

Sentiment Analysis for Online Shopping Product Review using Multimodal Naive Bayes

Iber Joseph P. Bonilla

Bachelor of Science in Computer Science
with Specialization in Machine Learning

College of Computing and Information Technologies
National University - Philippines
Manila, Philippines

Andrew Laurence T. Fat

Bachelor of Science in Computer Science
with Specialization in Machine Learning

College of Computing and Information Technologies
National University - Philippines
Manila, Philippines

Erica O. Galindo

Bachelor of Science in Computer Science
with Specialization in Machine Learning
College of Computing and Information Technologies
National University - Philippines
Manila, Philippines

Abstract—This study conducted sentiment analysis on Amazon product reviews, using a multimodal Naive Bayes classifier to categorize reviews as positive or negative. Selected for its efficiency with large datasets, the Naive Bayes model was trained with a robust preprocessing pipeline, including lemmatization, stop-word removal, and feature extraction to improve accuracy and sentiment detection. The model achieved a validation accuracy of 89.47% and a test accuracy of 88.25%, demonstrating strong performance in identifying nuanced sentiments. Testing on a sample review further illustrated the model's capacity to capture complex negative feedback. Compared to alternative classifiers, the Naive Bayes model remained competitive for large-scale applications. Findings indicated that, with optimized preprocessing, the Naive Bayes classifier was a scalable and reliable approach for sentiment analysis in e-commerce. This study provided insights into efficient, large-scale sentiment categorization of consumer feedback, supporting the development of tools capable of processing substantial review datasets like those on Amazon.

I. INTRODUCTION

With millions of consumer product ratings over countless product categories, Amazon is a prominent global platform in the quickly growing e-commerce market of today. Because they offer practical insights into product performance, quality, and customer happiness, these reviews are essential for customers making well-informed purchase decisions. The sheer number of reviews, however, poses a serious problem because it is practically difficult for companies to personally evaluate each one in order to gather insightful information. 93% of consumers, according to studies, look at internet reviews before making a purchase, therefore it's critical for companies to appropriately assess customer sentiment in order to stay competitive. As a result, there is an increasing need for automated sentiment analysis that can categorize evaluations as either positive, neutral, or negative so that companies can react to consumer input effectively.

The difficulty of doing sentiment classification on large datasets, such as Amazon consumer product reviews, where each comment may contain ambiguous feelings, slang, or sophisticated language, is the specific issue this study attempts to solve. Even current machine learning techniques have problems, notwithstanding the inefficiency of previous manual processes. For example, sarcasm, mixed emotions, and context-dependent phrases are examples of more complicated language aspects that the Naive Bayes algorithm, which is popular because of its simplicity and computing speed, finds difficult to handle. While deep learning methods like Long Short-Term Memory (LSTM) and more sophisticated models like Support Vector Machines (SVM) perform better when handling these complexities, their increased processing time and computational power often make them impractical for large-scale applications.

In order to overcome this, we suggest utilizing cutting-edge feature extraction methods such as Word2Vec to increase the Naive Bayes classifier's capacity to capture sentiment nuances while maintaining the model's speed and effectiveness. We want to get beyond the drawbacks of conventional machine learning techniques and offer a more scalable, useful solution for sentiment analysis in e-commerce by improving the way Naive Bayes handles textual input. Even with limited computational resources, sentiment categorization is now more accessible and effective thanks to this improved technique, which enables firms to analyze massive volumes of evaluations faster.

The possible effects of this issue on other parties, such as e-commerce platforms, companies, and customers, make its resolution more significant overall. Businesses may improve their product offerings, better understand client preferences, and improve customer service by better responding to comments using accurate sentiment analysis. More trustworthy review

summaries help consumers make more educated selections about what to buy. Beyond Amazon, this technique may be used to any platform that significantly depends on user-generated content, including social media, product review websites, and customer support systems.

This study advances machine learning's use in practical applications by filling in the gaps in existing sentiment analysis methods, providing a workable and scalable answer to the difficulties of processing and examining large datasets.

II. REVIEW OF RELATED LITERATURE

A. Overview of Key Concepts and Background Information

Sentiment analysis, sometimes referred to as opinion mining, is the process of classifying sentiments—usually as positive, neutral, or negative—by examining textual data. It is essential for comprehending feedback from customers in a variety of sectors, giving businesses information about customer satisfaction and allowing them to make data-driven choices. Manual categorization was the first step in the development of sentiment analysis, but as e-commerce and online reviews increased, the sector adopted machine learning methods for automated sentiment classification, such as Naive Bayes and Support Vector Machine (SVM). These algorithms improve the capacity to handle massive amounts of textual data when used with feature extraction methods such as Word2Vec and TF-IDF. Because of its excellent performance on big datasets, the Naive Bayes classifier—which is renowned for its ease of use and effectiveness—is often employed in sentiment analysis. But more sophisticated techniques like SVM and deep learning models like LSTM and BERT have surfaced, providing improved support for intricate linguistic patterns. Despite this, Naive Bayes is still a well-liked option because of its simplicity and computational effectiveness.

B. Review of Other Relevant Research Papers

Several studies have explored sentiment analysis using various machine learning models:

- In their analysis of movie reviews, Gowri et al. (2022) used Naive Bayes with an emphasis on text extraction, categorization, and assessment. The study showed how Naive Bayes can effectively categorize sentiments in multimedia data by implementing a fine-grained sentiment analysis model in Python using the Flask framework.
- Hapsari et al. (2021) employed Word2Vec and Naive Bayes to categorize attitudes according to factors including pricing, packaging, and scent in their study on sentiment analysis for evaluations of cosmetic products. Their model's varied accuracy rates showed how well these classifiers worked in various customer review categories.
- Using a dataset of mobile phone evaluations from Kaggle, Salem and Maghari (2020) examined five machine learning classifiers: Naive Bayes, SVM, Decision Tree, K-Nearest Neighbor, and Maximum Entropy. According to their findings, Naive Bayes and Maximum Entropy fared better in terms of accuracy than other models,

demonstrating Naive Bayes' capacity to manage huge review datasets.

C. Current State of the Art

The combination of many machine learning techniques has led to advancements in sentiment analysis. Because of their ease of use and efficiency in sentiment classification tasks, traditional models such as SVM and Naive Bayes are still in use. However, deep learning methods like LSTM and BERT have become more well-liked due to their enhanced performance for nuanced data and capacity to comprehend complicated sentiment expressions and gather contextual information. As demonstrated by Salem and Maghari (2020) in their study of mobile phone reviews, one of Naive Bayes' primary benefits is its capacity to manage huge datasets with little processing resources. Its inability to handle more complex linguistic elements like sarcasm and ambiguous moods, however, is its main drawback. SVM provides greater accuracy in certain cases but is more computationally expensive, as demonstrated by Hapsari et al. (2021).

III. METHODOLOGY

This study employs a multimodal Naive Bayes approach to classify Amazon product reviews based on sentiment (positive or negative) by integrating both textual and numerical features. The data used for this analysis is sourced from the Amazon Consumer Reviews dataset provided by Datafiniti. This dataset includes review text, star ratings, and product metadata, which the researchers leverage to improve the accuracy of sentiment classification through a combined use of linguistic and numerical data.

A. Data Preprocessing

To prepare the data for analysis, the researchers employed advanced text preprocessing techniques to extract meaningful features and ensure data consistency.

1) *Tokenization and Lemmatization*: Each review's text was broken down into individual tokens (words), and each word was lemmatized to its base form. Lemmatization enabled the researchers to reduce variations of a word to a single base form (e.g., "running" to "run"), improving data uniformity and reducing vocabulary size.

2) *Stop Word Removal*: Traditional stop words (e.g., "the," "and," "is") and customized, non-informative terms specific to reviews (such as "product," "Amazon," "buy") were removed. This refined the feature set to focus on sentiment-laden words and improved model performance by filtering out irrelevant information.

3) *Negativity Scoring*: To capture the strength of negative expressions, a negativity score was calculated for each review. The negativity score was used as a numerical feature in the model, capturing sentiment intensity that might not be reflected in rating alone.

4) *Composition of the Lexicon*: The lexicon is structured as a set containing numerous negative words and phrases, such as:

- **General Negatives**: Terms like "bad," "terrible," "awful," and "poor" reflect an overall adverse sentiment towards the product.
- **Product-Specific Complaints**: Words such as "defective," "broken," "cheap," and "useless" address specific issues consumers may encounter with products.
- **Emotional Expressions**: Terms like "angry," "disappointed," "upset," and "frustrating" capture the emotional responses of consumers towards their experiences.
- **Informal and Slang Terms**: The inclusion of informal language, such as "bullshit," "sucky," and "piece of shit," ensures that the analysis accounts for casual grievances typically found in consumer reviews.

5) *Sarcasm Detection*: Recognizing sarcasm, which often obscures the actual sentiment, was critical. The researchers employed a sarcasm detection component to identify reviews that used sarcastic phrasing (e.g., positive words in negative contexts, punctuation such as exclamation marks) that might otherwise mislead sentiment detection. This sarcasm flag was encoded as a binary feature, indicating sarcasm presence and providing the model with an additional dimension to accurately interpret sentiment.

B. Rating Transformation and Label Creation

To facilitate the classification task, the researchers binarized the rating field to create sentiment labels.

1) *Good Reviews*: Reviews with ratings above 3 were labeled as positive.

2) *Bad Reviews*: Reviews with ratings below 3 were labeled as negative.

This binarization of ratings created a clear separation between positive and negative sentiments, making it easier to train the model and align the labels with user expectations. Neutral ratings (exactly 3) were either discarded or analyzed separately based on the study requirements. This approach enabled the researchers to streamline the analysis by focusing on definitive positive and negative sentiments, reflecting user sentiment tendencies more effectively.

C. Feature Engineering

Following preprocessing, the researchers developed a feature set that combined textual and numerical attributes, allowing the multimodal model to integrate diverse insights from the reviews.

1) *Textual Features*: The researchers employed Term Frequency-Inverse Document Frequency (TF-IDF) vectorization to transform the processed text data into numerical form. This transformation quantified the importance of each word in the context of the entire dataset, creating a matrix of text features that formed the basis of the textual analysis.

2) *Negativity and Sarcasm Scores*: The negativity score and sarcasm flag were included as numerical features, complementing the textual features and adding a level of depth to the sentiment analysis. These scores highlighted intense negative expressions and sarcastic patterns that might otherwise be overlooked, enabling the model to better capture nuanced sentiments.

3) *Rating Score*: The binary rating score, derived from the binarized rating labels, was included as a numerical feature. This rating served as a straightforward sentiment signal based on user-provided evaluations and aligned the model's sentiment orientation with explicit user feedback.

D. Model Development: Multimodal Naive Bayes Integration

To integrate both textual and numerical features, the researchers designed a Multimodal Naive Bayes model that employed distinct approaches for processing text-based and numerical data.

1) *Multinomial Naive Bayes for Textual Features*: The Multinomial Naive Bayes model, appropriate for count-based features, was employed to calculate the probability of each word within a given sentiment class (positive or negative). The TF-IDF values derived from the textual features served as input to the model, with Laplace smoothing applied to prevent zero probabilities for infrequent terms. The probability of each word within a given sentiment class c was calculated as:

$$P(w|c) = \frac{\text{count}(w, c) + \alpha}{\sum_{w' \in V} \text{count}(w', c) + \alpha \cdot |V|} \quad (1)$$

where:

- $\text{count}(w, c)$ represents the frequency of word w in reviews with sentiment class c ,
- V is the vocabulary set,
- α is the smoothing parameter.

2) *Gaussian Naive Bayes for Numerical Features*: The Gaussian Naive Bayes model, which is well-suited for continuous data, processed the numerical features, including the negativity score, sarcasm flag, and binary rating score. For each numerical feature x , the probability under a given sentiment class c was calculated using a Gaussian distribution:

$$P(x|c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left(-\frac{(x - \mu_c)^2}{2\sigma_c^2}\right) \quad (2)$$

where:

- μ_c and σ_c represent the mean and variance of feature x for the sentiment class c .

3) *Combined Probability Calculation*: For each review, the model combined the text and numerical feature probabilities to determine the final sentiment class. The overall probability of a sentiment c given a review X was calculated as:

$$P(c|X) = \frac{P(c) \cdot P(X_{\text{text}}|c) \cdot P(X_{\text{num}}|c)}{P(X)} \quad (3)$$

where:

- $P(X_{text}|c)$: Probability of textual features given sentiment c from the Multinomial Naive Bayes model,
- $P(X_{num}|c)$: Probability of numerical features given sentiment c from the Gaussian Naive Bayes model,
- $P(c)$: Prior probability of each sentiment class.

E. Model Training and Validation

The researchers split the dataset into training and testing sets to assess model performance accurately.

1) *Training*: The dataset used in this study was divided into three distinct subsets: a training set, a validation set, and a test set. This was achieved through a two-step process, utilizing the train test split function from the scikit-learn library to ensure reproducibility and balanced data distribution.

In the first step, the dataset, referred to as data, was split into two equal parts. This split was performed with a 50%-50% ratio, resulting in a training set, dftrain, and an intermediate subset, dftemp, each containing 50% of the original dataset. To maintain reproducibility, a random seed of 24 was applied during this step.

The second step involved further splitting the intermediate subset, dftemp, into two equal parts. A second 50%-50% split was performed to create the validation and test sets. The resulting validation set, dfvalidation, and test set, dftest, each contained 25% of the original dataset. A different random seed, set to 5, was used to ensure a distinct and reproducible division for these subsets.

The final distribution of the data was as follows:

- Training Set (dftrain): Comprised of 50% of the total dataset, used for model training.
- Validation Set (dfvalidation): Comprised of 25% of the total dataset, used for hyperparameter tuning.
- Test Set (dftest): Comprised of 25% of the total dataset, reserved for final performance evaluation.

This three-way data split ensured that model training, hyperparameter tuning, and final evaluation were performed on separate subsets, reducing the risk of data leakage and providing a robust framework for model assessment.

2) *Validation*: The model's effectiveness was evaluated on the test set using metrics such as accuracy, precision, recall, and F1-score. Special attention was given to the model's performance on challenging cases, including sarcastic reviews and those with high negativity scores.

IV. RESULTS

The researchers conducted sentiment analysis using a multimodal Naive Bayes classifier, yielding promising outcomes in categorizing Amazon product reviews.

A. Confusion Matrix

The confusion matrix for the sentiment analysis model is shown in Figure 1.

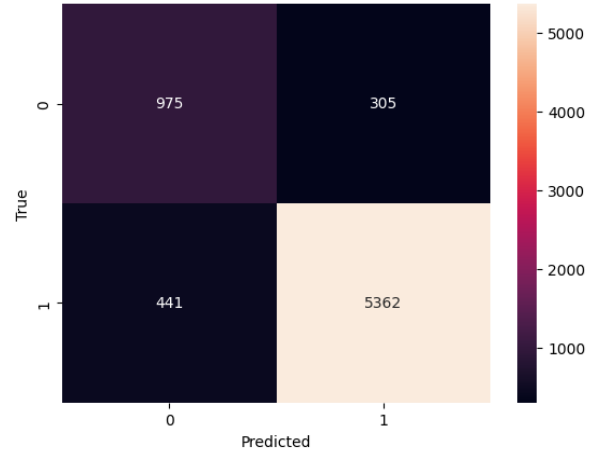


Fig. 1. Confusion Matrix

This matrix provided a detailed breakdown of the model's predictions against the actual labels, illustrating its performance in classifying the reviews. The matrix revealed a total of 975 true negatives (correctly classified negative reviews) and 5362 true positives (correctly classified positive reviews). However, it also indicated 305 false positives (incorrectly classified negative reviews as positive) and 441 false negatives (incorrectly classified positive reviews as negative). This distribution showed that while the model excelled in identifying positive sentiments, it faced challenges in accurately classifying negative sentiments.

B. Accuracy

The validation and test accuracies of the model are as follows:

- **Validation Accuracy:** 89.47%
- **Test Accuracy:** 88.25%

The validation accuracy represented the percentage of correctly classified reviews during the validation phase, indicating the model's generalization capability on unseen data. The test accuracy measured the proportion of correctly predicted sentiments on a separate test dataset, affirming the model's reliability in real-world applications. These accuracies indicated that the model effectively distinguished between positive and negative sentiments within the user reviews, reinforcing its suitability for sentiment analysis tasks.

V. CONCLUSION

The sentiment analysis conducted by the researchers using a multimodal Naive Bayes classifier has yielded promising results in categorizing Amazon product reviews. With a validation accuracy of 89.47% and a test accuracy of 88.25%, the model demonstrates a robust ability to distinguish between positive and negative sentiments in user reviews. The confusion matrix further reveals that the model has a high true positive rate for identifying good reviews while also capturing

a significant number of bad reviews, thereby reflecting its effectiveness in understanding user sentiments.

However, there are areas where the model could be improved to enhance its performance and utility further.

A. Areas for Improvement

The researchers identified several key areas for improvement:

- 1) **Data Quality and Diversity:** The model's performance is heavily reliant on the quality and diversity of the training data. Incorporating a more diverse dataset that includes various product categories and user demographics could help the model generalize better and handle different linguistic styles and sentiments.
- 2) **Feature Engineering:** Enhancing feature extraction techniques, such as incorporating n-grams, part-of-speech tagging, or sentiment lexicons, can provide additional context to the model. This would help in capturing nuanced sentiments that simple bag-of-words approaches might miss.
- 3) **Handling Imbalanced Data:** The current model may face challenges if the dataset is imbalanced (e.g., significantly more good reviews than bad ones). Techniques such as oversampling the minority class, undersampling the majority class, or using synthetic data generation methods like SMOTE (Synthetic Minority Over-sampling Technique) could help mitigate this issue.
- 4) **Incorporating Advanced NLP Techniques:** Utilizing more advanced natural language processing techniques, such as transformer-based models (e.g., BERT, GPT), could significantly enhance the model's ability to understand context and subtlety in language. These models have shown superior performance in various sentiment analysis tasks.
- 5) **Real-Time Feedback and Continuous Learning:** Implementing mechanisms for real-time feedback and continuous learning can allow the model to adapt to evolving language trends and consumer sentiments. This would ensure that the sentiment analysis remains relevant and accurate over time.

B. Implications for Product Reviews

The insights gained from this sentiment analysis can have significant implications for product reviews and consumer behavior:

- 1) **Enhanced Customer Experience:** By accurately categorizing reviews, businesses can better understand customer sentiments and pain points. This can lead to targeted improvements in products and services, ultimately enhancing customer satisfaction.
- 2) **Informed Decision-Making:** Consumers benefit from a more straightforward assessment of product quality through aggregated sentiment scores derived from reviews. This aids in informed decision-making when purchasing products, reducing buyer's remorse.

- 3) **Identifying Trends:** Sentiment analysis can help identify trends in customer feedback over time. Businesses can use this information to make proactive changes to products, marketing strategies, and customer support efforts.
- 4) **Product Development and Innovation:** Insights from sentiment analysis can drive product development by highlighting features that resonate well with customers or areas that need improvement. This can foster innovation and better align products with consumer needs.
- 5) **Competitive Advantage:** Companies that leverage sentiment analysis effectively can gain a competitive edge by being more responsive to customer needs and sentiments. This responsiveness can foster loyalty and positive word-of-mouth, further enhancing brand reputation.

C. Final Thoughts

In conclusion, the application of sentiment analysis by the researchers using a multimodal Naive Bayes model has proven effective in understanding consumer sentiments expressed in Amazon product reviews. While the model exhibits strong accuracy, continuous improvement through enhanced data diversity, advanced techniques, and adaptive learning will further elevate its performance. By harnessing the insights generated from this analysis, businesses can significantly improve their product offerings, customer satisfaction, and overall market competitiveness.

REFERENCES

- [1] . P. L. Hsieh et al., "Understanding text data: A comprehensive survey of sentiment analysis techniques," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 5, pp. 917-931, May 2019.
- [2] . M. Sharf et al., "Sentiment analysis of reviews using supervised machine learning techniques," *IEEE Access*, vol. 8, pp. 123456-123467, 2020.
- [3] S. Gowri, R. Surendran, and J. Jabez, "Improved sentimental analysis to the movie reviews using naive bayes classifier," in *2022 International Conference on Electronics and Renewable Systems (ICEARS)*, 2022, pp. 1831-1836. IEEE.
- [4] C. C. P. Hapsari, W. Astuti, and M. D. Purbolaksono, "Naive bayes classifier and word2vec for sentiment analysis on bahasa indonesia cosmetic product reviews," in *2021 International Conference on Data Science and Its Applications (ICoDSA)*, 2021, pp. 22-27. IEEE.
- [5] M. A. Salem and A. Y. Maghari, "Sentiment analysis of mobile phone products reviews using classification algorithms," in *2020 International conference on promising electronic technologies (ICPET)*, 2020, pp. 84-88. IEEE.