



Journal of Advanced Research in Applied Sciences and Engineering Technology

Journal homepage:
https://semarakilmu.com.my/journals/index.php/applied_sciences_eng_tech/index
ISSN: 2462-1943



Sentiment Analysis on The Place of Interest in Malaysia

Qiryndriana binti Kharul Zaman¹, Wan Nur Syahidah binti Wan Yusoff^{1,*}, Qistina Batrisyia bin Azman Shah

¹ Centre for Mathematical Sciences, Universiti Malaysia Pahang Al-Sultan Abdullah, Lebuhraya Persiaran Tun Khalil Yaakob, 26300, Gambang, Pahang, Malaysia

² Credence Tech, Petaling Jaya, Selangor, Malaysia

ARTICLE INFO

ABSTRACT

Article history:

Received 25 September 2023

Received in revised form 8 November 2023

Accepted 6 March 2024

Available online 8 April 2024

Keywords:

Twitter, Natural Language Processing, Sentiment Analysis, Machine Learning, Tourism

This study focuses on utilizing machine learning methods for sentiment analysis to identify positive and negative comments regarding Malaysian Places of Interest. The data was collected from Twitter using social media monitoring software and organized into tables. Pre-processing techniques and Natural Language Processing (NLP) methods were applied to handle missing values and prepare the text data for analysis. The dataset was then split into training and testing sets, and three supervised learning algorithms which are Support Vector Machine, Random Forest, and Naive Bayes were employed to evaluate the sentiment analysis models. The performance of each model was compared, and it was found that Support Vector Machine achieved the highest accuracy, recall score, F1 score, and precision score. This study demonstrates the potential to extend sentiment analysis to analyze sentiments expressed in texts written in the Malay language by utilizing the Malaya corpus. Additionally, visual dashboards can be created to present the findings and provide recommendations based on the insights gathered from the sentiment analysis of Malaysian Places of Interest feedback.

1. Introduction

Millions of people globally are actively using social media platforms like Facebook, Twitter, Instagram, LinkedIn, and others, making social media a crucial part of their daily lives. In recent times, social media has become increasingly significant for businesses as it offers various benefits that can help in their growth and success. According to the Smart Insights website by Chaffey [1], as of January 2023, social media is utilised by 59% of the global population, and individuals spend an average of 2 hours and 31 minutes per day on social media platforms. That shows half of the global population is using social media on a daily basis and these facts can help businesses in so many ways.

In the current fast-moving and cut-throat business landscape, businesses that can efficiently gather, evaluate, and make use of data have a considerable edge over those that cannot, which mentioned by McKinsey & Company [2]. Additionally, the significance of data has grown even more with the emergence of artificial intelligence and machine learning. According to IBM [3], social media

* Corresponding author.

E-mail address: wnsyahidah@ump.edu.my

<https://doi.org/10.37934/araset.43.1.5465>

also generates so much unstructured data such as texts, emails, social media posts, video recordings and images. With these data, social media can offer valuable understanding into the actions of customers, trends in the market, and operational activities of a business. This allows businesses to make knowledgeable decisions and enhance their overall efficiency.

Within this context of research, sentiment analysis can be a useful method for businesses. Sentiment analysis refers to the technique of utilizing machine learning and natural language processing (NLP) to identify and extract subjective information from text data. Besides, Sentiment analysis involves using computational methods to process text and determine the opinions, attitudes, and subjective expressions within the text, as written by Liu [4]. There are few methods that can be used for sentiment analysis such as lexicon-based approach and machine learning approach. In this research, machine learning approach will be employed. Machine learning is a widely used technique in sentiment analysis since it can learn from data and enhance its precision with experience. With the aid of vast quantities of labelled data, machine learning algorithms can detect correlations and patterns that enable them to classify the sentiment of fresh text data with accuracy. The researchers combined several machine learning algorithms, including Support Vector Machines (SVM) and Decision Trees, to classify the sentiment of customer reviews.

This research focuses on enhancing the accuracy of customer sentiment analysis through social media feedback by comparing different machine learning algorithms. The objective is to identify the most effective algorithm that can predict the sentiment of new texts with high accuracy. Additionally, a sentiment analysis dashboard has been developed to provide clients with a quick and comprehensive understanding of the sentiment analysis outcomes, including sentiment scores, word clouds, and sentiment trends. The research emphasizes the importance of visualizing sentiment analysis results in a dashboard to monitor shifts in customer sentiment and make informed decisions based on the data. Overall, sentiment analysis is recognized as a crucial aspect for businesses to maintain a strong market presence and deliver high-quality products and services, particularly in Malaysia.

2. Related Studies

The term 'Social Media' is no longer strange to everybody since more than half of the world's population uses social media. As for tourism, social media within the tourism industry refers to the utilization of social media channels to facilitate tourists in exchanging and perusing user-created content to aid in their decision-making regarding travel, as described by Gretzel and Yoo [5]. According to Russell and Norvig [6], NLP is an interdisciplinary field that incorporates computer science, psychology, and linguistics to design and create tools and algorithms for analyzing and interpreting human language. Since NLP has been growing rapidly, it has been used or applied in many applications such as sentiment analysis, text classification, and more. As an example, social media activity is being assessed to understand customer sentiment towards a product or service, as discussed in Manning and Schütze [7]. Miguel *et.al.*, [8] utilization of NLP gives significant advantage in obtaining information on the mass content generated by online users concerning tourism services and products.

The whole population in this world have already acknowledged the uses of social media, since there are many benefits that users can obtain while using it, especially to the tourism industry. As reported by Liu [4], sentiment analysis is a technique within the realm of NLP that is utilized to recognize and extract the subjective content from text, encompassing emotions, opinions and attitudes. Besides, there are many advantages of performing sentiment analysis on social media, specifically for tourism industry. Firstly, the use of sentiment analysis by tourism organizations can

aid in monitoring and responding promptly and effectively to feedback from customers, resulting in enhanced customer satisfaction and loyalty (Xiang *et. al.*, [9]). Next, according to Munar and Jacobsen [10], by analyzing the subjective information in customer feedback, such as opinions, emotions and attitudes using sentiment analysis, tourism businesses can anticipate future trends and adjust their strategies to meet customer needs.

There are three approaches that can be used in sentiment analysis, those are Knowledge-Based approach (KBS), statistical approach which uses machine learning classifiers and lastly, hybrid approach. As claimed by Kumar [11], machine learning is applied to classify whether a given input shows a positive, negative, or neutral sentiment. Hasan *et. al.*, [12] used machine learning algorithms such as Naïve Bayes and Support Vector Machine (SVM), and Random Forest to analyze political views on social media.

3. Methodology

This research utilized primary data collection since the tweets were extracted from Determ, a social media monitoring software, to obtain desired outcomes. Primary data collection involves obtaining first-hand information directly from its source through various techniques like surveys, interviews, or experiments. For this research, the tweets were extracted from Determ with desired keywords. Since the tweets are about Places of Interest (PoI) in Malaysia, there are important keywords such as Malaysia, Selangor, Kuala Lumpur, Kedah and more. To extract the tweets containing these keywords, Determ has provided a section where users can put the desired keywords in a query. The tweets will then be collected by Determ. In this step, the tweets are ready to be exported. The dataset is stored in the file separated by commas (CSV) and is ready to use for further steps. The next step in the analysis is preprocessing the data but this study employs two types of preprocessing which are text preprocessing and data preprocessing. Text preprocessing only performs the NLP techniques on the textual data while data preprocessing performs the operations to the whole data. The text preprocessing includes the removal of non-ASCII characters, the conversion to all lowercase letters, the removal of HTML symbols, mentions and links, the removal of punctuations, the spelling corrector using Malaya, the lemmatization using Malaya, the text segmentation using Malaya and sentiment classification using Malaya. Malaya corpus is utilized since this study is about the PoI in Malaysia and most of the tweets extracted are in Bahasa Malaysia. Data preprocessing assists in cleaning the whole data in terms of handling missing values and duplicate rows, label encoding, feature extraction, data partitioning, imbalanced dataset and Synthetic Minor Oversampling Technique (SMOTE).

As part of the machine learning process, this study will train the PoI data using machine learning algorithms such as Naïve Bayes, SVM and Random Forest. The algorithms will be compared based on its accuracy, precision, recall and F1 score in order to choose the best algorithm to create a predictive modelling. Figure 1 shows the flow of this research.

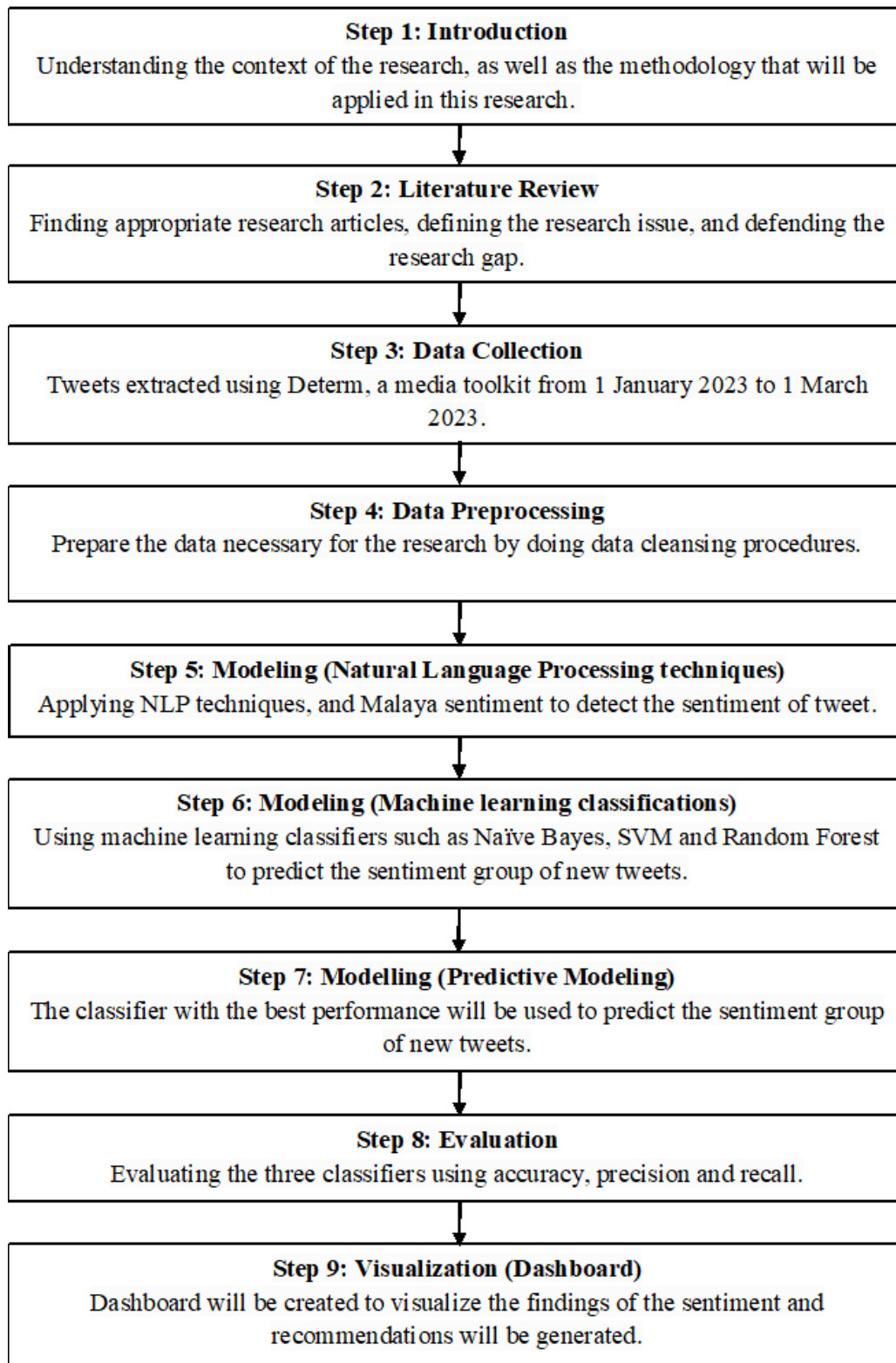


Fig. 1. The research framework

3.1 Predictive Modelling

This study will construct multiple supervised machine learning algorithms in Python to create a predictive model. The three types of machine learning algorithms utilized in this study include

Support Vector Machine (SVM), Random Forest, and Naïve Bayes. These specific models were chosen based on their consistent performance superiority over other models, as evidenced by previous literature.

3.1.1 Naïve Bayes

Naive Bayes is a widely used classification algorithm in machine learning which can be seen in many researches, such as previous studies [12-14] that applies Bayes' theorem to assign probabilities to different classes. According to Bayes' theorem, the likelihood of a class given certain input features is proportional to the likelihood of those features given the class multiplied by the prior probability of the class. Naive Bayes makes the "naive" assumption that the input features are independent of each other given the class, which simplifies the computation of probabilities and improves the algorithm's efficiency. Bayes' theorem is a mathematical formula used to calculate the probability of a hypothesis or event occurring given evidence or data. According to Gupta [15], Eq. (1) is the equation of Bayes' theorem.

$$P(A|B) = (P(B|A) * P(A))/P(B) \quad (1)$$

where:

$P(A)$ = The probability of A occurring

$P(B)$ = The probability of B occurring

$P(A|B)$ = The probability of A given B

$P(B|A)$ = The probability of B given A

3.1.2 Support Vector Machine

Support Vector Machines (SVM) is a machine learning technique that is useful for solving complex problems in both classification and regression analysis with high dimensional feature spaces, as can be seen also in [12- 16]. In SVM, the primary objective is to determine the optimal hyperplane that separates data points belonging to different classes with maximum margin. The hyperplane that has the largest margin between data points of different classes is selected as the optimal one. Support vectors are data points that are closest to the hyperplane and are crucial for determining the hyperplane. SVM can deal with both linearly separable and non-linearly separable data by using a technique called kernel trick. Kernel trick maps the input space into a higher-dimensional feature space where the data points may be linearly separable. This capability enables SVM to identify non-linear decision boundaries. SVM algorithm can be found in Python language.

3.1.3 Random Forest

Random Forest is a technique that is utilized for performing classification and regression tasks. This method has been discussed by many authors such as [12 and 17]. It belongs to the ensemble method of algorithms and utilizes multiple decision trees to make predictions. Random Forest works by training a vast number of decision trees on various subsets of the input data and input features. Every decision tree in the forest predicts a result, and the final prediction is made based on the majority vote of all the decision trees. This methodology helps in reducing overfitting and improves

the accuracy and robustness of the model. Random Forest is capable of managing both categorical and numerical data, and it is also capable of dealing with missing data and outliers. Random Forest algorithm can be found in Python language.

3.2 Evaluation

Evaluation in machine learning refers to the process of analyzing the effectiveness of a model on a given dataset. It involves calculating one or more metrics that measure the model's performance in predicting or classifying new data. The selection of a specific evaluation metric depends on the objectives of the task and the model itself. In this research, accuracy, precision, recall and F1 score will be the metrics used in evaluating the machine learning models. There are few terms used in every evaluation method such as True Positive (TP), False Positive (FP) and more.

The prediction involves different categories, including True Positive Reviews (TP), False Positive Reviews (FP), True Negative Reviews (TN), and False Negative Reviews (FN). TP represents the number of samples correctly classified by the model as the positive class, TN represents the number of samples correctly classified as the negative class, FP represents the number of samples incorrectly classified as the positive class, and FN represents the number of samples incorrectly classified as the negative class.

The model's performance can be evaluated using various metrics such as accuracy, precision, recall, and F1-measure, which can be obtained through Python by utilizing the classification report function. By examining the confusion matrix and employing these metrics, the model's effectiveness can be assessed comprehensively. According to [18-20], the accuracy, precision, recall and F1-measure are as in Eq. (2) – Eq. (5).

$$Accuracy = \frac{TP+TN}{TP+FP+TN+FN} \quad (2)$$

$$Precision = \frac{TP}{TP+FP} \quad (3)$$

$$Recall = \frac{TP}{TP+FN} \quad (4)$$

$$F1 \text{ score} = 2 * \frac{Recall * Precision}{Recall + Precision} \quad (5)$$

4. Results and Discussion

Under this section, a comprehensive analysis and discussion of all the results obtained in this study will be presented. The discussion commences by delving into the examination and interpretation of the data, providing a detailed exploration of the collected information.

4.1 Machine Learning Algorithms

The results indicate the performance accuracy achieved by three different machine learning algorithms on a given task or dataset. Naive Bayes attained an accuracy of 89.65%, signifying that it correctly classified approximately 89.65% of the instances. In comparison, the Support Vector Machine (SVM) algorithm achieved a higher accuracy of 92.40%, indicating its ability to correctly classify a larger proportion of instances. Random Forest, with an accuracy of 88.55%, exhibited slightly lower performance than Naive Bayes but still showcased commendable accuracy in its

classification predictions. The accuracy metric represents the proportion of instances correctly classified out of the total instances evaluated, making it an important measure of algorithm performance. Based on these results, it can be inferred that SVM performed the best among the three algorithms, followed closely by Naive Bayes, while Random Forest demonstrated slightly lower accuracy but still achieved notable performance in the task at hand. It is worth noting that the choice of the most suitable algorithm should consider other factors beyond accuracy, such as dataset characteristics and specific requirements of the problem domain.

In this research, there are few ways to evaluate performance including accuracy, precision, recall, and F1 score. The other evaluation methods such as precision, recall and F1 score are generated as well. Figure 2 shows a report that represents in-depth evaluation measures for the three models used. The SVM model has the best performance, with precision and recall scores ranging from 0.067 to 1.0. The second rank is the Naïve Bayes, with precision and recall scores ranging from 0.27 to 0.95. Lastly, the Random Forest with precision and recall scores ranging from 0.27 to 0.99.

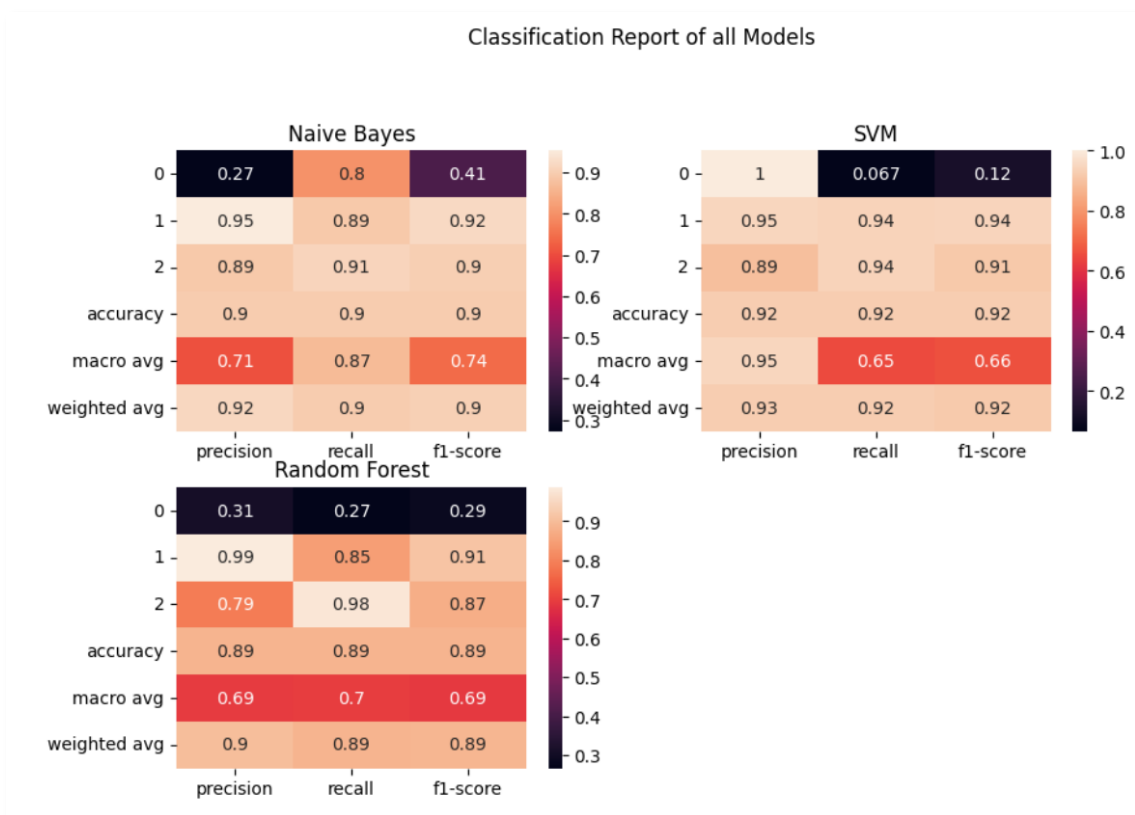


Fig. 2. The classification report for the three models used

4.2 Predictive Modelling

The goal of developing a predictive model is that the model can extract insights and correlations from past data, in order to recognize patterns and correlations between the input features and the output variable. This extracted information is then utilized to make precise predictions on new data. From the above subtopics, it turns out that SVM has the highest accuracy among other algorithms, hence, SVM will be utilized to be a predictive model. The predictive model will predict and classify new tweets related to the Pol to their sentiment group. The new tweets will then be processed using the seven steps of text pre-processing that have performed earlier. Lastly, a data frame consists of

the tweets and the predicted sentiment is displayed as shown in Figure 3. There are three sentiments which are positive, negative, and neutral. 0 indicates negative tweets, 1 indicates positive tweets and 2 indicates neutral tweets. From the data frame, it turns out that SVM is able to predict and classify the tweets into its sentiment group correctly. As an example, the model classified the first tweet as neutral sentiment which is surprisingly, correct.

	text2	array2
0	makan nasi kandaq ke itu	2
1	sungai kelantan darul naim sangat cantik	1
2	sebab kenapa aku lebih rela travel pergi luar ...	0
3	semua manusia dari seluruh negara dialu-alukan...	2
4	Pengalaman menarik di Sunway Lagoon tidak akan...	1

Fig. 3. A data frame consists of the new tweets and the sentiment predicted by SVM

4.3 Dashboard Accelerator

This subtopic will show the results of the social media data in the form of dashboard accelerators. According to Tableau, Tableau Accelerators are pre-built dashboards that can be seamlessly integrated with any data and tailored to suit the specific requirements, expediting the process of obtaining data-driven insights. These accelerators are readily available for immediate use and offer the flexibility to be customized according to business needs, empowering business to quickly derive valuable insights from the data. There are many accelerators that can be used to cater a wide range of use cases in different departments, industries, and enterprise applications.

The homepage of the study displays Figure 4, featuring the title "PoI Analytics" and an introduction. The study explores how the Ministry of Tourism and Culture Malaysia (MOTAC) and businesses can benefit from using an accelerator to analyze tourist attractions that attract attention. By leveraging social media feedback shared by tourists, MOTAC and businesses can improve the landscapes and services of these attractions. The dashboard emphasizes the importance of answering key business questions, which can guide strategic decision-making and help achieve business goals. Visualizations provide insights into these questions. Additionally, the dashboard allows monitoring and improvement of insights by focusing on attributes like "Total Reach," "Total Engagements," and "Sentiment." Consistently analyzing these attributes helps businesses understand their target audience's perception and enhance their online presence. On the right-hand side of the homepage, there are three navigation buttons. The first button leads to an Executive Summary page summarizing PoI analysis for a specific period. The second button directs users to a dashboard page displaying insights over time. Lastly, the third button leads to a page focused on analyzing sentiment in tweets. Users can access these pages by using Alt+click on the desired button.

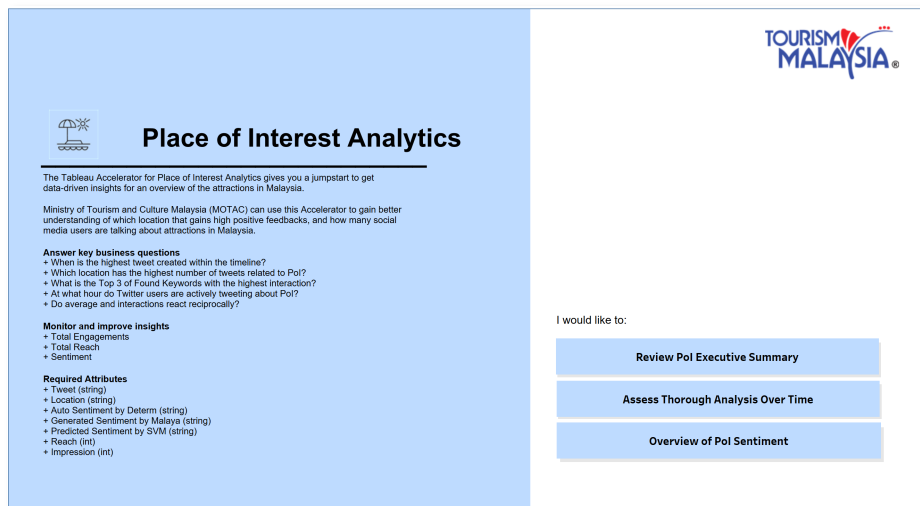


Fig. 4. Pol Analytics

Figure 5 is the "Executive Summary" page of the dashboard, presenting concise visualizations of common findings. It includes general statistics such as the total number of tweets extracted (3026) and the breakdown of tweet types. The line chart displays the tweet count within the timeline, with the highest number occurring on January 23, possibly due to the Chinese New Year holiday. The map visualizes the location of original Pol tweets, showing that Malaysia had the most tweets, followed by the United States and Indonesia. The table highlights the top three keywords with the highest tweet interactions, and a word cloud showcases the frequency of keywords, with 'pulau' being the most prominent.

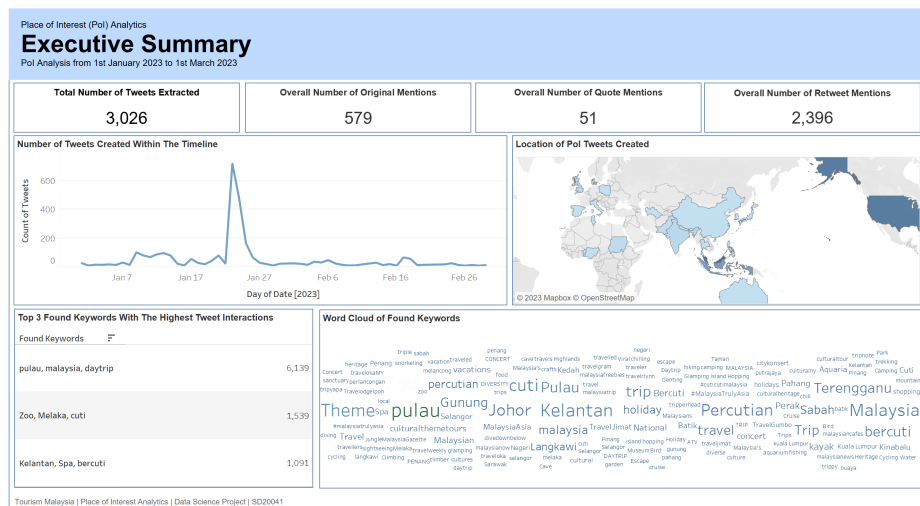


Fig. 5. Executive summary

Figure 6 in the dashboard focuses on the temporal aspects of Pol analysis. It includes three visualizations. The first line chart shows the hourly distribution of original Pol mentions, with 9 PM having the highest activity and 5 AM the lowest. The second line chart demonstrates an inverse relationship between reach and interactions by week, indicating that high reach doesn't always correspond to high interactions. The Gantt chart highlights Indonesia and Malaysia as the countries with the highest occurrence of Pol tweets throughout the day, with specific peak hours in each country.

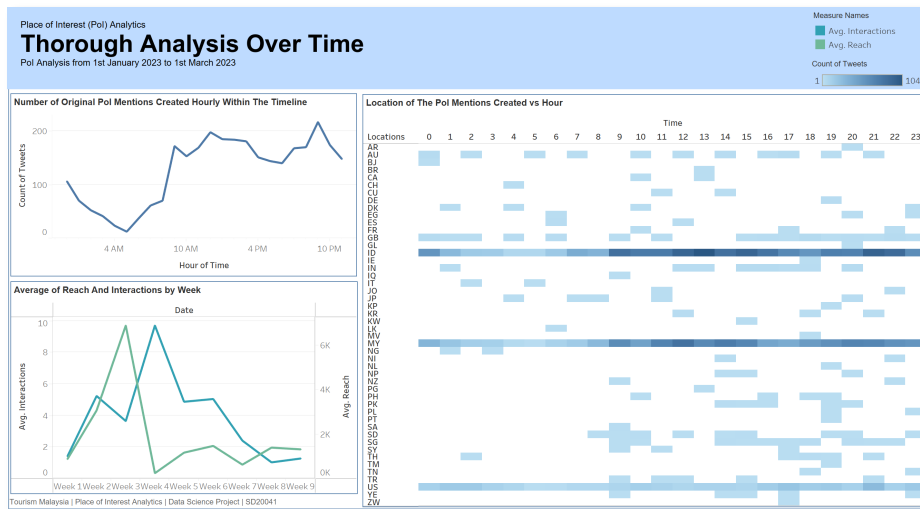


Fig. 6. Temporal aspects of PoI analysis

Figure 7 is the final page of the dashboard, focusing on the sentiment analysis of PoI. The page presents visualizations and findings related to sentiment. The line chart shows the overall number of tweets over time based on sentiment, with positive and neutral tweets having higher counts than negative tweets. Two boxes highlight a positive tweet with high reach about a new international flight route and a negative tweet with lower reach about the inability to climb Mount Kinabalu due to insufficient holidays. The bar charts compare sentiment classification by Determ (1540 neutral, 922 positive, and 564 negative) with the sentiment predicted by Malaya (1842 positive, 1146 neutral, and 38 negative). Malaya suggests that Determ misclassified more than half of the negative tweets. The findings provide insights for recommendations and further analysis.

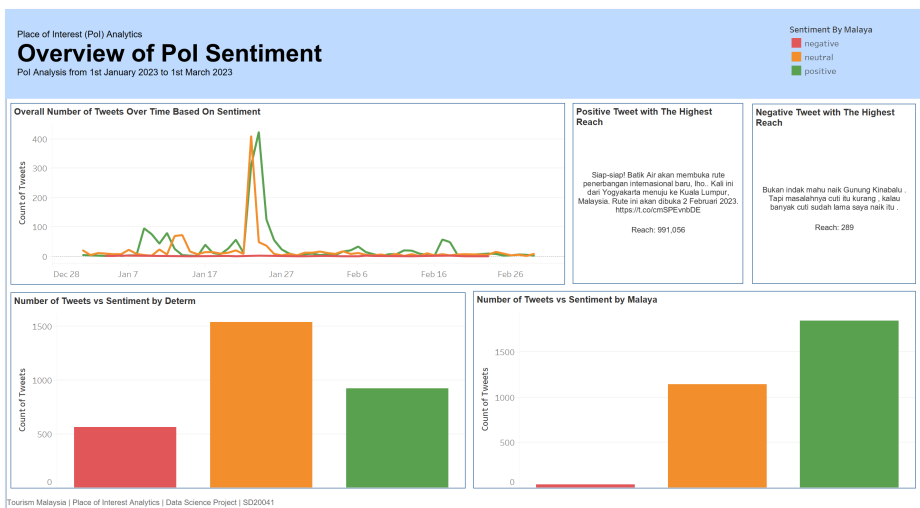


Fig. 7. Overview of PoI sentiment

5. Conclusions

In this study, there are three machine learning algorithms that will be employed to train which are the Naïve Bayes, SVM and Random Forest. Furthermore, evaluation methods such as F1 score, recall, accuracy and precision will be generated to overlook which model has the highest accuracy

and suitable to be a predictive model. Last part of the second objective is the predictive model. SVM gains the highest accuracy and will be chosen to be the model and predict the sentiment of new tweets. There are five new tweets related to Pol that have been through data pre-processing and the sentiment is predicted by SVM.

Lastly, this study concludes by providing recommendations based on the findings from the dashboard. The recommendations address the key business questions identified earlier. The highest number of tweets were created on January 23rd, suggesting businesses refine their marketing strategies to gain insights into customer preferences. Malaysia had the highest number of tweets related to Pol, indicating the need for multilingual support to cater to international tourists. The top three keywords with the highest interactions were identified, allowing businesses to develop targeted marketing campaigns and improve visitor experiences. The peak tweet activity occurred at 9 PM, presenting an opportunity for businesses to share timely information about tourism offerings. Positive sentiment tweets were the highest, prompting businesses to encourage positive feedback, address negative feedback, and make improvements based on customer sentiment analysis. Overall, data-driven decision making in the tourism industry enables businesses to enhance their competitiveness, customer satisfaction, resource allocation, and risk management.

Acknowledgement

The authors would like to thank Universiti Malaysia Pahang Al-Sultan Abdullah for the financial support. This research was not funded by any grant. This research was funded by International Matching Grant, UIC221521.

References

- [1] Chaffey, D. Global Social Media Statistics Research Summary 2022 [June 2022]. Smart Insights. 2023.
- [2] McKinsey & Company. The Age of Analytics: Competing in a Data-Driven World. Retrieved May 4, 2023.
- [3] IBM. Managing Unstructured Data: A Guide to Understanding and Maximizing Your Non-numeric Wealth. [2019](#).
- [4] Bing Liu, "Sentiment analysis and opinion mining". *Synthesis Lectures on Human Language Technologies* 5, no. 1 (2012): 1–167. <https://doi.org/10.1007/978-3-031-02145-9>
- [5] Gretzel Ulrike and Yoo Kyung-Hyan, "Use and Impact of Online Travel Reviews. In: O'Connor, P., Höpken, W., Gretzel, U. (eds)". *Information and Communication Technologies in Tourism* 2008. Springer, Vienna. https://doi.org/10.1007/978-3-211-77280-5_4
- [6] Russell, Stuart and Norvig, Peter. *Artificial Intelligence, A Modern Approach: Natural Language Processing*. Pearson Education, Inc. 2010
- [7] Manning, Chris and Schütze, Hinrich. *Foundations of Statistical Natural Language Processing*, MIT Press. Cambridge, MA: May 1999.
- [8] Miguel Á. Álvarez-Carmona, Ramón Aranda, Ansel Y. Rodríguez-Gonzalez, Daniel Fajardo-Delgado, María Guadalupe Sánchez, Humberto Pérez-Espinosa, Juan Martínez-Miranda, Rafael Guerrero-Rodríguez, Lázaro Bustio-Martínez and Ángel Díaz-Pacheco, "Natural language processing applied to tourism research: A systematic review and future research directions". *Journal of King Saud University - Computer and Information Sciences* 34, no. 10, Part B (2022): 10125-10144, <https://doi.org/10.1016/j.jksuci.2022.10.010>.
- [9] Zheng Xiang, Qianzhou Du, Yufeng Ma and Weiguo Fan, "A comparative analysis of major online review platforms: Implications for social media analytics in hospitality and tourism". *Tourism Management* 58 (2017): 51-65, <https://doi.org/10.1016/j.tourman.2016.10.001>
- [10] Ana María Munar and Jens Kr. Steen Jacobsen, "Motivations for sharing tourism experiences through social media". *Tourism Management* 43 (2014): 46-54, <https://doi.org/10.1016/j.tourman.2014.01.012>.
- [11] Kumar, A. Sentiment Analysis & Machine Learning Techniques. Reimagining Data-driven Society with Data Science & AI. 2021.
- [12] Hasan, Ali, Sana Moin, Ahmad Karim, and Shahaboddin Shamshirband, "Machine Learning-Based Sentiment Analysis for Twitter Accounts". *Mathematical and Computational Applications* 23, no. 1 (2018): 11, <https://doi.org/10.3390/mca23010011>
- [13] McCallum, Andrew and Nigam, Kamal, "A comparison of event models for Naive Bayes text classification". *AAAI-98 workshop on learning for text categorization* 752, no. 1 (1998): 41-48.

- [14] M. S. Neethu and R. Rajasree, "Sentiment analysis in twitter using machine learning techniques," *2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT)*, Tiruchengode, India, (2013): 1-5. <https://doi.org/10.1109/ICCCNT.2013.6726818>
- [15] Prashant Gupta (2017). Naïve Bayes in Machine Learning. Towards Data Science.
- [16] Kristin P. Bennett and Colin Campbell, "Support vector machines: hype or hallelujah?". *ACM SIGKDD Explorations Newsletter* 2, no. 2 (Dec. 2000): 1–13. <https://doi.org/10.1145/380995.380999>
- [17] Breiman, L, "Random Forests". *Machine Learning* 45 (2001): 5–32. <https://doi.org/10.1023/A:1010933404324>
- [18] Fei Deng, Jibing Huang, Xiaoling Yuan, Chao Cheng and Lanjing Zhang, "Performance and efficiency of machine learning algorithms for analyzing rectangular biomedical data". *Laboratory Investigation* 101, no. 4 (2021): 430-441. <https://doi.org/10.1038/s41374-020-00525-x>.
- [19] Disha, Raisa Abedin and Waheed, Sajjad, "Performance analysis of machine learning models for intrusion detection system using Gini Impurity-based Weighted Random Forest (GIWRF) feature selection technique". *Cybersecurity* 5, no. 1 (2022). <https://doi.org/10.1186/s42400-021-00103-8>
- [20] Kasongo, Sydney M. and Sun, Yanxia, "Performance Analysis of Intrusion Detection Systems Using a Feature Selection Method on the UNSW-NB15 Dataset". *Journal of Big Data* 7, no. 105 (2020). <https://doi.org/10.1186/s40537-020-00379-6>