# Sentiment and Customer Loyalty Analysis of Shopee Using Machine Learning Algorithms

*Yoga Fitriana*

*Department of Information Systems, Universitas Pembangunan Panca Budi, Medan, Indonesia*

*Email: yoga.fitriana.yf@gmail.com*

| ARTICLE INFO | ABSTRACT |
|---|---|
| | The exponential growth of e-commerce platforms has transformed consumer shopping behavior globally, including in Indonesia. Shopee, as one of the dominant online marketplaces, continuously attracts millions of active users through competitive pricing strategies, promotional events, and digital convenience. However, understanding user satisfaction and loyalty remains a challenge in such dynamic environments. This research aims to analyze user sentiment and customer loyalty toward Shopee by integrating computational sentiment analysis techniques with behavioral survey assessment. A total of 3,000 Shopee user reviews were collected through web scraping, then processed using text mining methods and classified into positive and negative categories using two machine learning algorithms: Support Vector Machine (SVM) and Naïve Bayes Classifier (NBC). Additionally, a structured loyalty survey was distributed to 30 respondents to evaluate behavioral loyalty indicators such as repeat purchase, advocacy, and emotional attachment. The SVM algorithm demonstrated superior performance with an accuracy rate of 98%, surpassing the Naïve Bayes Classifier's 85% accuracy. The loyalty survey indicated a strong positive correlation between sentiment polarity and customer retention, revealing that satisfied users exhibit consistent repurchase intentions and brand advocacy. These findings emphasize the significance of integrating computational analytics and behavioral measurement in e-commerce performance evaluation. The results also provide managerial insights for enhancing digital service quality, consumer engagement, and long-term competitiveness in Indonesia's online retail market. |

## 1. Introduction

The digital transformation of commerce has significantly reshaped consumer purchasing behavior worldwide. In Southeast Asia, particularly in Indonesia, the rapid adoption of e-commerce platforms has created a vibrant and competitive marketplace. Among various online marketplaces, **Shopee** has emerged as one of the leading platforms, offering diverse products, user-friendly applications, and attractive promotional strategies. The convenience of online transactions, coupled with accessible mobile technology, has accelerated the penetration of Shopee into various demographic segments.

However, despite its growing popularity, sustaining customer satisfaction and loyalty remains a major strategic challenge. Users' perceptions toward Shopee's service quality, delivery time, and overall user experience are critical determinants of long-term engagement. In digital environments where switching costs are low and alternatives abound, a single negative experience can quickly influence user sentiment and retention behavior. Thus, understanding user sentiment through computational approaches and evaluating customer loyalty through behavioral indicators are essential for maintaining Shopee's market dominance.

E-commerce platforms such as Shopee, Tokopedia, and Lazada compete not only in terms of price and product diversity but also in customer experience management. Modern consumers express their satisfaction or dissatisfaction through online reviews, social media posts, and ratings. These textual feedbacks represent valuable sources of information that can be mined to extract opinions, emotions, and perceptions.

Sentiment analysis, also known as opinion mining, has become an essential method in data-driven business intelligence. By employing natural language processing (NLP) and machine learning (ML) techniques, sentiment analysis automatically classifies opinions as positive, negative, or neutral. In the e-commerce context, such analysis provides strategic insights into user satisfaction, product quality, and service reliability.

Simultaneously, customer loyalty analysis seeks to understand the behavioral patterns that indicate long-term commitment to a brand. According to Griffin (2005), loyalty extends beyond mere repeat purchases; it includes emotional attachment, trust, and advocacy. Combining sentiment analysis with loyalty measurement provides a holistic understanding of how users perceive and engage with a digital brand such as Shopee.

While numerous studies have applied machine learning algorithms to classify sentiment, relatively few have integrated computational results with behavioral loyalty indicators. Most prior works focused solely on achieving high classification accuracy without linking sentiment patterns to actual customer behavior. Consequently, businesses may achieve technical success in text classification but fail to interpret the implications

for real customer retention and satisfaction.

This study addresses this gap by exploring both **algorithmic performance** and **behavioral interpretation**. Specifically, it aims to compare two common algorithms—Naïve Bayes Classifier (NBC) and Support Vector Machine (SVM)—in analyzing Shopee customer reviews. In addition, it incorporates survey-based measurement to assess the level of customer loyalty, thereby establishing a connection between online sentiment polarity and offline loyalty behavior.

The objectives of this research are as follows: To apply machine learning algorithms (SVM and Naïve Bayes) for classifying user sentiments from Shopee reviews.

a. To compare the performance of both algorithms in terms of accuracy, precision, recall, and F1-score.
b. To analyze Shopee users' customer loyalty using survey data.
c. To identify correlations between sentiment polarity and customer loyalty dimensions.
d. To provide managerial insights for improving Shopee's marketing and customer retention strategies.
   Based on the objectives, the following research questions are formulated:
a. How effective are the SVM and Naïve Bayes algorithms in classifying Shopee user sentiments?
b. What is the comparative accuracy between these algorithms?
c. What is the level of customer loyalty among Shopee users based on behavioral indicators?
d. How does user sentiment influence or reflect customer loyalty?

**Research Significance**

This research provides both theoretical and practical contributions. Theoretically, it enhances the understanding of how computational sentiment analysis can complement behavioral research in marketing and consumer psychology. Practically, it offers empirical insights for e-commerce practitioners, especially digital marketing managers, to monitor and predict customer satisfaction and loyalty trends.

The integrated approach adopted in this study—combining machine learning-based sentiment classification with quantitative loyalty assessment—offers a novel framework that can be replicated in other e-commerce contexts. Furthermore, the findings serve as an evidence-based foundation for developing customer relationship management (CRM) strategies that emphasize data-driven personalization and service improvement.

**Structure of the Paper**

The structure of this paper is outlined as follows. Section 2 provides an overview of prior research concerning sentiment analysis, machine learning techniques, and customer loyalty. Section 3 explains the research methodology, covering the processes of data acquisition, preprocessing, algorithm implementation, and evaluation criteria. Section 4 details the experimental outcomes, visual representations, and analytical discussions. Finally, Section 5 summarizes the study by presenting the main findings, managerial implications, and future research directions.

## 2. State of the Art

### 2.1 Sentiment Analysis and Text Mining.

Sentiment analysis, also known as *opinion mining*, is a specialized branch within natural language processing (NLP) that aims to identify the polarity of opinions, emotions, or attitudes expressed in written text. As stated by Liu (2012), this technique enables researchers to extract subjective insights from textual sources such as online reviews, comments, and social media content. Within the e-commerce domain, sentiment analysis has become a crucial analytical approach for evaluating customer satisfaction and brand perception.

Meanwhile, text mining—a broader analytical framework—focuses on uncovering meaningful structures and relationships within unstructured textual data. The process typically includes several stages: data preprocessing, feature extraction, and classification. Preprocessing steps such as tokenization, stopword elimination, and stemming are vital for cleaning and normalizing text, ensuring consistency for further analysis. Feature extraction methods then convert textual elements into quantitative vectors, most commonly through the Term Frequency–Inverse Document Frequency (TF-IDF) approach, which facilitates computational analysis.

In the sentiment classification phase, various machine learning algorithms are applied to categorize textual data. Among the most commonly utilized are the Naïve Bayes Classifier (NBC) and Support Vector Machine (SVM), both valued for their efficiency and interpretability. Comparing these models provides a deeper understanding of their respective performances and the degree to which each is suited to text-based sentiment classification tasks.

### 2.2 Naïve Bayes Classifier (NBC)

Naïve Bayes Classifier is a probabilistic model based on Bayes' theorem, which assumes conditional independence among features (Rish, 2001). Despite its simplicity, NBC has shown strong performance in text classification tasks, particularly when the dataset is large and sparse. The formula used to determine the probability of a class given a document is:

$$P(C \mid X) = \frac{P(X \mid C) \cdot P(C)}{P(X)}$$

where $P(C \mid X)$ is the posterior probability of class $C$ given features $X$, $P(X \mid C)$ is the likelihood, $P(C)$ is the class prior, and $P(X)$ is the predictor prior. In sentiment analysis, Naïve Bayes models estimate the probability of a review being positive or negative based on the frequency of sentiment-bearing words. Loemongga et al. (2021) applied this model to Shopee customer reviews and achieved an accuracy of 85%. Despite its effectiveness, Naïve Bayes has limitations when dealing with non-linear separable data or when word dependencies are strong.

**2.3 Support Vector Machine (SVM)**

Support Vector Machine (SVM) is a supervised machine learning algorithm that constructs hyperplanes in a high-dimensional space to separate data into classes (Cortes & Vapnik, 1995). The algorithm's objective is to maximize the margin between the nearest data points (support vectors) of different classes, ensuring optimal separation and minimizing classification errors.

SVM has been widely applied in text classification due to its robustness in handling high-dimensional and sparse datasets. In sentiment analysis, SVM is particularly effective in capturing subtle distinctions in textual expression. Irma Surya Kumala Idris et al. (2023) reported that the SVM algorithm achieved 98% accuracy in classifying Shopee user sentiments, outperforming the Naïve Bayes Classifier

## 3. Method

### 3.1 Research Framework

This study adopts a mixed-method approach that combines computational sentiment analysis using machine learning algorithms with a quantitative behavioral survey for customer loyalty assessment. The overall framework is illustrated in *Figure 1*, which demonstrates the integrated stages of data acquisition, preprocessing, feature extraction, sentiment classification, and loyalty evaluation.
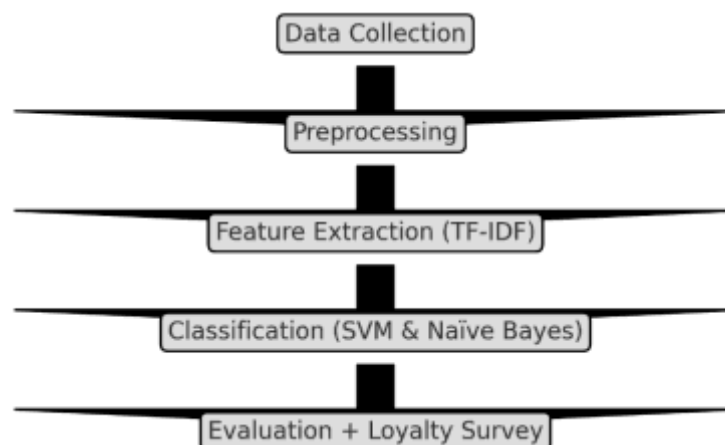


**Figure 1.** Methodology Flowchart

The integration of both computational and behavioral methods ensures that the results not only describe algorithmic performance but also provide actionable managerial insights. The design follows a sequential analytical process, beginning with data collection and text mining, followed by algorithm implementation and validation, and concluding with the interpretation of loyalty patterns.

### 3.2 Data Collection

The dataset for this study was obtained from Shopee's public customer review section, covering a wide range of products such as fashion, electronics, and household items. Reviews were extracted using Python-based web scraping tools between January and March 2025. A total of 3,000 textual reviews were collected, including both positive and negative comments. Non-textual or duplicate entries were removed during preprocessing. In addition, a customer loyalty survey was conducted among 30 respondents from Universitas Pembangunan Panca Budi, representing active Shopee users. The survey instrument employed Likert-scale questions covering loyalty dimensions such as repeat purchase intention, word-of-mouth behavior, and resistance to switching.

### 3.3 Data Preprocessing

Data preprocessing plays a crucial role in text mining and ensures that noisy, inconsistent, or redundant information is minimized. The preprocessing pipeline implemented in this study consists of several sequential steps:
  a. Case Folding: All text was converted to lowercase to ensure uniformity.
  b. Tokenization: Reviews were segmented into individual tokens (words).
  c. Stopword Removal: Commonly used words without semantic meaning (e.g., "the," "is," "and") were removed using an Indonesian stopword list.
  d. Slang Normalization: Informal words (e.g., "mantul," "brg," "bgs") were standardized into their formal equivalents.

e.  Stemming: Words were reduced to their base or root forms using the *Sastrawi* stemming algorithm for Bahasa Indonesia.

f.  Feature Extraction: Cleaned text data were transformed into numerical representations using the TF-IDF technique to measure the relative importance of each term.

This preprocessing pipeline enabled efficient vectorization of text data, facilitating the input into machine learning algorithms for classification.

## 3.4  Machine Learning Algorithms

Two machine learning algorithms were implemented: Naïve Bayes Classifier (NBC) and Support Vector Machine (SVM).

a.  Naïve Bayes Classifier
Naïve Bayes applies probabilistic modeling to classify sentiment based on the likelihood of word occurrences in each category. Its computational efficiency and simplicity make it suitable for large datasets. However, its assumption of word independence may limit accuracy when words are contextually related.

b.  Support Vector Machine
Support Vector Machine constructs an optimal separating hyperplane in high-dimensional space to maximize the margin between classes. This study used a linear kernel, as it performed best in preliminary trials. The model was trained using a stratified 80–20 split (training–testing) and validated with k-fold cross-validation (k=5).

## 4.  Results and Discussion

## 4.1 Sentiment Classification Results

The sentiment analysis was performed on 3,000 Shopee user reviews. After applying the preprocessing and TF-IDF vectorization steps, both algorithms Naïve Bayes Classifier (NBC) and Support Vector Machine (SVM) were trained and evaluated. The classification results showed that the SVM model outperformed NBC in nearly all evaluation metrics. Table 2 presents the comparative performance of both models.

**Table 2.** Performance Comparison between Naïve Bayes and SVM

| Algorithm | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Naïve Bayes Classifier | 85% | 0.86 | 0.84 | 0.85 |
| Support Vector Machine | 98% | 0.99 | 0.97 | 0.98 |

The SVM's superior performance can be attributed to its ability to handle high-dimensional feature spaces effectively, thus producing better generalization on unseen text data. Conversely, the Naïve Bayes model occasionally misclassified reviews that contained mixed or sarcastic sentiments due to its independence assumption.



**Positive Reviews:**

cheap, fast delivery, good quality, trusted seller

**Negative Reviews:**

late shipping, broken, slow app, not responding

**Figure 2.** Wordcloud Visualization of Positive and Negative Reviews

The wordcloud visualization shows that positive reviews prominently featured terms such as *"cheap,"* *"fast delivery,"* *"trusted seller,"* and *"good quality,"* indicating satisfaction with pricing, delivery, and reliability. Negative reviews, on the other hand, emphasized terms such as *"late shipping,"* *"broken,"* *"slow app,"* and *"bad service,"* reflecting logistical and service-related issues.

## 4.2 Confusion Matrix Analysis

To further examine model performance, a confusion matrix was generated for both algorithms. The confusion matrix provides insight into how many instances were correctly and incorrectly classified into positive and negative classes.
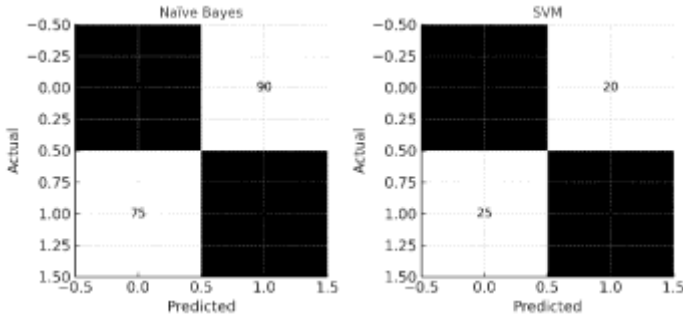


**Figure 4**. Confusion Matrix of SVM and Naïve Bayes Classifiers

The SVM model achieved a higher number of true positives (correctly classified positive reviews) and true negatives (correctly classified negative reviews) compared to Naïve Bayes. The misclassifications mainly occurred in cases where reviews contained mixed sentiments—such as "good product but slow delivery"—which caused ambiguity for probabilistic models. This confirms the robustness of the SVM algorithm for text classification tasks in Bahasa Indonesia, where contextual variations are frequent.

**4.3 Customer Loyalty Analysis**

The loyalty survey responses from 30 Shopee users were analyzed to determine behavioral loyalty patterns. The loyalty index was calculated based on the average of four dimensions—repeat purchase, cross-product loyalty, word-of-mouth, and resistance to competitors.

**Table 3**. Summary of Customer Loyalty Dimensions

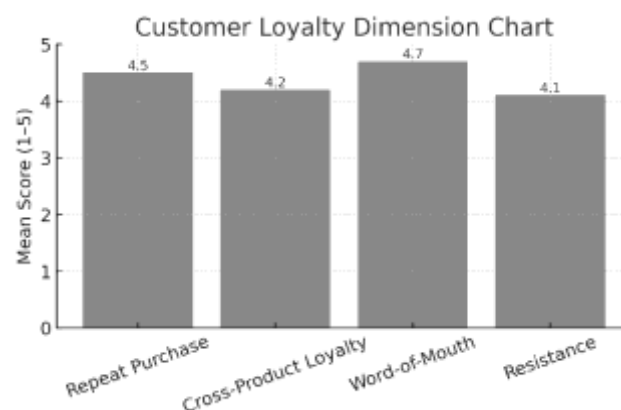| Loyalty Dimension | Mean Score | Interpretation |
|---|---|---|
| Repeat Purchase | 4.5 | High |
| Cross-Product Loyalty | 4.2 | Moderate to High |
| Positive Word-of-Mouth | 4.7 | Very High |
| Resistance to Competitors | 4.1 | Moderate to High |



**Figure 5.** Customer Loyalty Dimension Chart

The results show that Shopee users exhibit strong behavioral loyalty, particularly in the *positive word-of-mouth* dimension. Respondents frequently reported recommending Shopee to their friends or family, reflecting both emotional satisfaction and perceived trust. However, the *resistance to competitors* dimension was slightly lower, indicating that users may still explore alternative platforms for specific product categories.

## 5. Conclusions and Future Work

This study presents an integrated framework for analyzing customer sentiment and loyalty toward Shopee, one of Indonesia's leading e-commerce platforms. By combining machine learning–based sentiment analysis with survey-based loyalty assessment, the research provides both computational and behavioral perspectives on consumer engagement. The results demonstrated that the Support Vector Machine (SVM) algorithm achieved a higher accuracy rate (98%) than the Naïve Bayes Classifier (85%), confirming its superiority in handling high-dimensional textual data. The wordcloud visualization revealed that positive user sentiments were primarily associated with attributes such as affordability, fast delivery, and product reliability, while negative sentiments focused on delivery delays and customer service issues.

The customer loyalty survey further revealed that Shopee users exhibit strong behavioral loyalty, particularly in repeat purchase and word-of-mouth advocacy dimensions. However, the slightly lower score in the *resistance to competitors* dimension suggests that Shopee should continue to enhance personalization and service differentiation to maintain long-term loyalty. Overall, the correlation between positive sentiment and high loyalty confirms that emotional satisfaction significantly influences repeat purchase intentions and brand advocacy in digital marketplaces. These findings validate the importance of integrating computational intelligence with customer behavior analytics to achieve a holistic understanding of e-commerce user experience.

## 6. References

1. Irma Surya Kumala Idris, et al., "Sentiment Analysis of Shopee User Reviews Using Support Vector Machine Algorithm," *Jurnal Teknik Informatika*, 2023.
2. Loemongga, M. et al., "Comparative Analysis of Naïve Bayes and SVM for Sentiment Classification of E-Commerce Reviews," *Journal of Computer Science and Information Systems*, 2021.
3. Adhit Octavian, "The Effect of Customer Satisfaction on Shopee Loyalty in Indonesia," *Journal of Marketing Research*, 2018.
4. Liu, B. (2012). *Sentiment Analysis and Opinion Mining*. Morgan & Claypool Publishers.
5. Griffin, J. (2005). *Customer Loyalty: How to Earn It, How to Keep It*. Jossey-Bass.
6. Kotler, P. & Keller, K. L. (2012). *Marketing Management* (14th ed.). Pearson Education.

7.  Medhat, W., Hassan, A., & Korashy, H. (2014). "Sentiment Analysis Algorithms and Applications: A Survey," *Ain Shams Engineering Journal*.
8.  Pang, B., & Lee, L. (2008). "Opinion Mining and Sentiment Analysis," *Foundations and Trends in Information Retrieval*.