

Improvement of User Sentiment Classification Model for the Indomaret Poinku Application Using the Naive Bayes Method

Sofyan Hidayat^{1*}, Nining Rahaningsih², Raditya Damar Dana³, Mulyawan⁴

¹Informatics Engineering, STMIK IKMI Cirebon

²Accounting Computer, STMIK IKMI Cirebon

³Informatics Management, STMIK IKMI Cirebon

⁴Software Engineering, STMIK IKMI Cirebon

sofyanhidayat1710@gmail.com^{1*}, niningr157@yahoo.co.id², radith_damar@yahoo.com³, mulyawan00@gmail.com⁴

Abstract

The Indomaret Poinku application provides a platform for users to give reviews related to products and services. With the increasing number of user reviews, an effective method is needed to automatically analyze opinions. This research aims to improve the sentiment analysis model on the Indomaret Poinku application using the Naive Bayes algorithm. The selection of this algorithm is based on its simplicity and effectiveness in text classification. To improve the model's performance, this research also applies preprocessing techniques such as cleaning, case folding, tokenizing, normalization, stopword removal, and stemming, as well as the feature selection technique Information Gain. The research method involves stages of collecting review data from the Google Play Store, manually labeling the data, and analyzing the data using TF-IDF numerical representation. The Multinomial Naive Bayes model was trained and tested using evaluation metrics such as accuracy, precision, recall, and F1-score. The evaluation results show that the developed model is capable of achieving an accuracy of 75.5%, a precision of 78%, a recall of 75%, with an average F1-score of 70.4%. Further analysis shows that features such as "good" and "great" have a significant influence in sentiment classification. The results of this study reveal that the enhancement of the Naive Bayes model through feature selection and optimization of the Preprocessing process is capable of improving sentiment classification accuracy. These findings contribute to application developers in understanding user opinions, which can be used to improve the quality of services and products.

Keywords: Sentiment Analysis, Naive Bayes, Information Gain, Indomaret Poinku, TF-IDF

1. Introduction

The development of information technology has driven the increased use of mobile applications, including retail applications like Indomaret Poinku. As a platform that allows users to review products and services, an efficient approach is needed to automatically analyze user opinions. In this context, machine learning-based sentiment analysis becomes an important method for understanding user perceptions.

Naive Bayes is one of the algorithms widely used in sentiment analysis due to its simplicity in text classification and its ability to provide relevant results [1], [2]. However, the main challenge in this research is the high volume of reviews that vary in form and structure, necessitating the optimization of data processing and feature selection to improve the model's accuracy. Previous studies have shown that preprocessing techniques such as stemming and stop word removal can improve the performance of the Naive Bayes model in sentiment analysis [3]. In addition, model evaluation using metrics such as accuracy, precision, and F1-score becomes an important factor in assessing the effectiveness of the approach used [4].

Several studies have shown the effectiveness of Naive Bayes in various sentiment analysis applications, including the issue of electricity tariff increases on Twitter, which successfully identified predominantly negative sentiment [3]. However, in the context of mobile applications, this model can still be improved with more optimal data processing techniques [5]. Therefore, this research aims to improve the sentiment analysis model in the Indomaret Poinku application by developing better data processing techniques and evaluating the model's performance using appropriate metrics.

The results of this research are expected to contribute to enhancing the understanding of user opinions and assist application developers in improving the quality of the services offered. With the optimization of the Naive Bayes method, this research also aims to fill the gap in the literature regarding the improvement of sentiment analysis performance in mobile applications.

2. Research Method

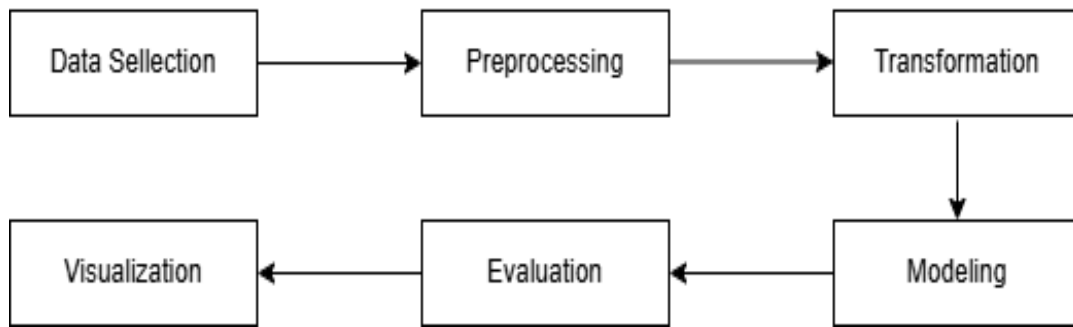


Figure 1: Stages of Knowledge Discovery in Databases (KDD)

In the image above, due to the collected data encompassing data diversity and unstructured data elements, this research method uses Knowledge Discovery in Database (KDD) [6]. Specifically for text data that requires adjustments, data preprocessing stages are necessary to organize the data and facilitate its understanding. In addition, the ability to select features that can be incorporated into the available Knowledge Discovery in Database (KDD) methods to accelerate the analysis process and minimize the amount of irrelevant data [7].

2.1. Data Selection

Data selection is the initial stage of this research, which aims to select and collect reviews from the Indomaret Poinku application available on the Google Play Store. Data collection was carried out using scraping techniques with Google Colaboratory. In this process, a total of 1,000 reviews were successfully collected and then manually labeled by an Indonesian language teacher with sentiment categories of Positive, Neutral, and Negative.

2.2. Preprocessing

Dataset preprocessing is the initial processing stage aimed at transforming raw data into a more structured form, ready for use in the classification process. At this stage, there are several important steps taken, cleaning, casefolding, tokenizing, normalization, stopword removal, and stemming to return words to their base form. This process aims to ensure that the data used is of high quality and relevant in sentiment analysis.

2.3. Transformation

Data transformation is part of the Knowledge Discovery in Database (KDD) process that serves to modify or convert reviews to make them easier to process and analyze. This process is carried out after the pre-processing stage and can include various techniques, such as word weighting using Term Frequency-Inverse Document Frequency (TF-IDF) and the application of Information Gain. With this transformation stage, the collected reviews become more structured, thereby improving the accuracy of data analysis and interpretation.

2.4. Modeling

The modeling stage in KDD aims to build a model capable of recognizing patterns or relationships in review data. At this stage, the processed and transformed data is used as training data to build a sentiment classification model. This research uses the Multinomial Naïve Bayes algorithm to build a sentiment prediction model based on the collected reviews.

2.5. Evaluation

Model evaluation in KDD is a stage to measure the performance of the model that has been built in the modeling step. The model that has been created will be tested using the test dataset that has been prepared beforehand. The model's prediction results will be compared with the expected results to assess its accuracy. Several metrics used in this evaluation include accuracy, precision, recall, and F1-score. In addition, the confusion matrix is also used to measure the extent to which the model can correctly classify sentiment. This evaluation is important to assess the effectiveness of the model in recognizing sentiment patterns in Indomaret Poinku reviews.

2.6. Visualization

The visualization stage is the final step in KDD, aimed at presenting the analysis results in a more easily understandable form. After the model successfully builds the sentiment classification, the data will be presented in the form of graphs or diagrams to facilitate interpretation. In this study, visualization is done using a word cloud, which displays the distribution of words that frequently appear in the reviews. With the visualization, the analysis results can be more easily understood and used as a basis for decision-making.

3. Result and discussion

3.1. Data Selection

The data used in this research was obtained from user reviews of the Indomaret Poinku application on the Google Play Store. The reviews collected are in Indonesian and were obtained using data scraping techniques, totaling 1,000 reviews. After the data collection process is complete, the data is then stored in XLSX format.

3.2. Preprocessing

The next stage is data preprocessing, which is the process of preparing data before it is used in machine learning models or algorithms. This stage aims to clean, format, and adjust raw data to make it ready for modeling and analysis. Here are the pre-processing stages applied in this research.

3.2.1. Cleaning

Cleaning is the process of removing special characters, punctuation, and other symbols that do not provide significant meaning in sentiment analysis.

Table 1: Cleaning

Appearance	
Before	After
Kuncinya setiap minggu di update terus katanya mbk indomaret 😊	Kuncinya setiap minggu di update terus katanya mbk indomaret

3.2.2. Casefolding

Case folding converts the entire text to lowercase to ensure consistency. This step ensures that words with the same meaning are not considered different just because of capitalization differences.

Table 2: Casefolding

Appearance	
Before	After
Kuncinya setiap minggu di update terus katanya mbk indomaret	kuncinya setiap minggu di update terus katanya mbk indomaret

3.2.3. Tokenizing

Tokenizing is the process of breaking down text into separate words. Each word in the review can be analyzed as a separate unit, making it easier for the model to understand each word in its context.

Table 3: Tokenizing

Appearance	
Before	After
kuncinya setiap minggu di update terus katanya mbk indomaret	['kuncinya', 'setiap', 'minggu', 'di', 'update', 'terus', 'katanya', 'mbk', 'indomaret']

3.2.4. Normalization

Normalization aims to standardize text by removing numbers or other elements that are not relevant for analysis. For example, if there are numbers that do not have significant meaning, those numbers will be removed.

Table 4: Normalization

Appearance	
Before	After
['aplikasi', 'ngga', 'bisa', 'berjalan', 'di', 'android', '14']	['aplikasi', 'ngga', 'bisa', 'berjalan', 'di', 'android', '']

3.2.5. Stopword removal

Stopword removal is the process of removing common words that frequently appear but do not have specific meaning in sentiment analysis, such as "dan", "di", "yang", and so on. By removing these words, the model can focus more on words that have significant meaning.

Table 5: Stopword removal

Appearance	
Before	After
['kuncinya', 'setiap', 'minggu', 'di', 'update', 'terus', 'katanya', 'mbk', 'indomaret']	['kuncinya', 'minggu', 'update', 'mbk', 'indomaret']

3.2.6. Stemming

Stemming is the process of converting words into their base or root forms. This step helps reduce word variations that have the same meaning but different forms, such as "kuncinya" becoming "kunci".

Table 6: Steaming

Appearance	
Before	After
['kuncinya', 'minggu', 'update', 'mbk', 'indomaret']	['kunci', 'minggu', 'update', 'mbk', 'indomaret']

4. Conclusion

Based on the research conducted, the enhancement of the Multinomial Naïve Bayes model for sentiment analysis on the Indomaret Poinku application successfully improved classification performance. The developed model achieved an accuracy of 75.5%, precision of 78%, recall of 75%, and an average F1-score of 70.4%. This model has proven to be more optimal compared to the model without feature selection, as demonstrated by the effectiveness of the evaluation metrics used. Furthermore, the application of data processing techniques such as cleaning, case folding, tokenizing, normalization, stopword removal, and stemming, as well as the use of Information Gain in feature selection, significantly contributes to the improvement of the model's accuracy and performance. This technique allows the model to focus more on words that have high relevance to sentiment, such as "good" and "great," which have a significant impact on sentiment analysis.

Furthermore, the optimized model provides clearer insights into user opinions regarding the Indomaret Poinku application. Positive sentiment is generally related to the ease and practicality of the application, while negative sentiment is more often associated with technical issues such as login problems or promotions. These findings can be utilized by application developers to enhance services and overall user experience.

This research shows that optimizing the Naïve Bayes model through data processing and feature selection can significantly improve sentiment classification accuracy. The results obtained provide a valuable contribution to the development of sentiment analysis systems based on user reviews.

References

- [1] A. Sasmita, G. A. Pradnyana, and D. G. H. Divayana, "Sistem Analisis Sentimen Untuk Evaluasi Kinerja Dosen dengan Metode Naïve Bayes," *JST (Jurnal Sains dan Teknol.*, vol. 11, no. 2, pp. 451–462, Sep. 2022, doi: 10.23887/JSTUNDIKSHA.V11I2.44384.
- [2] N. Syafitri Kustanto, N. Gusriani, P. Studi S-, F. Mipa, U. K. Padjadjaran Jl Raya Bandung Sumedang, and J. Sumedang, "Analisis Sentimen dengan Metode Klasifikasi Naïve Bayes Dan Seleksi Fitur Information Gain," *Search (Informatic, Sci. Entrep. Appl. Art. Res. Humanism)*, vol. 21, no. 2, pp. 134–144, Nov. 2022, doi: 10.37278/INSEARCH.V21I2.524.
- [3] A. Kusuma and A. Nugroho, "Analisa Sentimen Pada Twitter Terhadap Kenaikan Tarif Dasar Listrik Dengan Metode Naïve Bayes," *J. Ilm. Teknol. Inf. Asia*, vol. 15, no. 2, pp. 137–146, Dec. 2021, doi: 10.32815/JITIKA.V15I2.557.
- [4] F. Fitriani, E. Utami, and A. D. Hartanto, "ANALISIS SENTIMEN MASYARAKAT TERHADAP PELAKSANAAN P3K GURU DENGAN ALGORITMA NAIVE BAYES DAN DECISION TREE," *Tek. Teknol. Inf. dan Multimed.*, vol. 3, no. 1, pp. 23–30, Jun. 2022, doi: 10.46764/TEKNIMEDIA.V3I1.53.
- [5] D. Pratmanto, F. Fandi, D. Imaniawan, V. Maarif, P. Studi, and T. Komputer, "Analisis Sentimen Pada Ulasan Pengguna Aplikasi Identitas Kependudukan Digital Dengan Metode Naive Bayes Dan K-Nearest," *Comput. J. Comput. Sci. Inf. Syst.*, vol. 7, no. 2, pp. 155–166, Dec. 2023, doi: 10.24912/COMPUTATIO.V7I2.26322.
- [6] and I. M. M. Raffi, A. Suharso, "Analisis Sentimen Ulasan Aplikasi Binar Pada Google Play Store Menggunakan Algoritma Naïve Bayes," *J. Inf. Technol. Comput. Sci.*, vol. 6, no. 1, pp. 1–7, 2023.
- [7] S. Widaningsih, "Perbandingan Metode Data Mining Untuk Prediksi Nilai Dan Waktu Kelulusan Mahasiswa Prodi Teknik Informatika Dengan Algoritma C4,5, Naïve Bayes, Knn Dan Svm," *J. Tekno Insentif*, vol. 13, pp. 16–25, 2019.