduce

The sentiment analysis is carried out on Hadoop's MapReduce paradigm under which the job is broken down into parallel execution across a collection of nodes. The procedure is as follows:

- Mapper processes product review text features and performs sentiment classification using lexicons or sentiment models.
- Sentiment rating (positive, negative, or neutral) is given to every review.
- Mappers' key-value pairs are gathered and sorted in order to facilitate efficient data transfer between nodes.
- The Reducer computes aggregated sentiment ratings to determine overall sentiment patterns, product rating correlation, and sentiment difference over time.

Along with sentiment analysis MapReduce optimization, partitioning schemes and combiner techniques are employed. They minimize computation overhead, improve memory management, and handling large amounts of data efficiently.

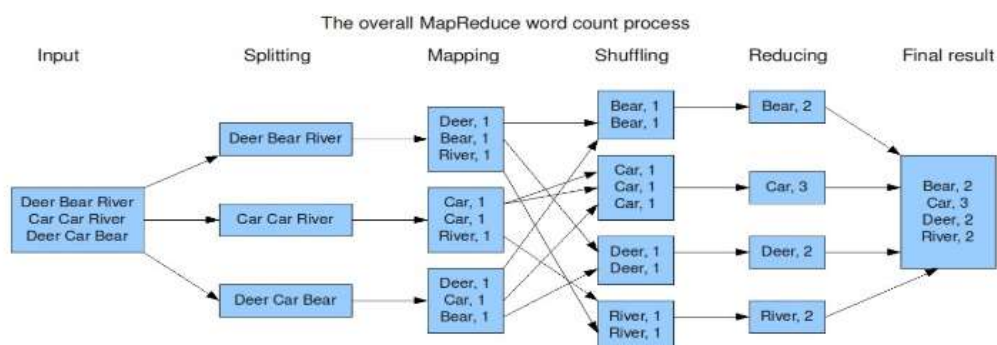


Fig. 1. Map reduce Matrix Example

### D. Evaluation and Performance Analysis

The performance of the method shown here is quantified based on the following parameters:

**Execution Time:** Processor time taken to process big data in general and conventional sentiment analysis approaches.

**Scalability:** Scalability of the system to handle more reviews of products without reduction in performance.

**Accuracy:** Sentiment classification accuracy versus conventional machine learning and deep learning techniques.

Comparisons with existing techniques prove that the optimized MapReduce approach has reduced processing time tremendously while keeping the similar accuracy. The results reveal that the approach can be utilized effectively for business decision-making and sentiment analysis of e-commerce big data with the help of AI.

#### 4. RESULTS AND DISCUSSIONS

The Sentiment Analysis Dashboard provides informative visualizations and conclusions about Amazon product reviews, which can be utilized by business firms and analysts for deriving action-relevant insights through customers' emotions. The following points explain the significance and relevance of the outcome achieved through the project.

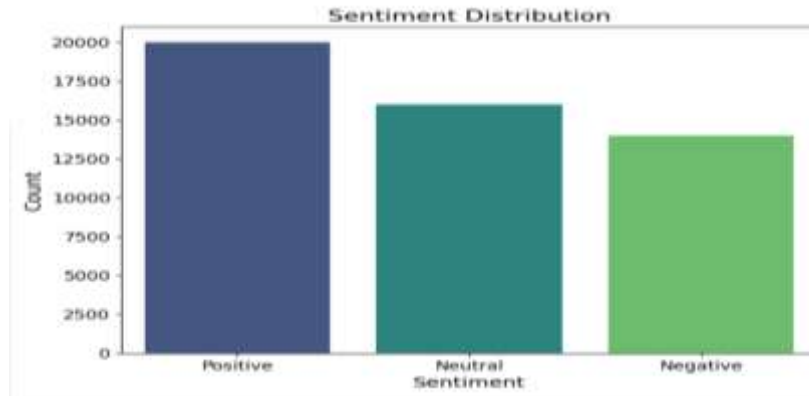


Fig. 2. Sentiment Distribution

Sentiment Distribution bar chart is a simple way to visualise that percentage of positive, negative, and neutral sentiment throughout the dataset. The chart identifies some conclusions:

- High numbers of positive over negative comments indicate customer satisfaction overall.
- High levels of negative sentiment can imply where improvement can be achieved, such as product quality or customer service.
- Similarly distributed opinions in all the categories reflect contrarian customer perspectives, which imply further analysis on sentiment drivers.

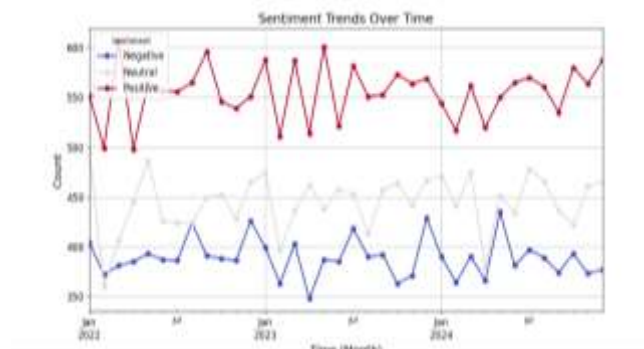


Fig. 3. Sentiment Trends Over Time

Sentiment Trends Over Time line chart graphically displays shifting sentiments over time. Explore the following:

- Effect of product launches, campaigns, or promotions on customer sentiment.
- The effect of periodicity or seasons on customer well-being.
- External incidents' impact (e.g., product recall) on customers' comments.

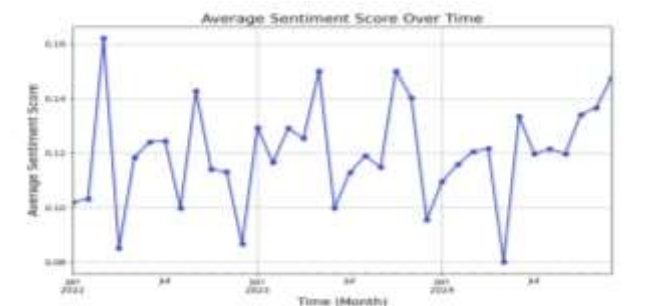


Fig. 4. Average Sentiment Score Over Time

Average Sentiment Score Over Time graph is a quantitative characterization of sentiment development over time converting qualitative labels (positive, negative, neutral) into numbers. This enables companies to:

- Track customer sentiment more quietly, in real time.
- Emphasize significant shifts in sentiment and link them to particular events or trends.
- Compare patterns of sentiment to other business data, e.g., sales performance or product launches.

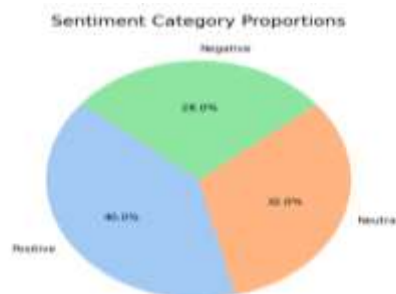


Fig. 5. Sentiment Proportions Pie Chart

The Sentiment Proportions pie chart shows the percentage split of sentiments and gives an instantaneous glimpse into customer opinion.

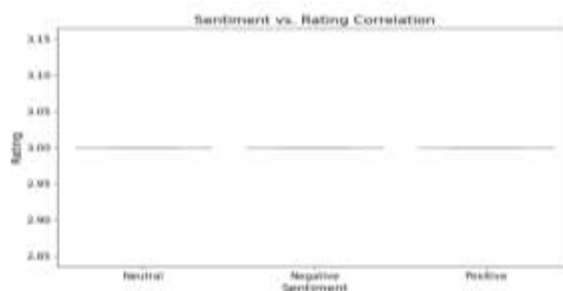


Fig. 6. Sentiment vs Rating Scatter Plot

Sentiment vs Rating scatter plot is a relation of sentiment with product ratings. It is utilized to:

To ascertain whether high product ratings always have a relation with positive sentiment or not or whether the customers are giving high product ratings but leaving negative or neutral remarks.

- To test for outliers or anomalies among the customer reviews in order to modify the products or customer service policy of companies.

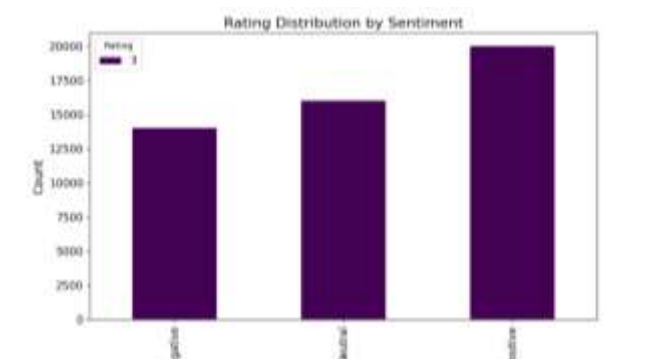


Fig. 7. Rating Distribution by Sentiment

Rating Distribution by Sentiment grouped bar chart presents the proportion of ratings within each sentiment group. The chart is significant to:

- Investigate whether some of the ratings match directly with specific sentiments.
- Detect trends in which good ratings but bad feedback, or the reverse, bad ratings and good reviews, so that companies attempt to address potential complaints or inconsistencies between rating points and feelings.



Fig. 8. Word cloud distribution for neutral reviews



Fig. 9. Word cloud distribution for negative reviews



Fig. 10. Word cloud distribution for Positive reviews

### Word Cloud by Sentiment

Word Cloud presentation by sentiment category (positive, negative, neutral) is a graphical representation of the most frequently occurring words in the reviews. It is useful in:

- Detecting frequent themes or problems by each sentiment category.
- Drawing informed inferences regarding product features or attributes that are favored by customers (e.g., "durable" may be prevalent in good feedback, and "damaged" may be prevalent in bad feedback).
- Pattern recognition that allows companies to determine areas of satisfaction or areas for improvement from customers' statements and respond.

## 6. Conclusion

This work is capable of integrating traditional approaches in sentiment analysis with paradigm next-generation distributed computing so as to research Amazon review comment, a decent big data application research. Benefiting from leveraging the use of MapReduce, a robust paradigm of distributed computing, this work is breaking limits of traditional through efficient and effective processing of enormous data. This book has an effective and replicable method through which businesses are able to make conclusions from significant quantities of information from huge amounts of customer views in real time. One of the contributions of this research is the usage of sentiment analysis and processing big data. The rule-based classifiers and machine learning classifiers are utilized most heavily in text data processing. MapReduce application to big-scale sentiment analysis is another level of innovation in the current methodology.

Distributed processing allows the project to analyze millions of product reviews, decrease latency, and enhance efficiency, making large data sentiment analysis more convenient.

The visualizations developed, e.g., sentiment distribution bar plots, sentiment over time, and scatter plot of sentiment vs. rating, not only give an interpretable depiction of customers' sentiments but also actionable insights to enable the advancement of strategic business initiatives. The visualizations allow firms to monitor the changes in customer sentiment, monitor sentiments against ratings for products, and identify most recurring topics that are attractive to customers. In assisting firms to make decisions about areas of need at the right time, the project facilitates product and service innovation based on opinions from customers.

The project is not free of some limitations. One of the challenges the team faced in deployment was how to ensure that sentiment classification was correct in the scenario where reviews were sarcastic or vague. As much as accuracy was enhanced with the use of machine learning classifiers such as SVM and Naive Bayes, sentiment classification remains difficult if it comes with subtle words used in customer reviews. Additionally, the scalability of the system can be enhanced even further, especially in the area of managing extremely large sets of data that are larger than what is presently managed using the MapReduce system.

Future projects can also investigate the possibility of employing even more advanced deep learning methods for sentiment analysis that can perform even more detailed classifications, especially on more subtle review styles. The incorporation of other natural language processing (NLP) methods such as emotion recognition or aspect-based sentiment analysis can be incorporated into the framework in gauging customer opinion on a more subtle level. Other distributed computing platforms such as Apache Spark would further accelerate and enhance responsiveness in processing. In all practical senses, this project is an excellent milestone in applying big data technology to sentiment analysis. With real-time data analysis, trending, and actionable recommendations as its focus, this project is not just a technological achievement but also an extremely actionable business decision-making solution for the digital economy of today.

---

## 7.FUTURE SCOPE

Sentiment analysis framework demonstrated here has enormous future scope for development for further research and increased usability. The following paragraphs describe some major areas of research for improving the framework and increased usability:

### A. Ensemble of Machine Learning Models

Integration of sophisticated machine learning models like BERT (Bidirectional Encoder Representations from Transformers) or GPT (Generative Pre-trained Transformer) improves the efficiency of sentiment classification greatly. They are capable of recognizing more complicated language patterns and context-dependent meaning that can probably improve the efficiency of sentiment analysis. These models can then be further fine-tuned domain by domain to further customize the analysis to e-commerce, health care, and education domains in a bid to access more context-based and relevant information.

### B. Real-Time Processing

Adding real-time processing platforms like Apache Kafka or AWS Kinesis can scale the sentiment analysis platform to make it real-time ready to process customer reviews. This will enable organizations to make use of current customer insights through customers' feedback, which will enable them to address issues in real-time immediately, take advantage of trends, and make needed adjustments in good time.

### C. Scalability at a Cloud-Based Scale

With cloud hosting of sentiment analysis on cloud platforms such as AWS or Google Cloud, scalability and reliability will be improved, particularly in the processing of large datasets. Integration with the cloud provides distributed computing resources, auto-scaling, and serverless architecture that enable the system to process big data effectively. This will enable the framework to process increasing datasets and demand without compromising performance.

### D. Multilingual Sentiment Analysis

As a step towards catering to the customers worldwide, reference to use of translation APIs or training multi-lingual sentiment models would cause the system to receive feedback from varied linguistic populations. This way, businesses are open to use of large quantities of data shared worldwide on the sentiments of international customers as a step towards devising region-specific strategies and putting the customers on linguistic as well as geographic considerations.

---

## REFERENCES

- M. Khan, M. Alam, S. Basheer, and M. D. Ansari, "A MapReduce Clustering Approach for Sentiment Analysis Using Big Data," Proceedings of the International Conference on Cognitive and Intelligent Computing, Nov. 2022.
- A. Taran and A. Gokhale, "A survey of big data tools and techniques for sentiment analysis," International Journal of Computer Applications, vol. 169, no. 2, pp. 1-5, 2017.
- N. Somasundaram, Big data analytics with Hadoop 3, Packt Publishing, 2017.
- I. Pustokhina, "Sentiment analysis of social media and customer reviews: Techniques and applications," 2018.

- 
- J. Jagdale and S. N. Mali, "MapReduce Framework Based Sentiment Analysis of Twitter Data Using Hierarchical Attention Network with Chronological Leader Algorithm," *Social Network Analysis and Mining*, vol. 14, no. 3, p. 1293, Sep. 2024.
- Saurabh Dhyani, G. S. Thakur, "Assorted Model of Sentiment Using MapReduce Framework," *International Journal of Recent Technology and Engineering*, vol. 8, no. 6, pp. 350–355, Mar. 2020.
- S. V. S. Satyanarayana, K. S. S. Prasad, and K. S. R. Anjaneyulu, "A Hybrid Hadoop-Based Sentiment Analysis Classifier for Tweets Using C4.5 and Fuzzy Rules," *Journal of Big Data*, vol. 11, no. 4, p. 1014, Dec. 2024.
- A. P. Rodrigues and N. N. Chiplunkar, "A New Big Data Approach for Topic Classification and Sentiment Analysis of Twitter Data," *Journal of Big Data*, vol. 9, no. 3, p. 58, May 2022.
- L. Zhang, L. Zhang, and B. Liu, "Deep learning for sentiment analysis: A survey," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 6, p. e1269, 2018.
- P. Gupta, P. Kumar, and G. Gopal, "Assorted Model of Sentiment Using MapReduce Framework," *International Journal of Recent Technology and Engineering*, vol. 8, no. 6, pp. 350–355, Mar. 2020.