

Sentiment Analysis of User Comments on Machine Learning Posts: A Data-driven Approach Using R*

Joseph Chung

Apr 2, 2023

Abstract

In this study, we present a comprehensive data analysis paper that utilize a sentiment model in R to analyze user comments from a dataset containing post_ids and comments related to the topic of Machine Learning. We aim to investigate sentiment polarity of user comments and its distribution across different posts, identifying trends and patterns that may provide valuable insights into the Machine Learning community's perception and engagement.

Table of contents

1	1. Introduction	2
1.1	1.1 Background and motivation	2
1.2	1.2 Research questions	3
1.3	1.3 Scope and limitations	3
1.4	2.2 Sentiment analysis techniques in R	4
1.5	2.3 Applications of sentiment analysis in Machine Learning research	4
2	3. Data Collection and Preprocessing	6
2.1	3.1. Description of the dataset	6
2.2	3.2. Data cleaning and preprocessing	6
2.2.1	3.2.1. Handling missing values	6
2.2.2	3.2.2. Text normalization	6
2.2.3	3.2.3. Tokenization	6
2.2.4	3.2.4. Slopword removal	6
2.2.5	3.2.5. Stemming/Lemmatization	6
3	4. Methodology	6
3.1	4.1. Sentiment analysis model selection	6
3.2	4.2. Model training and validation	6
3.3	4.3. Model implementation in R	6
3.4	4.4. Evaluation metrics	6

*Code and data are available at: <https://github.com/UtopianYoungChung/Sentiment-Analysis-Using-R.git>

4	5. Results	6
4.1	5.1. Sentiment polarity distribution	6
4.2	5.2. Correlation between post attributes and sentiment	6
4.3	5.3. Identification of trends and patterns	6
4.4	5.4. Comparison with previous studies	6
5	6. Discussion	6
5.1	6.1. Interpretation of results	6
5.2	6.2. Implications for Machine Learning community	6
5.3	6.3. Limitations of the study	6
5.4	6.4. Future research directions	6
6	7. Conclusion	6
7	References	6
8	Appendices	6
8.1	A. R and Python code for data preprocessing	6
8.2	B. R code for model training and validation	6
	C. R code for sentiment analysis and visualization	6

List of Figures

List of Tables

1 1. Introduction

The rapid growth of online forums and communities discussing various topics, including Machine Learning (ML), has led to an abundance of user-generated content in the form of text data. Analyzing this data can reveal valuable insights into user perceptions, opinions, and sentiment towards specific subjects. Sentiment analysis, or opinion mining, is a natural language processing technique used to determine the sentiment expressed in a piece of text, such as positive, negative, or neutral. This paper focuses on employing sentiment analysis techniques to understand the sentiment polarity of user comments related to the topic of Machine Learning.

1.1 1.1 Background and motivation

Machine Learning has become a popular field of study and research due to its vast applications in various industries. As a result, online communities have emerged where people discuss ML topics, share knowledge, and express their opinions. Analyzing the sentiment of these user comments can provide insights into the perception of the ML community, identify common issues or concerns, and uncover emerging trends. This understanding could be valuable for researchers, educators, and industry professionals who aim to improve their work based on community feedback.

1.2 1.2 Research questions

This study aims to address the following research questions:

1. What is the overall sentiment polarity (positive, negative, or neutral) of user comments related to Machine Learning?
2. How is the sentiment polarity distributed across different posts in the dataset?
3. Are there any correlations between post attributes (e.g., post length, engagement) and sentiment polarity?
4. Can any trends or patterns in sentiment polarity be identified over time or in relation to specific subtopics within Machine Learning?

1.3 1.3 Scope and limitations

The scope of this study is limited to the analysis of a dataset containing post IDs and comments related to the topic of Machine Learning. The results and interpretations are based solely on the data provided and may not be generalizable to other online communities or topics. The sentiment analysis model employed in this study is subject to the limitations inherent in natural language processing techniques, such as context ambiguity and handling of sarcasm or irony. Additionally, the study does not account for demographic information or biases present in the dataset, which could potentially influence the sentiment analysis results.

2. Literature Review {#sec-literature review} ## 2.1 Sentiment analysis in online communities Over the past decade, sentiment analysis has become increasingly popular as a tool to mine user opinions and emotions from online platforms, such as social media, blogs, and forums (**citePang?**). Numerous studies have focused on the application of sentiment analysis in various domains, such as politics (**citeTumasjan?**), finance (**citeBollen?**), and consumer reviews (**citeDave?**). These studies demonstrate the potential of sentiment analysis as a valuable technique for understanding user-generated content and informing decision-making processes.

1.4 2.2 Sentiment analysis techniques in R

R is a widely used programming language for statistical analysis and data visualization, with several libraries and packages available for natural language processing and sentiment analysis (R Core Team 2020). Some of the most popular sentiment analysis packages in R include tidytext (**citeSilge?**), syuzhet (**citeJockers?**), and sentiment (**citeRinker?**). These packages provide various algorithms and techniques for sentiment analysis, such as lexicon-based approaches, machine learning models, and deep learning methods. Researchers have employed these techniques to analyze sentiment in diverse contexts, ranging from Twitter data (**citeWijeratne?**) to news articles (**citeKearney?**).

1.5 2.3 Applications of sentiment analysis in Machine Learning research

Sentiment analysis has been employed in Machine Learning research to understand user perceptions, opinions, and trends related to the field. For example, Aliza (**citeAliza?**) analyzed the sentiment of tweets related to Machine Learning conferences, revealing insights into the community's responses to specific presentations and events. In another study, Wang (**citeWang?**) investigated the sentiment of user comments on popular Machine Learning platforms, such as GitHub and Stack Overflow, to identify trends and patterns in the community's engagement with specific algorithms and techniques. These studies highlight the potential of sentiment analysis to uncover valuable insights into the Machine Learning community and inform future research directions.

The literature demonstrates the growing interest in sentiment analysis as a valuable tool for understanding user-generated content in various domains. R has emerged as a popular platform sentiment analysis due to its extensive libraries and packages. However, there remains a need for more comprehensive studies that investigate the sentiment of user comments related to Machine Learning, particularly in the context of online communities where users share their knowledge and opinions on the topic. This study aims to address this gap by employing sentiment analysis techniques in R to analyze a dataset of post IDs and comments related to Machine Learning.

2 3. Data Collection and Preprocessing

2.1 3.1. Description of the dataset

2.2 3.2. Data cleaning and preprocessing

2.2.1 3.2.1. Handling missing values

2.2.2 3.2.2. Text normalization

2.2.3 3.2.3. Tokenization

2.2.4 3.2.4. Stopword removal

2.2.5 3.2.5. Stemming/Lemmatization

3 4. Methodology

3.1 4.1. Sentiment analysis model selection

3.2 4.2. Model training and validation

3.3 4.3. Model implementation in R

3.4 4.4. Evaluation metrics

4 5. Results

4.1 5.1. Sentiment polarity distribution

4.2 5.2. Correlation between post attributes and sentiment

4.3 5.3. Identification of trends and patterns

4.4 5.4. Comparison with previous studies

5 6. Discussion

5.1 6.1. Interpretation of results

5.2 6.2. Implications for Machine Learning community

5.3 6.3. Limitations of the study

5.4 6.4. Future research directions

6 7. Conclusion

7 References

8 Appendices

8.1 A. R and Python code for data preprocessing