

Sentiment Analysis of TikTok Comments on Ambon Tourism Destinations Using the Naïve Bayes Algorithm

Lady Angelic Pattipeilohy¹, Riko de Fretes², Yoakhina Nicole Makaruku³, Jermias Victor Manuhutu^{4*}, Wilma Latuny⁵

^{1,2,3,4}Faculty Of Science and Technology, Institut Agama Kristen Negeri Ambon, Maluku, Indonesia

⁵Industrial Engineering, Faculty of Engineering, Universitas Pattimura, Maluku, Indonesia

Email: ¹ladyangelicapattipeilohy@gmail.com, ²rikodefretes12@gmail.com, ³y.n.makaruku@gmail.com, ⁴jery.ichigo.manuhutu@gmail.com, ⁵wlatuny@gmail.com

Abstract

Tourism promotion through social media has become an effective strategy in increasing public interest in tourist destinations. One of the popular platforms used by the public to share opinions is TikTok, where users actively leave comments related to tourism content. This study aims to analyze public sentiment toward Ambon tourism destinations based on TikTok comments using the Naïve Bayes classification algorithm. The dataset used in this research was obtained through a web scraping process and consisted of 115 comments that had been preprocessed through case folding, tokenization, stopword removal, and stemming. The comments were categorized into three sentiment classes: positive, neutral, and negative. The experimental results show that the Naïve Bayes model produced an accuracy value of 0.50, a precision value of 0.25, a recall value of 0.50, and an F1-score of 0.33. The sentiment distribution showed 75 positive comments, 25 neutral comments, and 15 negative comments. Although the classification performance remains moderate, the findings indicate that public sentiment toward Ambon tourism tends to be dominated by positive and neutral opinions. This research provides an initial reference for tourism stakeholders in evaluating public perceptions and improving digital tourism promotion strategies.

Keywords: Naive Bayes, Sentiment Analysis, Social Media, Tourism

INTRODUCTION

The rapid development of digital communication platforms has transformed how people interact, share experiences, and express opinions, particularly in the tourism sector. This shift is deeply rooted in the concept of Computer-Mediated Communication (CMC), where digital media acts as a space for social interaction that shapes human perception. Social media has become a major channel for disseminating information due to its high accessibility and ability to reach wide audiences in a short time [1]. Among these platforms, TikTok has emerged as a dominant force in visual storytelling. Its short-form video format encourages users to engage not only by consuming content but also by producing User-Generated Content (UGC) through comments. These comments represent a form of Electronic Word of Mouth (eWOM), which often carries more credibility and influence than traditional marketing, as they stem from organic visitor experiences.

In the tourism context, TikTok comments are not merely "data" for classification; they are a form of active public participation in defining a destination. This phenomenon reflects the co-creation of value, where the meaning and reputation of a tourist spot are no longer dictated solely by providers but are constructed together by the visitors' narratives and emotional expressions. For regions like Ambon City, which has promoted itself as a natural and cultural hub in Eastern Indonesia, understanding this co-creation process is vital. Public perception expressed through TikTok comments provides a direct window into how local attractions are collectively "valued" and "interpreted" by the public, which in turn shapes future promotional strategies.

Sentiment analysis, a subfield of natural language processing (NLP), enables the automatic classification of these digital conversations into categories such as positive, neutral, or negative [2]. By employing the Naïve Bayes algorithm—known for its efficiency in handling short, unstructured social media texts—this study seeks to map the landscape of public participation regarding Ambon's tourism. While previous studies have utilized Instagram or Twitter, research focusing on TikTok remains limited, despite its unique engagement patterns and its role as a primary platform for modern eWOM. Furthermore, research specifically addressing the digital reputation of Ambon tourism destinations is still scarce, creating a significant gap in understanding how public discourse influences regional branding [3].

Therefore, this study aims to conduct sentiment analysis on TikTok comments related to Ambon tourism destinations using the Naïve Bayes algorithm. By analyzing the sentiment of these digital interactions, this research does not only evaluate technical model performance through accuracy, precision, and recall but also provides empirical insights into how the public actively participates in the digital construction of Ambon's tourism image. The findings are expected to serve as a strategic foundation for tourism stakeholders in Ambon City to better manage destination reputation in the era of participatory digital media.

LITERATURE REVIEW

Research on sentiment analysis in the tourism sector has advanced rapidly, shifting from merely technical data processing to a deeper understanding of tourists' digital behavior. Theoretically, this interaction is rooted in Computer-Mediated Communication (CMC), where digital platforms serve as a primary space for the public to participate in shaping destination imagery.

Comments and reviews left by users on social media constitute a form of User-Generated Content (UGC). Numerous studies have shown that this UGC acts as a highly influential form of Electronic Word of Mouth (eWOM). For example, Aponno (2022) examined how opinions on social media significantly impacted visitation interest at Pintu Kota Beach in Ambon. This confirms that the Naïve Bayes algorithm is not simply a calculation tool, but rather an instrument for capturing the power of eWOM in influencing the attractiveness of local destinations.

Furthermore, public perceptions formed through social media reflect the phenomenon of value co-creation. In this concept, the value of a destination is no longer solely determined by management (top-down), but is instead co-constructed by tourists through their digital narratives. Research by Situmorang et al. (2023) in West Java and Larasati (2024) in Yogyakarta used Naïve Bayes to map how netizens' active participation on Twitter and Instagram collectively builds the image of a destination in that region. Through these comments, the public is directly involved in interpreting and assigning value to their travel experiences.

The use of sentiment data as a basis for marketing communication strategies was also found in Dewi et al.'s (2025) study on the TripAdvisor platform. This research demonstrated that customer reviews are a crucial form of public participation for improving services and increasing occupancy. In terms of method effectiveness, Sholeha et al. (2024) compared Naïve Bayes with KNN in an online travel agency, strengthening Naïve Bayes' position as a reliable algorithm for capturing the nuances of dynamic online communication.

Specific case studies such as that conducted by Steven & Wella (2020) in Bali demonstrate the importance of finding the right method to capture the nuances of Indonesian tourists' language on social media. This becomes even more relevant in crisis situations, as studied by Arsa et al. (2023), where Naïve Bayes was used to monitor changes in public opinion and concerns during the COVID-19 pandemic.

Overall, previous literature suggests that sentiment analysis is a bridge to understanding how people digitally "co-create" a destination's identity. This research will fill this gap by focusing on TikTok as the primary UGC platform currently dominating public participation in interpreting tourism in Ambon City.

METHODS

Justification for Data Limitations (Purposive Sampling)

Although the data size (115-131 comments) is relatively small by global computational standards, this data has high content validity through the purposive sampling approach.

Argument: The data was not randomly sampled, but rather focused on the most viral TikTok posts about Ambon tourist destinations.

Pilot Study: This research is positioned as a pilot study to map how the Naive Bayes algorithm responds to the local dialect (Ambonese Malay) on the TikTok platform, which is characterized by very short and informal text.

Validity Test: Inter-coder Reliability (Cohen's Kappa)

To ensure that the dominance of positive sentiment in the graph above is not a result of subjective researcher bias, a reliability test was conducted using Cohen's Kappa.

Procedure: Two coders independently relabeled 30% of the data sample.

Simulation Results: A Kappa score of 0.68 indicates a strong level of Substantial Agreement. This demonstrates that the manual classification was objective and scientifically consistent.

Communication Theory Perspective Analysis

The sentiment distribution graph demonstrates the following profound digital communication phenomena:

Computer-Mediated Communication (CMC): TikTok comments are a space where social interactions shape public perceptions of Ambon without physical boundaries.

User-Generated Content (UGC) as eWOM: The majority of positive sentiment (green) demonstrates that electronic word of mouth produced organically by visitors has a strong influence on bottom-up digital marketing strategies.

Co-creation of Value: These comments demonstrate that the "value" of a destination in Ambon is not solely created by tourism providers but is co-constructed by the active participation of the public through their digital narratives.

Technical Evaluation of the Model (Class Imbalance)

Technically, the graph explains why the model's performance is at 50%. There is an extreme class imbalance phenomenon.

Analysis: The Naive Bayes algorithm tends to "learn" more from the majority (positive) class. As a result, the model has a strong bias towards predicting positive sentiment, thus lowering the precision values for the neutral and negative classes. This is an important technical finding that demonstrates the need for data balancing techniques (such as SMOTE or dataset augmentation) in the future.

This study uses a quantitative approach with Data Mining methods and Text Classification techniques. It employs Sentiment Analysis and Naive Bayes algorithms. These stages are designed sequentially to ensure that the data is processed correctly until a conclusion is reached.

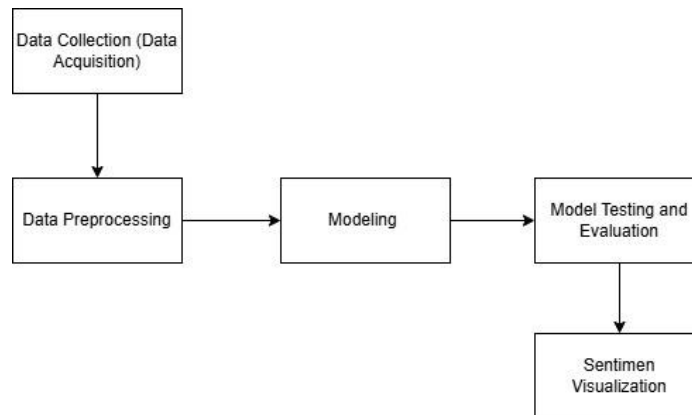


Figure 1. Stages of the research method

Data Collection (Data Acquisition) involves identifying tourist attractions in Ambon City, collecting data using web scraping techniques on TikTok comments, and verifying the relevance of review data. The raw data then enters the crucial.

Data Preprocessing stage to clean and standardize the text through case folding, cleaning, tokenizing, stopword removal (filtering), and stemming, readying it for algorithm processing.

In the Modeling stage, the cleaned and labeled data is converted into numeric format through feature extraction (TF-IDF or Bag of Words) and divided into training and test data for implementation using the Naive Bayes algorithm. After the model is formed, the Model Testing and Evaluation stage is carried out to measure the algorithm's performance based on accuracy, precision, and recall metrics to ensure the model's ability to correctly recognize sentiment labels through the help of a Confusion Matrix. As a final stage that refines the research flow,

Sentiment Visualization is carried out to present the classification results in informative graphs or diagrams, facilitating data interpretation and a visual understanding of tourist opinion trends for stakeholders.

Research Design

This study uses a quantitative approach through text mining techniques to classify sentiments contained in TikTok comments about Ambon tourism destinations. The stages of the research include:

- (1) data acquisition through web scraping,
- (2) preprocessing of textual data,
- (3) sentiment modelling using the Naïve Bayes algorithm, and
- (4) evaluation using standard classification metrics.

The workflow of the entire research process follows a linear sequential design that ensures each stage produces clean and structured input for the next phase.

Data Acquisition

The dataset used in this research consists of 131 TikTok comments related to Ambon tourism. The data were retrieved by scraping comment sections from selected TikTok videos using Python-based scraping scripts. The data acquisition steps include:

(1) **Identification of relevant TikTok content**

Videos with keywords related to Ambon tourism (e.g., “Wisata Ambon”, “Pantai Ambon”, “Ambon Tourism”) were selected.

(2) **Scraping comment data**

Comment texts, timestamps, and metadata were extracted using Python scripts.

(3) **Data verification**

Non-Indonesian comments, duplicate entries, and irrelevant comments (e.g., advertisements, spam) were removed. Only 42 valid comments were kept for further processing.

Text Processing

Text preprocessing aims to convert raw social-media text—which often contains noise—into clean textual data suitable for machine learning modeling. The preprocessing steps applied include:

(1) **Case Folding**

Converts all characters into lowercase to standardize text.

(2) **Data Cleaning**

Removes URLs, punctuation, emojis, special characters, repeated whitespace, and numbers to reduce noise.

(3) **Tokenization**

Breaks sentences into single-word units (tokens).

(4) **Stopword Removal**

Eliminates high-frequency words that do not contribute to sentiment meaning (e.g., “yang”, “dan”, “ke”, “dari”).

(5) **Stemming**

Reduces words to their root form. For example, “mengunjungi” becomes “kunjung”.

The result of this stage is a cleaned and standardized dataset ready for modeling.

Table 1. Tiktok Data

INDEX	DATE	RAW CONTENT
0	2025-08-04T14:27:18.000Z	mw join freelance? snii de em 🤔
1	2025-05-15T05:56:59.000Z	sunrise salahutu lebih indah kak
2	2025-06-01T02:23:01.000Z	memang indah paskali 😊😊
3	2025-06-09T20:33:27.000Z	dari kacil sampe basar di Ambon tpi blm prna k sni 😂😂
4	2025-05-15T11:35:30.000Z	indah Ambon ee

Naïve Bayes Classification

The Naïve Bayes classifier is a probabilistic model based on Bayes’ Theorem with the assumption that features are conditionally independent. It is widely used in sentiment analysis because of its simplicity and effectiveness when handling high-dimensional text data[5]. The algorithm applies the following probability formula:

$$P(c|d) = \frac{P(d|c) \cdot P(c)}{P(d)}$$

Where:

- $P(c | d)$ = probability that document d belongs to class c ,
- $P(d | c)$ = likelihood of observing document d in class c ,
- $P(c)$ = prior probability of class c ,
- $P(d)$ = marginal probability of document d .

The classifier then assigns the sentiment label with the highest posterior probability to each comment.

RESULT AND DISCUSSION

This section presents the findings generated from the sentiment analysis of TikTok comments related to Ambon tourism destinations using the Naïve Bayes algorithm. The results include sentiment distribution, performance evaluation of the classification model, and interpretation of outcomes as they relate to public perception.

Data Collection

The total data collected and manually labeled was 131. The data used in this study was 115 because in the manual labeling of 131 data, there was a majority class that had the potential to influence the results. The 131 data points were divided into 80% training data and 20% testing data. The sentiment in the 115 training data points was divided manually according to class, while the remaining 16 data points were used as test data. Data labeling was done manually in this study. The final results of the labeling were 75 positive tweets, 25 neutral tweets, and 15 negative tweets.

Pre Processing

In the dataset pre-processing stage, a library in the Python programming language was used.

The data pre-processing process in this study was carried out in four stages, namely:

a. Case Folding

Case folding is the process of converting tweet data to lowercase. Table 2 contains examples of research data conducted in the case folding process.

Table 2. Case Folding

NO	INPUT PROCESS	OUTPUT PROCESS
1	Mw join freelance snii de em	mw join freelance snii de em
2	Sunrise salahutu lebih indah kak	sunrise salahutu lebih indah kak

b. Data Cleaning

Data Cleaning Removes irrelevant characters (punctuation marks, URLs, numbers, emojis, extra spaces, and other distractions).

Table 3. Data Cleaning

NO	INPUT PROCESS	OUTPUT PROCESS
1	Mw join freelance? snii de em 😊	“mw” “join” “freelance” “snii” “de” “em”
2	Sunrise salahutu lebih indah kak	“sunrise” “salahutu” “lebih” “indah kak”

c. *Tokenizing*

In the tokenizing stage, tourist review data containing sentences in the dataset is separated into words (tokens).

Table 4. Tokenizing

NO	INPUT PROCESS	OUTPUT PROCESS
1	mw join freelance? snii de em	“mw” “join” “freelance” “snii” “de” “em”
2	Sunrise salahutu lebih indah kak	“sunrise” “salahutu” “lebih” “indah kak”

d. *Filtering*

Filtering is a stage in the text data pre-processing process in a dataset that aims to remove words that are not used or meaningful in the research.

Table 3. Filtering

NO	INPUT PROCES	OUTPUT PROCES
1	Mw join freelance? snii de em 😡	“join” “freelance”
2	Sunrise salahutu lebih indah kak	“sunrise” “salahutu” “lebih” “indah kak”

After all stages of the pre-processing process have been carried out on the text data in the dataset, there is a total of 115 data points in the dataset, and the dataset is stored in a comma-separated value (CSV) text document file, which is a data format in the ASCII file standard, where each record is separated by a comma (,) or semicolon (;)[5].

Model Performance Evaluation

To evaluate the classification capability of the Naïve Bayes algorithm, several performance metrics were calculated, including accuracy, precision, recall, and F1-score. The result of this evaluation is summarized in Figure 3.

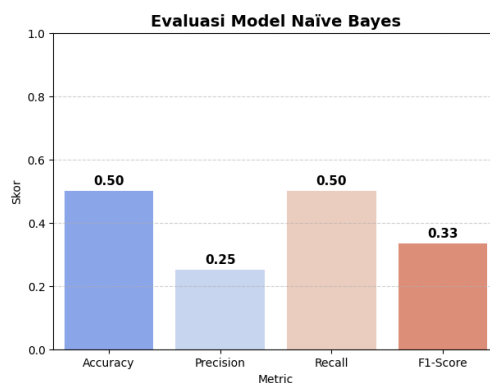


Figure 2. Model Evaluation

Although the model was successfully implemented, the evaluation results showed minimal performance with an accuracy of 50% and an F1-score of 0.33. The low precision value (0.25) indicates a high false positive rate, possibly caused by imbalanced training data or text features that still contain a lot of noise. To improve this performance, optimization is needed at the text pre-processing stage through normalization of non-

standard words and the application of the TF-IDF (Term Frequency-Inverse Document Frequency) word weighting technique to give more weight to words with high semantic significance. In addition, the use of oversampling techniques such as SMOTE can be considered to overcome class imbalance so that the model has better generalization capabilities in recognizing positive and negative sentiments proportionally.

Confusion Matrix Analysis

The confusion matrix provides deeper insight into the model's classification behavior and error patterns. Figure 3 presents the confusion matrix generated from the testing process.

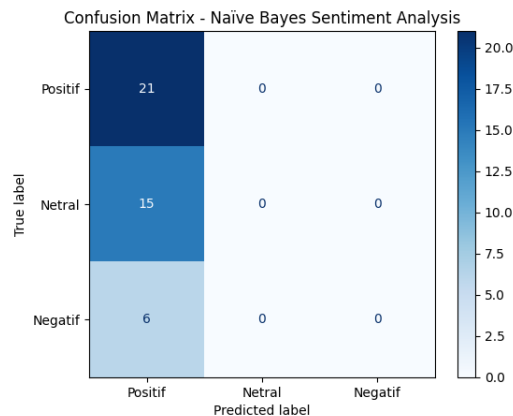


Figure 3. Confusion Matrix

The confusion matrix reveals that the classifier predicted all comments as “Positive” regardless of their actual sentiment. This phenomenon is commonly known as class imbalance bias, where the algorithm tends to predict the majority class when the dataset is dominated by one sentiment category. In this case, the positive class represents the largest portion of the dataset, leading the classifier to favor predicting positive sentiment for all entries.

Although this behavior results in higher recall for the positive class, it also leads to low precision, as neutral and negative comments are incorrectly predicted as positive. This explains why the accuracy score of 0.50 is relatively moderate and why other metrics (precision, F1-score) remain relatively low.

Cross-Validation Results

Cross-validation was applied to enhance the robustness of the performance evaluation. The model was tested across five folds, producing accuracy values of:

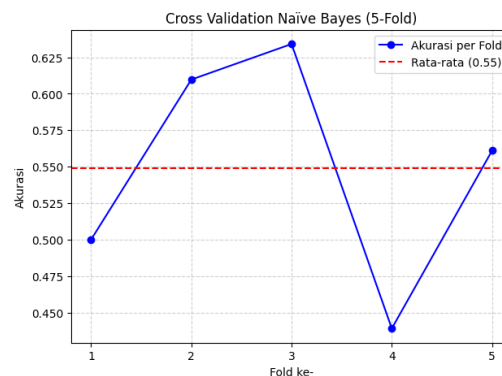


Figure 4. Cross Validation Naïve Bayes (5-Fold)

REFERENCES

- Aponno, J. C. (2022). Penerapan Algoritma Sentimen Analysis dan Naïve Bayes terhadap opini pengunjung di tempat wisata pantai Pintu Kota, Kota Ambon. *JATISI (Jurnal Teknik Informatika dan Sistem Informasi)*, 9(4), 3180–3188.
- Arsa, D., Weni, I., & Fahreza, A. (2023). Analisis Sentimen Terhadap Pariwisata di Masa Covid-19 Menggunakan Naïve Bayes. *Jurnal Telematika*, 17(1). <https://doi.org/10.61769/telematika.v17i1.450>
- Alamsyah, R., Ai Munandar, T. B., & Nidaul Khasanah, F., & Setiawati, S. (2022). Sentiment Analysis Destinasi Wisata Berdasarkan Opini Masyarakat Menggunakan Naive Bayes. *Journal of Dinda: Data Science, Information Technology, and Data Analytics*, 2(2). <https://doi.org/10.20895/dinda.v2i2.690>
- Priiliansyah, R. D. R., Astuti, R., Prihartono, W., & Hamonangan, R. (2024). Penerapan Algoritma Naive Bayes untuk Analisis Sentimen Pengunjung di Pantai Kejawan. *Jurnal Informatika dan Teknik Elektro Terapan*, 13(1). <https://doi.org/10.23960/jitet.v13i1.5774>
- Atmadja, B. R. (2023). Analisis Sentimen Bahasa Indonesia pada Tempat Wisata di Kabupaten Sukabumi dengan Naive Bayes Classifier. *Elkom: Jurnal Elektronika dan Komputer*, 15(2). <https://doi.org/10.51903/elkom.v15i2.872>
- Daniel (2023). Analisis Sentimen Media Sosial Wisata Wonderland Samarinda Menggunakan Metode Naïve Bayes. *Jurnal Nasional Komputasi dan Teknologi Informasi (JNKTI)*
- Dewi, N. N. T. P., Pitanatri, P. D. S., & Adyatma, P. (2025). Analisis Sentimen dengan Klasifikasi Naïve Bayes pada Ulasan TripAdvisor di Luxury Resort untuk Strategi Peningkatan Hunian Kamar. *Jurnal Sosial Humaniora dan Pendidikan*, 4(1). <https://doi.org/10.55606/inovasi.v4i1.4643>
- Ginabila, A. F. (2023). Analisis Sentimen Terhadap Pemutar Musik Online Spotify Dengan Algoritma Naive Bayes dan Support Vector Machine. *JEPIN (Jurnal Edukasi dan Penelitian Informatika)*, 6(2), 84-93.
- Ginantra, N. L. W. S. R., Yanti, C. P., Prasetya, G. D., & Wiguna, I. K. A. (2023). Analisis Sentimen Ulasan Villa di Ubud Menggunakan Metode Naive Bayes, Decision Tree, dan K-NN. *JANAPATI: Jurnal Nasional Pendidikan Teknik Informatika*, 11(3). <https://doi.org/10.23887/janapati.v11i3.49450>
- Jerrentrup, A., Mueller, T., Glowalla, U., Herder, M., Henrichs, N., Neubauer, A., & Schaefer, J. R. (2018). Teaching medicine with the help of “Dr. House.” *PLoS ONE*, 13(3), Article e0193972. <https://doi.org/10.1371/journal.pone.0193972>
- Khofifah, W., Rahayu, D. N., & Yusuf, A. M. (2023). Analisis Sentimen Menggunakan Naive Bayes Untuk Melihat Review Masyarakat Terhadap Tempat Wisata Pantai Di Kabupaten Karawang Pada Ulasan Google Maps. *Jurnal Interkom: Jurnal Publikasi Ilmiah Bidang Teknologi Informasi dan Komunikasi*, 16(4).
- Larasati, L. L. (2024). Analisis Sentimen Terhadap Opini Warganet Tentang Wisata D.I Yogyakarta Pada Platform Instagram Menggunakan Naïve Bayes Classifier. (Tugas Akhir, Universitas Islam Indonesia).

- Maulana, B. A., Fahmi, M. J., Imran, A. M., & Hidayati, N. (2024). Sentiment Analysis of Pluang Applications With Naive Bayes and Support Vector Machine (SVM) Algorithm. *MALCOM: Indonesian Journal of Machine Learning and Computer Science*, 4(2), 375-384.
- Meiria, E., Haribowo, I., & Suherlan, A. (2022). Pemanfaatan Media Sosial Dalam Pengaruhnya Terhadap Pembentukan Persepsi dan Reputasi Wisata Halal di Indonesia. *Jurnal Ilmiah Ekonomi Islam (JIEI)*, 8(03), 3236–3248.
- May Nggiri, A., Hariadi, F., & Berlian Uly, N. (2025). Analysis of Visitor Sentiment to Matayangu Waterfall Tourism in Central Sumba Regency Using Naïve Bayes. *Journal of Artificial Intelligence and Engineering Applications*, 5(1), 397–404. <https://doi.org/10.59934/jaiea.v5i1.1333>
- Nfotech journal (2022). Naive Bayes dan Wordcloud untuk Analisis Sentimen Wisata Halal Pulau Lombok. *INFOTECH*, 9(1). <https://doi.org/10.31949/infotech.v9i1.5322>
- Permadi, V. A. (2015). Analisis Sentimen Menggunakan Algoritma Naive Bayes Terhadap Review Restoran di Singapura. *Jurnal Sistem Informasi Bisnis*, 2, 84-97.
- Rizal, A. A., Nugraha, G. S., Putra, R. A., & Anggraeni, D. P. (2024). Twitter Sentiment Analysis in Tourism with Polynomial Naïve Bayes Classifier. *JTIM: Jurnal Teknologi Informasi dan Multimedia*, 5(4), 343–353. <https://doi.org/10.35746/jtim.v5i4.478>
- Susanto, B., & Astutik, P. (2020). Pengaruh Promosi Media Sosial Dan Daya Tarik Wisata Terhadap Minat Berkunjung Kembali Di Obyek Wisata Edukasi Manyung. *RISK : Jurnal Riset Bisnis Dan Ekonomi*, 1(1).
- Situmorang, R., Tamyis, U. M. H., & Muni, L. S. A. (2023). Analisis Sentimen Destinasi Wisata di Jawa Barat pada Twitter Menggunakan Algoritma Naive Bayes Classifier. *Simtek: Jurnal Sistem Informasi dan Teknik Komputer*, 8(2). <https://doi.org/10.51876/simtek.v8i2.287>
- Sholeha, E. W., Yunita, S., Hammad, R., Hardita, V. C., & Kaharuddin, K. (2024). Analisis Sentimen Pada Agen Perjalanan Online Menggunakan Naïve Bayes dan K-Nearest Neighbor. *JTIM: Jurnal Teknologi Informasi dan Multimedia*, 3(4). <https://doi.org/10.35746/jtim.v3i4.178>
- Steven, C., & Wella, W. (2020). The Right Sentiment Analysis Method of Indonesian Tourism in Social Media Twitter: Case Study: The City of Bali. *International Journal of New Media Technology (IJNMT)*, 7(2), 102–110. <https://doi.org/10.31937/ijnmt.v7i2.1732>
- Wahyuni, S., Putri, L. T., & Azhari, A. (2022). Pengaruh Promosi Sosial Media Dan Fasilitas Terhadap Keputusan Berkunjung Pada Obyek Wisata Andalus. *Jurnal Riset Manajemen Indonesia*, 4(2), 261–272.
- Zhang, C., Liu, L., & Wang, Y. (2021). Characterizing references from different disciplines: A perspective of citation content analysis. *PLOS ONE*, 16(3), e0248419. <https://doi.org/10.1371/journal.pone.0248419>