



Identification Digital Tourist Preferences using Sentiment Analysis: dealing in post-pandemic covid-19

Lutfia Septiningrum¹, Riska Yanu fa'rifah², Fandi Achmad³.

¹Department Public Administration, Universitas Negeri Yogyakarta, DI Yogyakarta, Indonesia

²Department System Information, Telkom University, Bandung, Indonesia

³Department Industrial Engineering, Telkom University, Bandung, Indonesia

ARTICLE INFO

Article history:

Received 10 October 2022

Received in revised form 21
November 2022

Accepted 4 Desember 2022

ABSTRACT

This aims to determine the post-pandemic impact on tourism in Rembang based on the results of community study using social media review. To achieve this goal, this study conducted a sentiment analysis of public comments regarding travel preferences after the Covid-19 pandemic using the SVM classification method for digital tourism. The results accuracy value of 83% and an AUC value of 82.3% show level of confidence on community comments using several digital platforms has a good category. The classification of community preferences is stated in two groups, first is the people statement post-pandemic digital tourism must be developed, and the second statement that they have decided to come directly to tourist attractions after Covid-19 in Rembang regency. Research suggestions for local governments to continue to develop digital tourism and improve real tourism facilities, where the pathway will increase tourist visits and the economy of the community around tourist attractions.

Keyword:

Covid-19, Digital Tourism,
Sentiment Analysis, SVM

INTRODUCTION

The COVID-19 pandemic is a nightmare for all sectors, including tourism. Tourism is one of the highest contributors to economic growth Indonesia. Based on (Statistics of Rembang Regency, 2022) and (Statistics Indonesia, n.d.) the country's foreign exchange income from the tourism sector from 2013 to 2022 decrease drastically in 2020 and 2021, as shown in Figure 1.

¹lutfiaseptiningrum@uny.ac.id

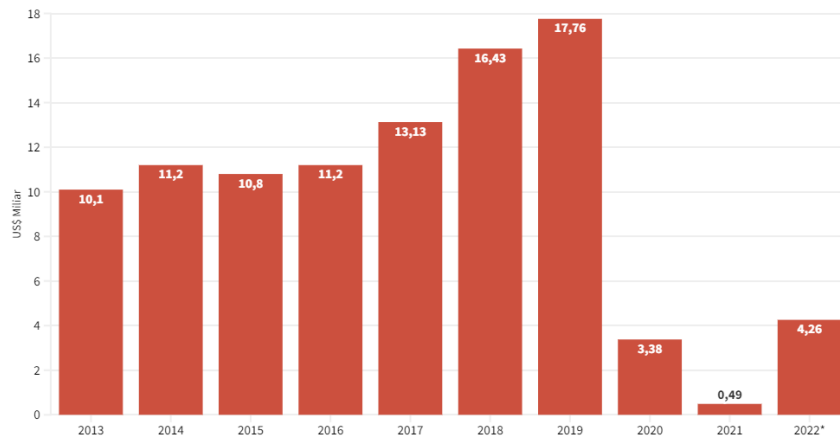


Figure 1. Foreign Exchange Income in the Tourism Sector

The decrease in foreign exchange income was a result of the decrease in foreign tourists visiting Indonesia by 58% in 2020. It continues to decrease until 2021, which is down by 78% from 2020 (Septiningrum & Pramuditya Soesanto, 2022),(Ramanathan & Meyyappan, 2019a). The decrease in the number of foreign tourists occurred in all regions of Indonesia, including Rembang.

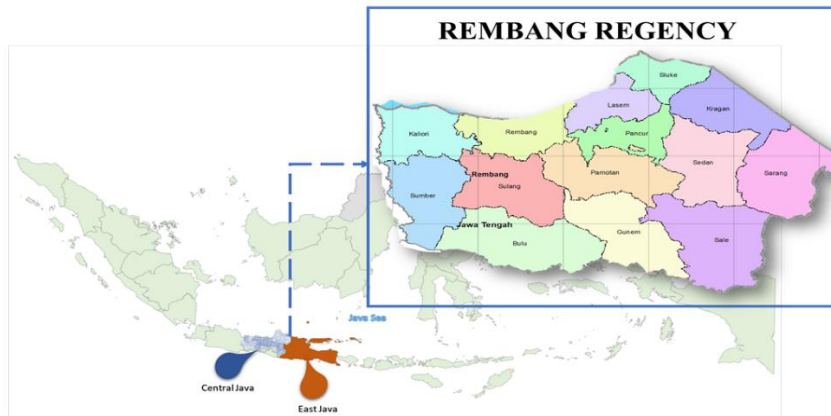


Figure 2. Administrative map of Rembang Regency Vs Indonesia

Based on Figure 2, Rembang is one of the regencies in Central Java Province, which is located on the northern coast of the island of Java. So many marine tourism objects are found. Apart from maritime arts, culture, religion, agro-tourism, and history have also developed in Rembang. During the COVID-19 pandemic, regional income from the tourism sector decreased because the number of tourists has decreased drastically. Based on information from (Septiningrum & Pramuditya Soesanto, 2022), in 2020, the number of domestic tourists decreased by 33% from 2019, and foreign tourists decreased by 75%. Detailed information can be seen in Table I.

Table 1 The Number Of Tourists In Rembang For 3 Years

Year	Tourists	
	Domestic	Foreign
2018	79652	597
2019	79848	313
2020	53456	75
2021	4756	29
2022	68294	376

The decrease in the number of tourists also resulted in a decrease in income, one of which was a decrease in revenue from hospitality. In 2022, the room occupancy rate (TPK) decreased by 7.83 points from 2018, and in 2022 it decreased by 1.15 points from 2018 [5].

The Rembang Regency Government has made efforts to restore the existence of tourism during COVID-19 by optimizing the digitization of promotions through the Tourism Information Center and digital tourism content through YouTube, which is managed by the Tourism Office (Rachman et al., 2022). Based on the Minister of Tourism and Creative Economy Regulation Number 9 of 2021 concerning Guidelines for Sustainable Tourism Destinations, digital tourism will boost the regional economy. So, using the digital media has provided benefits for the wider community who want to visit tourist objects during COVID-19 through the Tourism Information Center or enjoy tourist attractions on the YouTube channel. The government has provided a policy that has an impact on increasing the number of visitors to tourist attractions because concerns about COVID-19 are starting to subside. To find out how big the impact of the post-covid-19 pandemic on tourism in Rembang regency, this study will classify the results of community reviews on tourism after the pandemic period ends. One method that can be used to classify the results of reviews is sentiment analysis (Gu et al., 2018; Ramanathan & Meyyappan, 2019b). The methods used for classification is Support Vector Machine (SVM) (Gupta et al., 2019; Han et al., 2019; Krishna et al., 2017). Based on research (Lighthart et al., 2021) and (Ramanathan & Meyyappan, 2019b) SVM can provide a good classification of text mining. By using the Support Vector Machine (SVM), researchers can find out the classification results from reviews based on aspects commented on by the public relating to the impact of the post-covid-19 pandemic on tourism.

METHODS

In this study, a system was built to classify community reviews about the impact of the post-covid-19 pandemic on tourism. The method used for classification is Support Vector Machine (SVM). The processes involved in implementing this system include the data review reading stage, the feature extraction stage with text processing, the classification stage, and the validation and evaluation stage.

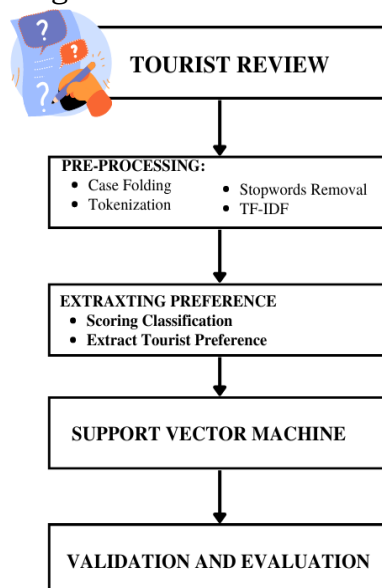


Figure 3. The method of extracting tourism preference.

The first stage is the stage of reading the opinion of visitors to tourist attractions due to the influence of social media. The dataset consists of documents containing one review of each document. After all the documents in the dataset have been read and collected, the

feature extraction stage is carried out with text processing. In this second stage, a collection of opinion texts will be processed so that the system will get features accompanied by their respective weight values (Yang et al., 2020),(Birjali et al., 2021),(Gu et al., 2018).

This stage begins with the case folding process. In this study, every letter of the alphabet in every text that has an uppercase form will be converted into a lowercase form. The next stage is the tokenization process. The text will be broken down into words called tokens. To separate text, the system uses a separator element in the form of elements other than the letters of the alphabet and a hyphen (-). The result of this tokenization is a collection of tokens from each text. This collection of tokens will then enter the stopwords removal process. In this process, each token will be matched with a word in the stopwords list. If they match, the token will be deleted. This process aims to make the computing process run faster by reducing unnecessary features. The results of the stopwords removal process will be used for feature selection. The feature of the word taken is the feature of the sentiment word lexicon (Jardim & Mora, 2022) The selected features are then used for classification. The next step is to calculate the weight of each selected feature using the term frequency-inverse document frequency (TF-IDF) method. The next stage after performing feature extraction with text processing is classification. There are two methods used in the classification, namely the SVM method. The SVM method is carried out according to Figure 4.

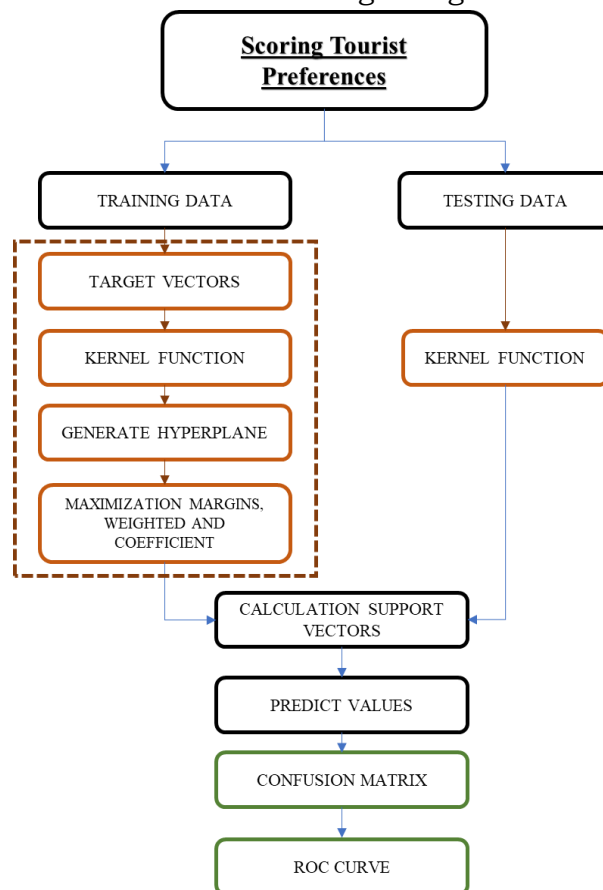


Figure 4. The method of classification tourism preference.

The next stage is the validation and evaluation stage. Validation is used to measure the accuracy of a model that is built based on a certain dataset. The method for this validation process is k-fold cross-validation (Gupta et al., 2019). This method divides the dataset into k equal parts of data. Then one part of the data is used as testing data, and the rest is used as training data. The final accuracy value is obtained from the average value of the accuracy

value in all iterations. The next stage is the evaluation process using the confusion matrix method. From this method, we can measure system performance by calculating accuracy, precision, specificity, and sensitivity.

RESULT AND DISCUSSIONS

A. Preprocessing

The explanation of the design process aims to provide a detailed description of each flow of method implementation in the opinion classification application of tourist attractions due to social media. In this step, the preprocessing technique used is Unigram where 1-character chunks are taken from a string. Blanks are added at the beginning and end of a string to determine the beginning and end of a string. For example, a string "TEXT" after adding the beginning and end with "_" instead of blank will get Uni-Gram T, E, X, T. This flow will be used in the implementation phase.

Phase Tourist Review: At this stage, the system reads documents that have been collected from tourist opinions related to tourist attractions in Rembang. The documents to be used consist of the same number of category 1 documents and category 2 documents. Each document contains text which is an opinion given the impact of the post-covid-19 pandemic on tourism in Rembang, Central Java, Indonesia, which was collected from the interview process and from tourist reviews on social media and web-based travel agents. The results of reading the document can be seen in Figure 5.

```
Sangat membantu
Pengaruhnya sangat besar.
Sangat berpengaruh untuk promosi tempat wisata, agar lebih banyak lagi orang yang tahu tentang
Tidak tahu
Sangat berpengaruh karena sosial media saat ini, Dapat menampilkan wisata diluar daerah saya
Saya penasaran dengan promosi instagram rumah merah sehingga saya kesini
Sangat berpengaruh untuk wisatawan yg mau berkunjung ke batik lasem
Sangat berpengaruh untuk media marketing batik lasem
Sangat berpengaruh untuk media marketing batik lasem
Sangat berpengaruh untuk media marketing pemerintah daerah
Sangat berpengaruh untuk menambah wisatawan luar kota
Bagus. Murah
Bagus. Masih perlu dikelola dengan baik lagi. Masih banyak sampah
Sangat berpengaruh
berpengaruh, bisa dilihat banyak orang jadi meningkatkan kemungkinan orang-orang ingin berwis
Bagus, tentunya akan mengenalkan kepada yang tidak tahu seperti saya. tp saya belum tahu apa r
Sangat signifikan terhadap perkembangan jumlah pengunjung setiap harinya
Cukup bagus dan gencar promonya di socmed
sosial media yg dikelola dg baik pasti akan sangat berpengaruh thd jumlah pengunjung di suatu
berpengaruh sekali bisa lebih mengenalkan kuliner rembang ke wisatawan luar daerah
berpengaruh sekali bisa lebih mengenalkan kuliner rembang ke wisatawan luar daerah
berpengaruh sekali, kadang penasaran dengan yang ada di sosial media
berpengaruh sekali, kelihatan estetik kalau di sosial media jadi penasaran
```

Figure 1. Text scrapping results.

Phase Case Folding: Case folding is the first process in the text pre-processing stage. Case folding aims to maximize further text processing. Tokenization is the second process after case folding. The result of this tokenization is that each document has a collection of tokens resulting from text splitting. The results of the case folding, and tokenization processes can be seen in Figure 6.

```
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
Rembang Wisata Digital Pantai Mangrove Batik Lasem Tuyuhan Estetik Pendapatn Daerah Viral Posting Terjangkau Murah Aktif
```

Figure 2. Case folding result

Phase Stopforward: The stopwords removal process is carried out after the tokenization process for each text using Indonesia Language. Each token will be matched with the words included in the stopwords word. If the token matches any of the stopwords

words, the token will be deleted. Some words that are included in Indonesian stopwords can be seen in Table 2.

Table 2 Example of Stop Forward Words Done

Words		
sehingga	setiap	aktif
yang	hanya	orang2
di	kepada	punya
ke	tapi	ya
paling	sangat	karena
kesana	agar	langsung
tidak	lebih	ada

Phase Weighting: Term Frequency-Inverse Document Frequency (TF-IDF) is a weighting method that is widely used. Term Frequency (TF) is the number of occurrences of the word in each document. Inverse Document Frequency (IDF) is the inverse value of Document Frequency (DF). IDF serves to reduce the weight of a term if its occurrence is spread throughout the document.

Table 3 Example Of Text Cloud Keyword Extraction Results

Keywords	Weight
Rembang	0,09
Digital	0,06
Wisata	0,05
Setuju	0,03
Pedagang	0,03

Of all the words taken using text scrapping, 5157 words were obtained. Table 3 shows 5 words that often appear in public reviews using social media and web-based travel agents related to tourism in Rembang, Central Java, Indonesia. After the weighting has been carried out, the next step is to classify the results of word occurrences using SVM so that people's preferences are obtained related to the form of tourism potential after the covid-19 pandemic. The classification carried out in this study uses the Support Vector Machine.

B. Support Vector Machine

Support vector machine is a machine learning method that aims to classify, in this study using SVM Cost 1 type, and Regression loss epsilon 0.1 while the kernel used is sigmoid. There are two categories of reviews/reviews produced, namely digital tourism, that need to be developed and choosing to go directly to tourist attractions after the covid-19 pandemic specificity and sensitivity. The four criteria can be calculated based on the confusion matrix obtained from the analysis. The confusion matrix in this study represents the results of the classification of public opinion/reviews regarding the impact of the post-covid-19 pandemic on tourism in Rembang, as follows.

Table 4 Confusion Matrix Of Reviews Related To Post-Pandemic

		Predicted			Total
		Digital tourism needs to be developed	to be	Direct travel	
Actual	Digital tourism needs to be developed	592		349	941
	Direct travel	712		697	1409
Total		1304		1046	2350

Based on the confusion matrix above, it is known that the AUC generated is 82.3%, the accuracy is 83%. The specificity is 85.5%, the sensitivity is 81.5%, and the precision is 85.1%. The sentiment analysis and classification are then made, where the research steps are in Figure 3 and Figure 4. The first step begins with sentiment analysis, namely taking reviews on digital media, then preprocessing the data consisting of case-folding, tokenization, stop word removal, and weighting. The result of this process is tourist preferences for tourism in Rembang regencyt, Central Java, Indonesia.

From the results of preprocessing text, data will be obtained using formal/standard words, the data will be used as text which will be processed to the next stage, namely keyword extraction. Data before and after preprocessing for #*wisatarembang* and #*wisatadigitalrembang* can be seen in tables 5 and 6. Of the 5157 keywords found containing #*wisatarembang*, #*wisatadigitalrembang*, not all keywords related to tourist preferences, only a few related, including Rembang, Digital, Wisata, Agree, and Traders. The word cheap has Rembang, and Digital has the highest weight because it appears in 384 reviews.

The results of the sentiment analysis obtained 2 categories, namely (1) people stated that digital tourism should be developed after the pandemic. When viewed from the word cloud, some words that emerged were Cheap, Affordable, Enjoyed, Income, and Regional. (2) the community states direct tourism. When viewed from the word cloud, some words that appear are Viral, Photo, Sentimental, Scenic, and Beautiful.

The results of sentiment analysis using the SVM classification method on community reviews about the post-pandemic impact are included in the good category because the AUC value obtained is above 0.8. This is in accordance with (Gorunescu, 2011), which says that if the resulting AUC value is between 0.8 – 0.9, then the resulting classification is included in the good category. Likewise, the accuracy obtained is 83%. These results also show that SVM gives good classification results. Calculation of the results of accuracy, precision, and recall with the value of k at different k-fold cross validation. Then an analysis will be carried out on the test results that have been obtained. The value of k in the k-fold cross validation that was tested was 2, namely k equal to 2 and k equal to 5. The number of features, the number of iterations and the number of documents tested were 5157 texts. The accuracy value obtained is the increase in the value of k affects the resulting accuracy, precision, and recall values. The accuracy value has increased from 79% to 83%. The precision value also increased from 73% to 75%. However, the recall value decreased from 65% to 69%. When sensitivity and specificity are connected, an ROC curve is formed, as shown in the figure below.

for digital and conventional tourism in Rembang Regency. The word that is used with the most frequency appears to be visual dominance and occupies the image in a proportion that exceeds the other words. In Figure 8 there is a dominance of the conjunction "Rembang" and the visualization of the words "Digital" and "Tourism" repeatedly. There are also the words "Forest", "Museum", and "Karang Jahe" with quite large visual proportions which indicate the large role of the core part in this wordcloud image. So, it can be seen that the community review supports digital tourism in Rembang Regency as an alternative form of public entertainment amid certain conditions. This research is also supported by data taken primarily from visitors to tourist places in Rembang. Data was taken from March 2022 to June 2022 as many as 300 visitors spread over 10 tourism places consisting of six natural tourism places and four local handicraft industries. The demographic details of the respondents are shown in Figure 9.

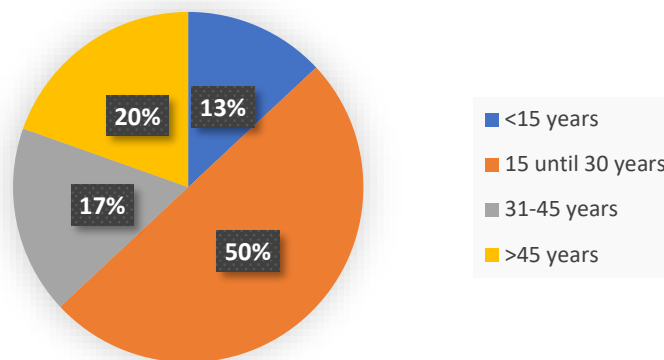


Figure 5. Descriptive statistics of respondent.

From the results of data collection, we know that majority tourist is around 15-30 years, it means that 50% of tourist have social media and have ability to make a review of tourism experience digital platform in Rembang regency.

CONCLUSION

The results of the sentiment analysis on the community review regarding the post-pandemic impact obtained 2 categories, namely (1) the community stated that digital tourism should be developed after the covid-19 pandemic, and (2) the community stated that they chose direct tourism after the Covid-19 case had slowed. Sentiment analysis results are based on the classification process using SVM and obtained an accuracy of 83% and AUC of 82.3%. The results of sentiment analysis in this study are in a good category. Research suggestions for Rembang local governments to optimize continue develop digital tourism and continue to improve real tourism facilities and infrastructure in Rembang regency, where the pathway will increase tourist visits and the economy of the community around tourist attractions. After this study we want to make suggestions policy in the next research.

ACKNOWLEDGEMENTS

We express our gratitude to the Rembang city government for creating content on various social media about local tourism, and to the Telkom University and Universitas Negeri Yogyakarta, Indonesia for providing support for researchers to always develop. The completion of this research was also supported by a lecturer in the Research Group Enterprise System Solution, Telkom University who helped in data collection.

REFERENCES

- Birjali, M., Kasri, M., & Beni-Hssane, A. (2021). A comprehensive survey on sentiment analysis: Approaches, challenges and trends. *Knowledge-Based Systems*, 226. <https://doi.org/10.1016/j.knosys.2021.107134>
- Gorunescu, F. (2011). Classification Performance Evaluation. In F. Gorunescu (Ed.), *Data Mining: Concepts, Models and Techniques* (pp. 319–330). Springer Berlin Heidelberg. https://doi.org/10.1007/978-3-642-19721-5_6
- Gu, Y. H., Yoo, S. J., Jiang, Z., Lee, Y. J., Piao, Z., Yin, H., & Jeon, S. (2018). Sentiment analysis and visualization of Chinese tourism blogs and reviews. *2018 International Conference on Electronics, Information, and Communication (ICEIC)*, 1–4. <https://doi.org/10.23919/ELINFOCOM.2018.8330589>
- Gupta, A., Tyagi, P., Choudhury, T., & Shamoan, M. (2019). Sentiment Analysis Using Support Vector Machine. *2019 International Conference on Contemporary Computing and Informatics (IC3I)*, 49–53. <https://doi.org/10.1109/IC3I46837.2019.9055645>
- Han, K., Chiu, C.-C., & Chien, W. (2019). The Application of Support Vector Machine (SVM) on the Sentiment Analysis of Internet Posts. *2019 IEEE Eurasia Conference on IOT, Communication and Engineering (ECICE)*, 154–155. <https://doi.org/10.1109/ECICE47484.2019.8942736>
- Jardim, S., & Mora, C. (2022). Customer reviews sentiment-based analysis and clustering for market-oriented tourism services and products development or positioning. *Procedia Computer Science*, 196, 199–206. <https://doi.org/10.1016/j.procs.2021.12.006>
- Krishna, M. H., Rahamathulla, K., & Akbar, A. (2017). A feature based approach for sentiment analysis using SVM and coreference resolution. *2017 International Conference on Inventive Communication and Computational Technologies (ICICCT)*, 397–399. <https://doi.org/10.1109/ICICCT.2017.7975227>
- Lighthart, A., Catal, C., & Tekinerdogan, B. (2021). Systematic reviews in sentiment analysis: a tertiary study. *Artificial Intelligence Review*, 54(7). <https://doi.org/10.1007/s10462-021-09973-3>
- Rachman, F. H., Imamah, & Rintyarna, B. S. (2022). Sentiment Analysis of Madura Tourism in New Normal Era using Text Blob and KNN with Hyperparameter Tuning. *2021 International Seminar on Machine Learning, Optimization, and Data Science (ISMODE)*, 23–27. <https://doi.org/10.1109/ISMODE53584.2022.9742894>
- Ramanathan, V., & Meyyappan, T. (2019a). Twitter Text Mining for Sentiment Analysis on People's Feedback about Oman Tourism. *2019 4th MEC International Conference on Big Data and Smart City (ICBDSC)*, 1–5. <https://doi.org/10.1109/ICBDSC.2019.8645596>
- Ramanathan, V., & Meyyappan, T. (2019b). Twitter Text Mining for Sentiment Analysis on People's Feedback about Oman Tourism. *2019 4th MEC International Conference on Big Data and Smart City (ICBDSC)*, 1–5. <https://doi.org/10.1109/ICBDSC.2019.8645596>
- Septiningrum, L., & Pramuditya Soesanto, R. (2022). Tourism Itinerary Design: User Experience Approach. *2022 International Conference Advancement in Data Science, E-Learning and Information Systems (ICADEIS)*, 01–05. <https://doi.org/10.1109/ICADEIS56544.2022.10037382>
- Statistics Indonesia. (n.d.). *Total Foreign Exchange of Tourism Sector (Billion US \$), 2016-2018*.
- Statistics of Rembang Regency. (2022, February). *Arrivals of Tourist (People) 2018-2020*. Statistics of Rembang Regency.
- Yang, L., Li, Y., Wang, J., & Sherratt, R. S. (2020). Sentiment Analysis for E-Commerce Product Reviews in Chinese Based on Sentiment Lexicon and Deep Learning. *IEEE Access*, 8. <https://doi.org/10.1109/ACCESS.2020.2969854>