

An Efficient Data Preprocessing Techniques for Sentiment Analysis Using MOOC Dataset in Machine Learning

¹S. Daisy Fatima Mary and ²Dr. G. Mageswary,

¹Guest Lecturer, PG Department of Computer Science, Government Arts and Science College, Thiruvannainallur, India

²Assistant Professor, Department of Computer Science, PSPT MGR Government Arts and Science College, Sirkali, Puthur, India

Abstract: Sentiment analysis is a significant field of study and a highly sought-after discipline that focuses on discerning the feelings, viewpoints, and emotions conveyed within a piece of text. Sentiment Analysis enables the extraction of valuable insights from textual data sourced from platforms such as Facebook, Twitter, Amazon, and more. Social media data are frequently unstructured and challenging to manage due to their diverse formats and complex nature. This paper employed the MOOC dataset for sentiment analysis. This research paper aims to provide essential information on how to preprocess reviews in order to determine sentiment and analyze whether they are positive or negative and neutral. Various text preprocessing techniques are applied to improve the efficiency of text classification includes stemming, lemmatization, tokenization, removing emoticons, removing stopwords, and spelling correction are applied to the unstructured text. Data preprocessing can be done with the help of Natural Language Processing (NLP) is a vital component in sentiment analysis, as it helps in preprocessing text data, extracting features, classifying sentiment, understanding language nuances, and presenting results in an interpretable manner. The accuracy of the text is calculated before and after preprocessing. Results proved that the accuracy of algorithm was significantly improved after applying the preprocessing steps. This research work demonstrates the impact of text preprocessing on the accuracy, particularly highlighting the improvement in machine learning algorithms. Proper preprocessing techniques contribute to improved prediction accuracy and reduced computational time, ultimately leading to better outcomes in various applications.

Keywords: Machine learning, Sentiment analysis, data Pre-processing, Natural Language Processing.

I. INTRODUCTION

Sentiment analysis is a crucial research domain, focuses on evaluating users' opinions on various aspects, including individual and collective perspectives. It employs data mining techniques to extract and analyze sentiments from textual data. By identifying and understanding these emotions, sentiment analysis plays a vital role in assisting decision-making processes related to products and services, ultimately improving their quality and customer satisfaction. Focusing on public opinions and sentiments it contributes to the improvement and overall decision-making processes. The proposed work involves examining trending MOOC information and gathering diverse user data to enhance the quality of Course. Data challenges, including noisy, missing, wrong, or inconsistent data, can significantly impact analysis, pattern recognition, and decision-making processes. To address these issues in mixed and unstructured data, preprocessing is crucial to transform them into structured or ordered representations. Data preparation holds immense importance

for several reasons: ensuring data quality and database reliability, facilitating analysis processes, enabling the application of algorithms to handle noisy and missing data, and ultimately enhancing the accuracy and effectiveness of data models. High-quality data is essential for generating reliable and valuable insights. Sentiment analysis, which primarily evaluates users' opinions from social media content, classifies text into neutral, negative, and positive categories. Researchers have employed diverse methods to train and classify MOOC datasets, resulting in varying performance outcomes. These approaches aim to enhance the accuracy and effectiveness of sentiment analysis in understanding public sentiment and opinions, thereby providing valuable insights for decision-making and improving Course for the welfare of the student.

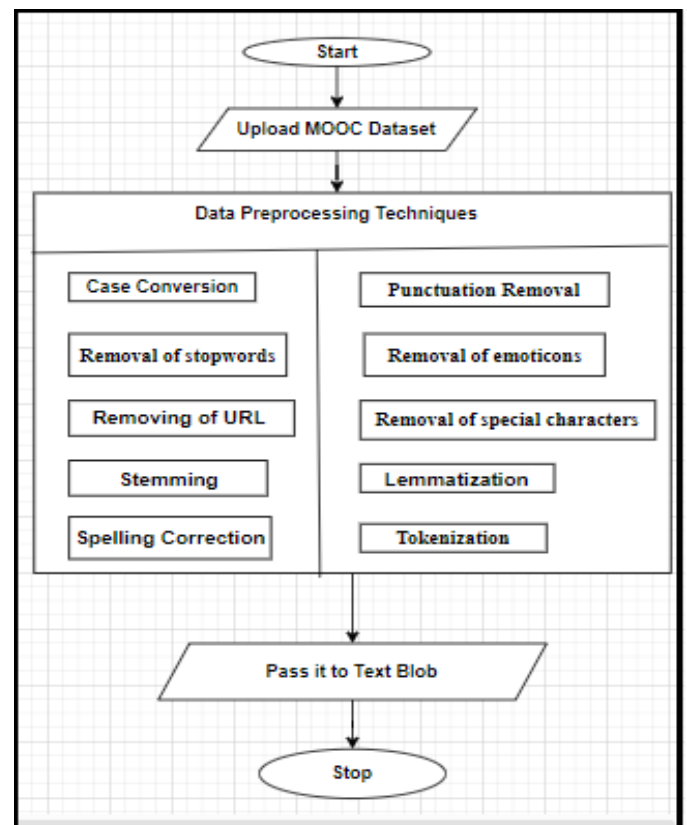


Figure 1: Architecture of MOOC Pre-Processing

With the increasing popularity of online learning, MOOCs have gained significant attention in recent years. MOOCs, or Massive Open Online Courses, provide individuals with the opportunity to expand their knowledge and skills from the comfort of their own homes. These courses cover a wide range of topics, from computer science to arts and humanities, making them accessible to a diverse audience worldwide. The importance of reviewing MOOC courses, it is crucial to understand the impact it can have on one's learning journey.

Reviews are insightful from other learners who have completed the course and can provide valuable feedback on the quality of the material and the effectiveness of the instructors. These reviews is a potential for students can gain a better understanding of what to expect, helping them make informed decisions before enrolling in a course.

The flow of experimental process is described in the form of flow chart figure 1 which describes the uploading of dataset and the application of preprocessing techniques to the unstructured data. Preprocessing methods are crucial in extracting accurate sentiment from large, unstructured data sets, particularly those sourced from social media platforms. These techniques help clean and organize information, making it more manageable for decision-making processes. However, when dealing with massive datasets, it's essential to choose methods that are computationally efficient while maintaining high accuracy. Natural Language Processing (NLP) plays a crucial role in sentiment analysis, which is the process of identifying, extracting, and classifying opinions, emotions, and attitudes expressed in textual data. NLP techniques are employed to analyze and interpret human language in a more comprehensive and sophisticated manner. This enables businesses, researchers, and individuals to gain valuable insights from textual data and make informed decisions based on public opinion and sentiment.

II. RELATED WORK

In [1] data preprocessing plays a crucial role in ensuring accurate results. It involves various techniques and data augmentation methods that help improve the quality of sentiment analysis models. Tokenization, stopword removal, stemming, and lemmatization are some of the preprocessing techniques used to standardize text data for better analysis. These methods help in generating more data points for the models to learn from, reducing over fitting, and ultimately achieving higher accuracy and generalizability. Data preprocessing is essential for advancing sentiment analysis research and applications. By utilizing these techniques, we can better capture the sentiments expressed in text data and improve the overall performance of sentiment analysis models.

In [6] Akrivi Krouska , “et al.” when dealing with large datasets containing unstructured text, applying various preprocessing methods can significantly impact the final sentiment classification. The first step in sentiment analysis is data cleaning, which involves removing irrelevant information such as stopwords and punctuation marks. Also, a proper normalization procedure can help by converting all text to lowercase, ensuring consistency in the analysis. Applying preprocessing techniques like stemming or lemmatization helps reduce the dimensionality of the dataset, leading to more accurate sentiment predictions. However, over-processing data can also adversely affect the results, so finding the right balance is crucial. It is clear that preprocessing techniques are essential in enhancing the accuracy of Twitter Sentiment Analysis. By carefully selecting and implementing these methods, researchers can improve the quality of sentiment classification results significantly.

Matthew J. Denny et al.[8] proposed unsupervised learning methods for political science text data research and preprocessing decisions in their work. Their approach involves introducing statistical procedures and software that help in understanding the sensitivity of preprocessing decisions. This allows researchers to characterize variability changes in preprocessing and make informed decisions based on the comparative analysis of data preprocessing in different

specifications across relative documents. The primary objective of this method is to minimize the risk for researchers by providing a theoretical foundation for preprocessing decisions. By analyzing the sensitivity of the results in preprocessing, researchers can better understand the impact of their choices on the overall analysis and make more informed decisions.

In [3] Machine learning algorithms have become a vital tool in sentiment analysis. The accuracy of these algorithms heavily depends on the preprocessing steps applied to the data before training the model. In article, it explores how different preprocessing techniques can affect the overall performance of machine learning models in sentiment analysis. Tokenization is a fundamental step in preprocessing text data. It involves breaking down the text into individual tokens or words. This step helps in building a vocabulary and analyzing the text at a granular level. Stopwords are common words that do not contribute much to the overall meaning of the text. Removing stopwords can help in reducing noise and improving the efficiency of machine learning models. Stemming and lemmatization are techniques used to normalize words to their base or root form. This process helps in reducing the complexity of the text data and improving the accuracy of machine learning models.

After conducting an extensive review of the scientific literature, it has been observed that there is a lack of extensive comparison of sentiment polarity classification methods for MOOC text. To address this gap, it is essential to focus on the role of text preprocessing in sentiment analysis and evaluate the impact of feature selection and representation on classification performance. Text preprocessing plays a vital role in sentiment analysis, as it helps improve the accuracy and efficiency of sentiment polarity classification. By investigating the role of text preprocessing in sentiment analysis and evaluating the impact of feature selection and representation on classification performance, researchers can gain insights into improving sentiment polarity classification methods for MOOC data. This will contribute to the advancement of sentiment analysis techniques and help in better understanding the public opinion and emotions expressed on social media platforms like Twitter.

III. METHODOLOGY

Data preprocessing is a crucial step in transforming raw data into a more manageable and understandable format, making it easier to work with and analyze. This process plays a significant role in improving the performance of Data Mining algorithms. In the context of the MOOC (Massive Open Online Course) dataset, Figure 1 demonstrates various data preprocessing techniques that can be applied to enhance the quality of the data for analysis.

3.1 DATASET

A MOOC (Massive Open Online Course) dataset typically refers to a collection of data generated from massive open online courses, which are online courses available to a large number of participants worldwide. A MOOC review dataset typically consists of a collection of reviews, ratings, and feedback provided by learners who have participated in various MOOCs. These datasets are valuable for researchers, educators, and course developers, as they offer insights into the effectiveness of MOOCs, the preferences of learners, and areas that need improvement. A MOOC review dataset may contain the following information:

3.1.1 Review Text

The main content of the MOOC data set is review text, where learners share their opinions, experiences, and feedback about the course. This textual data can be analyzed using natural language processing techniques to extract sentiment, topics, and key phrases.

3.1.2 Rating

A numerical score assigned by the learner to rate the course's overall quality, difficulty, or usefulness. Ratings can be used to calculate average scores for different aspects of the course and compare them across different courses. The figure 2 bar graph shows the review rating for the MOOC course by the students. This graph is illustrated using Python language in pycharm community edition.

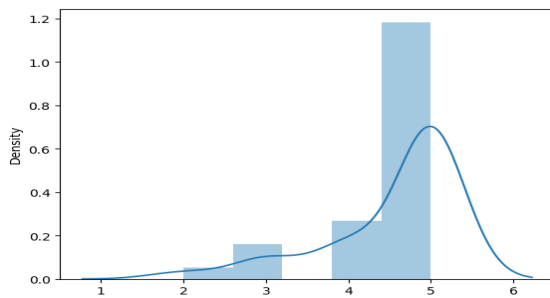


Figure 2: MOOC Review Rating

3.1.3 Course Information

Details about the MOOC, such as its title, provider, subject, level (beginner, intermediate, or advanced), and duration. This information can help identify patterns in learners' preferences and course effectiveness based on these factors.

3.1.4 Learner Information

Demographic details about the reviewer, such as age, gender, location, and educational background. This information can help understand the target audience for a particular course and identify any potential disparities in learners' experiences.

3.1.5 Timestamps

The date and time when the review was submitted can help track changes in learners' perceptions over time and identify any trends in course quality or popularity.

3.1.6 Additional Feedback

Other relevant information provided by learners, such as suggestions for improvement, specific aspects they liked or disliked, and their learning outcomes. By analyzing a MOOC review dataset, researchers and educators can gain valuable insights into the effectiveness of online learning platforms, identify areas for improvement, and develop strategies to enhance the learning experience for future participants. This research work considers the Review and Rating fields for sentiment analyzing.

3.2. Preprocessing Techniques

Data preparation is a crucial process involving various techniques to ensure that the data is properly organized and ready to serve as input for a Decision Making (DM) algorithm. This step helps in enhancing the efficiency and accuracy of the algorithm by addressing issues such as missing values, inconsistencies, and data formatting. Proper data preparation

can lead to better insights and more reliable results in decision-making processes.

3.2.1 Data Preparation

The importance of data preparation as a set of techniques that help initialize the data appropriately for use as input in Data Mining (DM) algorithms. In the context of this paper, the data is collected from Massive Open Online Courses (MOOCs) and focuses on the reviews of online courses. The primary goal of data preparation in this context is to ensure that the data is suitable for analysis and to improve the performance of DM algorithms. The MOOC data related to online course reviews is properly prepared for analysis. This will ultimately lead to more accurate and reliable insights, improving the performance of Data Mining algorithms in understanding the public opinion and experiences of learners in the context of online education.

3.2.2 Case conversion

It is an essential preprocessing step in natural language processing (NLP) and text analysis, where the goal is to convert all words in a text either to lower case or upper case. This process helps to standardize the text, making it easier to compare and process. Converting words to either lower case or upper case removes the difference between "Text" and "text," ensuring consistency and uniformity in the textual data.

3.2.3 Removal of stop words

In English, stop words are: the, is, at, which, on and so on. They have lexical content and the presence of these words can fail the required results. We have filtered out stopwords from our dataset as they are conventionally high in frequency and are not giving any useful information. Delete Stopwords are common words in a language that carry little to no semantic value, such as "the," "is," "at," "which," "on," and so on. In English, stopwords comprise approximately 20-30% of the words in a typical text corpus. Although they have lexical content, their presence in a text often does not contribute significantly to its meaning or context.

Removing stopwords from a dataset is a common practice in natural language processing (NLP) and text analysis, as they can sometimes confuse machine learning algorithms and fail to provide relevant information for analysis. As mentioned by Baradad and Mugabushaka (2015), stopwords are conventionally high in frequency and may not contribute much to the overall understanding of a text.

3.2.4 Removal of emoticons

Emoticons are visual representations of emotions and sentiments, often used in digital communication to convey the user's feelings more effectively. They can indicate happiness, sadness, anger, or other emotions and help in expressing the user's attitude towards a particular subject where they help in conveying emotions more accurately and efficiently. However, as mentioned by Khan et al. (2017), the use of emoticons can sometimes be complex, particularly when sarcasm is involved. Sarcasm is a linguistic approach where the intended meaning is opposite to the literal meaning of the words used. It can be challenging to identify whether a user is being sarcastic or genuinely expressing their feelings through emoticons, especially in short messages like tweets. In this research, it is decided to remove emoticons from the training dataset to avoid any complications in understanding the user's sentiment. This step can help simplify the text analysis process and ensure that the machine learning algorithms focus on the textual content without being influenced by the emoticons.

3.2.5 Removing of URL

URLs, or Uniform Resource Locators, are the addresses used to access websites and other online resources. While they serve as crucial links to external information, URLs can sometimes complicate the process of sentiment analysis in text. This is because URLs can contain words or phrases that may influence the perceived sentiment of a statement, as shown in the example you provided. In the given review, the presence of the URL "http://ecstasy.com" along with the words "not working" and "excited" can create ambiguity in determining the overall sentiment of the tweet. The word "not working" suggests a negative sentiment, while "excited" indicates a positive sentiment. However, considering the URL as a whole, it might be challenging to classify the sentiment as clearly positive or negative.

To avoid such complications and ensure accurate sentiment analysis, it is common practice to remove URLs from the text being analyzed. By doing so, the focus remains on the actual text content, and the machine learning algorithms can more effectively identify and classify the sentiment without being influenced by the URLs.

3.2.6 Stemming and Lemmatization

Stemming and lemmatization are two common techniques used in natural language processing (NLP) to preprocess text data. Both methods help simplify text by reducing words to their base or root form, making it easier to analyze and understand the underlying meaning. Stemming involves removing inflectional affixes from words to obtain their base form. This technique is particularly useful for languages like English and Spanish, which have a regular morphological structure. For instance, the words "playing" and "plays" can be stemmed to their base form "play," allowing for more efficient analysis and comparison of words with similar meanings.

On the other hand, lemmatization goes a step further by considering the morphological analysis of words and reducing them to their dictionary form, citation form, or canonical form. This process takes into account the word's context and its relationships with other words in a sentence. Lemmatization is more accurate than stemming, as it preserves the original word's meaning and context, making it a preferred choice for tasks like sentiment analysis, where maintaining the context is crucial. Both stemming and lemmatization are essential data preprocessing techniques in NLP. Stemming focuses on removing inflectional affixes to simplify words, while lemmatization considers the morphological analysis to reduce words to their proper dictionary form, preserving the original meaning and context. These techniques are particularly useful in tasks like text analysis and sentiment analysis, as they help in simplifying the data and making it easier to work with and understand.

3.2.7 Removal of special characters

Special characters, such as square brackets ([]), curly braces ({}), and parentheses (()) are often removed during text preprocessing, particularly when dealing with tasks like sentiment analysis or polarity assignment. These characters can create logical inconsistencies and incompatibilities in the analysis process, making it difficult to accurately interpret the sentiment expressed in the text. By removing these special characters, the focus remains on the actual content of the text, allowing for a more straightforward analysis of the sentiment. This is especially important when using machine learning algorithms or natural language processing techniques, as these

characters can interfere with the algorithms' ability to understand and classify the text accurately. Removing special characters during text preprocessing helps to resolve potential logical and compatibility issues, ensuring a more accurate and efficient sentiment analysis or polarity assignment process.

3.2.8 Punctuation Removal

Punctuation marks, such as commas (,) and colons (:), are essential elements in written language that help convey meaning and structure within a text. However, in certain textual analysis tasks, these punctuation marks may not directly contribute to the meaning being analyzed and can be considered as noise. In such cases, it is appropriate to remove them from the input text to simplify the analysis process. The removal can simplify the analysis process in certain tasks. Nevertheless, it is crucial to consider the potential impact of punctuation on the analysis and proceed accordingly.

3.2.9 Spelling Correction

Correcting spelling errors in a text is crucial for accurate analysis and understanding. With the help of automated spell-checking tools and language models, it is possible to correct incorrect words based on the selection of more probable words. These automated tools use various techniques such as contextual analysis, language models, and machine learning algorithms to identify and suggest the most likely correct word for a given misspelled word. By replacing incorrect words with their correct counterparts, the overall readability and comprehensibility of the text are significantly improved. The correcting spelling errors in a text using automated selection of more probable words can greatly enhance the accuracy and effectiveness of various textual analysis tasks. This not only improves the readability of the text but also ensures that the analysis is based on the correct interpretation of the content.

3.2.10 Tokenization

Tokenization is a crucial preprocessing step in natural language processing (NLP) and text analysis. It involves breaking down a piece of text into smaller units called tokens, which can be words, numbers, or even punctuation marks. These tokens serve as building blocks for further analysis and processing of the text. Word Tokenization: This is the most common method, where the text is split into individual words. For example, the sentence "I love playing guitar" would be tokenized as ['I', 'love', 'playing', 'guitar']. By breaking down the text into smaller, manageable units, tokenization allows for more efficient and accurate analysis of the text's content.

V. EXPERIMENTS AND RESULTS

In this research work the MOOC dataset is collected from Kaggle repository and retrieved specifically English language reviews using particular keywords. To perform sentiment analysis MOOC dataset involves several steps. Import Required Libraries by importing the necessary libraries for sentiment analysis, such as NLTK, and TextBlob. Load the preprocessed MOOC dataset for Preprocessing includes tokenized, stemmed or lemmatized, punctuation-free, emoji-removed, spelling-corrected, stop-word-removed, and case-converted text.

The various data preprocessing techniques is applied to the MOOC dataset and the comparison of the unstructured data before and after preprocessing is illustrated in the Table:

Table :1 Data Before Pre-processing

S.NO	Data Pre-processing Method
1	0 good and interestng
2	1 This clas is very helpfl to me. Currently, I'm...
3	2 like!Prof and TAs are helpful and the discussi...
4	3 Easy to follow and includes a lot basic and im...
5	4 Realy nice teacher!I could got the point eazli...

Before pre-processing, this data in the MOOC dataset might be unstructured, contain inconsistent values, and may not be in a format suitable for machine learning algorithms. The primary goal of pre-processing is to clean, transform, and structure the data to make it more accurate, reliable, and easier to analyze.

Table 2: Data After Pre-processing techniques

S.NO	DATA PRE-PROCESSING METHOD	AFTER PRE-PROCESSING
1	Spelling Correction	0 good and interesting 1 This class is very helpful to me. Currently, I'... 2 like!Prof and was are helpful and the discussi... 3 Easy to follow and includes a lot basic and im... 4 Really nice teacher!I could got the point earl...
2	Case Conversion	0 good and interesting 1 this class is very helpful to me. currently, i... 2 like!prof and tas are helpful and the discussi... 3 easy to follow and includes a lot basic and im... 4 really nice teacher!i could got the point eazli...
3	Punctuation Removal	0 good and interesting 1 this class is very helpful to me currently im ... 2 likeprof and tas are helpful and the discussio... 3 easy to follow and includes a lot basic and im... 4 really nice teacheri could got the point eazli...
4	Removal of stopwords	0 good interesting 1 class helpful currently im still learning clas... 2 likeprof tas helpful discussion among students... 3 easy follow includes lot basic important techn... 4 really nice teacheri could got point eazliy v
5	Removal of special characters	0 good interesting 1 class helpful currently im still learning clas... 2 likeprof tas helpful discussion among students... 3 easy follow includes lot basic important techn... 4 really nice teacheri could got point eazliy v

6	Stemming and Lemmatization	0 good interest 1 class help current im still learn class make l... 2 likeprof help discuss among student quit activ... 3 easi follow includ lot basic import techniqu u... 4 realli nice teacheri could got point earli v
7	Frequency word removal	[('course', 21), ('content', 8), ('excellent', 6), ('good', 5), ('interesting', 5), ('great', 5), ('lot', 4), ('knowledge', 4), ('one', 4), ('peer', 4)]
8	Tokenization	0 [good, interest] 1 [class, help, current, im, still, learn, class... 2 [likeprof, help, discuss, among, student, quit... 3 [easi, follow, includ, lot, basic, import, tec... 4 [realli, nice, teacheri, could, get, point, ea...

In the pre-processing of MOOC course review datasets the experiments and results may involve various techniques tokenization, stop-word removal, stemming, and lemmatization can be applied to normalize and simplify the text. This can improve the efficiency of subsequent natural language processing tasks. Experimental results from pre-processing MOOC datasets will vary depending on the specific techniques applied and the goals of the analysis. Generally, pre-processing aims to improve data quality, facilitate efficient machine learning, and enhance the overall understanding of the underlying patterns and relationships within the data.

The average polarity of sentiment analysis refers to the overall sentiment score obtained from analyzing a collection of texts using sentiment analysis techniques. The average polarity score provides a quantitative measure of the overall sentiment expressed in the dataset. A positive average polarity indicates that the majority of texts have a positive sentiment, while a negative average polarity suggests that the majority of texts have a negative sentiment. A neutral average polarity would mean that the dataset has an equal distribution of positive, negative, and neutral sentiments.

Average polarity:

Sentiment(polarity=0.5333333333333333, subjectivity=0.6166666666666667)

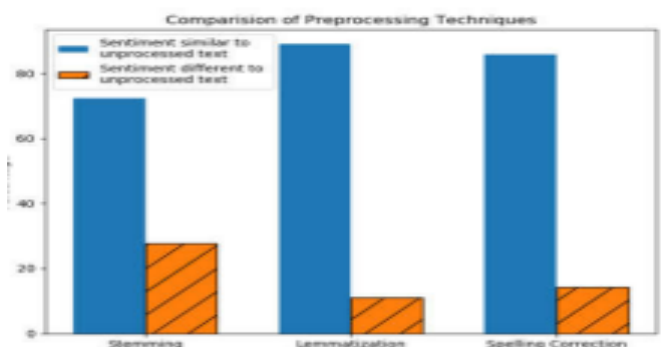


Figure 3: Comparison before and after preprocessing

Before preprocessing, raw data often contains issues that can hinder the effectiveness of machine learning models and analysis. In summary, comparing the dataset before and after preprocessing highlights the improvements in data quality, feature engineering, and analysis capabilities that result from addressing issues such as missing values, irrelevant information, inconsistencies, noise, and unstructured data.

VI. CONCLUSION

In this research paper it explores the impact of various preprocessing techniques on sentiment analysis performance across three datasets of MOOC, with one general and the others focused on specific topics. The study evaluates the effectiveness of data preprocessing techniques. This work emphasizes the importance of preprocessing in improving the overall quality of textual data and highlights the potential for more accurate sentiment analysis in various applications, such as social media monitoring, customer feedback analysis, and market trend predictions.

In conclusion, the effectiveness of preprocessing techniques in cleaning and preparing text data for analysis has been demonstrated in various studies. While some research has focused on basic programming for preprocessing, others have employed standard methods like tokenization and stemming to achieve better results. In our study, we have successfully combined various preprocessing techniques, including stemming, tokenization, special character and punctuation removal, and simple programming. This combination has proven to be more efficient in the cleaning process, ensuring a higher quality dataset for further analysis. By leveraging a comprehensive set of preprocessing techniques, we can better understand and extract valuable insights from textual data, such as MOOC reviews, to inform improvements in online learning experiences.

VII. FURTHER ENHANCEMENT

In this research the preprocessing techniques is applied to the text data from the MOOC dataset has been done in further the dataset will be Splitted into Training and Testing Sets: Divide the dataset into training and testing sets to evaluate the performance of your sentiment analysis model. A common split is 80% for training and 20% for testing. Select a suitable sentiment analysis tool or algorithm for your analysis. Some popular options include NLTK's SentimentIntensityAnalyzer, TextBlob's sentiment polarity will be utilized. Perform Sentiment Analysis on Training Data Apply the chosen sentiment analysis tool to the training set to assign sentiment scores (positive, negative, or neutral) to each review. This step will also help to fine-tune the tool if necessary. Finally Evaluate the Model and Calculate the accuracy, precision, recall, and F1-score of your sentiment analysis model on the testing set to assess its performance. The Python libraries like scikit-learn or pandas are used to perform these calculations. Analyze Sentiment Distribution is used to Visualize the sentiment distribution across the MOOC reviews using a bar chart, pie chart, or word cloud. This will help you understand the overall sentiment of the learners towards the MOOCs and identify any trends or patterns.

References

- [1] Huu-Thanh Duong^{1*} and Tram-Anh Nguyen-Thi, "A review: preprocessing techniques and data augmentation for sentiment analysis", Ho Chi Minh City Open University, Ho Chi Minh City, Vietnam
- [2] "Effective Text Data Preprocessing Technique for Sentiment Analysis in Social Media Data"

- [3] Saqib Alam, Nianmin Yao, "The impact of preprocessing steps on the accuracy of machine learning algorithms in sentiment analysis"
- [4] A. Naresh, P. Venkata Krishna, "An efficient approach for sentiment analysis using machine learning algorithm"
- [5] Anuja P Jain, Asst. Prof Padma Dandannavar "Application of Machine Learning Techniques to Sentiment Analysis".
- [6] Akrivi Krouska, Christos Troussas, Maria Virvou "The effect of preprocessing techniques on Twitter Sentiment Analysis".
- [7] T. Nikil Prakash¹, Dr. A. Aloysius "Data Preprocessing In Sentiment Analysis Using Twitter Data"
- [8] Matthew J. Denny and Arthur Spirling "Text Preprocessing For Unsupervised Learning: Why It Matters, When It Misleads, And What To Do About It" DOI: <https://doi.org/10.1017/pan.2017.44>, Year: 2017
- [9] Jaspreet Singh, Gurvinder Singh and Rajinder Singh, "Optimization of sentiment analysis using machine learning classifiers"
- [10] E. M. Badr, Mustafa Abdul Salam, Mahmoud Ali, Hagar Ahmed "Social Media Sentiment Analysis using Machine Learning and Optimization Techniques"
- [11] Neethu M S, Rajasree R, "Sentiment Analysis in Twitter using Machine Learning Techniques"
- [12] Shreyas Wankhede, Ranjit Patil, Sagar Sonawane, Prof. Ashwini Save "Data Preprocessing for Efficient Sentimental Analysis"
- [13] Shamantha Rai B, Sweekriti M Shetty, "Sentiment Analysis Using Machine Learning Classifiers: Evaluation of Performance"
- [14] Rajkumar S. Jagdale, Vishal S. Shirsat and Sachin N. Deshmukh, "Sentiment Analysis on Product Reviews Using Machine Learning Techniques"