

# Analysis of Sentiment on Aquaria KLCC Reviews based on TripAdvisor

Norazmiera Ayunie binti Azman  
Masters of Data Science (2022972987)  
College of Computing, Informatics and  
Mathematics  
Universiti Teknologi MARA (UiTM)  
Shah Alam, Malaysia  
norazmieraayunie@gmail.com

**Abstract**— Sentiment analysis, an application of natural language processing (NLP), is widely used to comprehend people's opinions and emotions towards various subjects. In this study, sentiment analysis of reviews related to Aquaria KLCC, sourced from TripAdvisor is used. By employing the CRISP-DM methodology, this study follows the structure such as project understanding, data understanding, data pre-processing, modeling, and evaluation. The dataset was self-extracted from TripAdvisor using WebHarvy and underwent pre-processing with Python and Orange. This study deployed nine diverse machine learning models consisting of Support Vector Machine, Logistic Regression and Random Forest with various parameters to predict sentiments expressed in the reviews, and compare the performance metrics. This study revealed that the Support Vector Machine (SVM) model with a linear kernel outperformed all other models. Recognizing the potential limitations, future research involving expert data labelers specialized in sentiment analysis or NLP is needed to prevent the biases in sentiment labeling. Moreover, further exploration of deep learning techniques, such as Long short-term memory networks (LSTM) and Artificial Neural Networks (ANN) is needed to enhance the analysis and achieve more accurate sentiment predictions.

**Keywords**—Sentiment Analysis, Machine Learning, Web-scraping, Data Mining

## I. INTRODUCTION

In recent years, social media has experienced rapid growth, with more people worldwide using it as a medium to connect. Social media acts as the primary source of information, especially in the tourism sector [1]. People tend to look up on social media for some reviews, experiences and thoughts on certain tourist places. A high rating reviews change the image perception of the place compared to the low rating view. However, the low rating does impact the image of the place. [2]

Among many areas of natural language processing (NLP), sentiment analysis is one of the most widely used. A sentiment analysis study analyzes large amounts of textual data to understand what people think about a particular topic [3]. It is also a valuable tool that enables businesses to comprehend the connection between human emotions or opinions and natural language text. It helps to assess how people view a particular object, which holds significant importance to the creators or producers of the object. It is used to determine the emotional tone of a text, such as positive, negative or neutral [4]. This information can be used to understand customer sentiment, monitor brand reputation, and make better business decisions. In addition, exploring beyond just sentiment can uncover valuable contextual information. For instance, identifying the

reasons behind the positive or negative sentiment can be immensely beneficial in refining marketing strategies, improving products or services, and lead the business.

## II. PROBLEM STATEMENT

As a popular tourist attraction, Aquaria KLCC receives numerous review and feedback from visitors worldwide, offering a valuable source of information for both potential visitors and the management of Aquaria KLCC. The challenge lies in efficiently analyzing this large amount of textual data to understand visitor sentiments and opinions accurately. This data-driven approach can help identify areas for improvement, understand visitor preferences, and enhance the overall visitor experience. Thus, overall sentiment analysis of Aquaria KLCC using various machine learning models need to be done by conducting this study. Hence, insights into the strengths and weaknesses of the aquarium's based on the past visitors can be obtained.

## III. RESEARCH QUESTIONS AND RESEARCH OBJECTIVE

### A. Research Questions

- How to extract and pre-process the review data of Aquaria KLCC that is suitable for sentiment analysis?
- What machine learning techniques can be used to model?
- How can the performance of different techniques be systematically evaluated and compared?

### B. Research Objectives

- To extract and preprocess review data of Aquaria KLCC to create a suitable dataset for sentiment analysis.
- To explore and select appropriate machine learning techniques for sentiment analysis that can effectively model the review data of Aquaria KLCC.
- To systematically evaluate and compare the performance of different machine learning techniques for sentiment analysis on the Aquaria KLCC review dataset.

## IV. METHODOLOGY

In this study, CRISP-DM workflow is applied, ensuring a structured and systematic approach to sentiment analysis of Aquaria KLCC reviews.



Figure 1 CRISP-DM Phases [5]

### A. Project Understanding

Aquaria KLCC is an impressive and modern aquarium situated in Kuala Lumpur, Malaysia. It is conveniently located below the Kuala Lumpur Convention Centre, close to the famous Petronas Twin Towers. Covering an extensive area of 60,000 square feet, the aquarium offers a captivating experience with approximately 5,000 exhibits featuring both land-based and aquatic creatures.

TripAdvisor, being a popular website for tourists worldwide, provides a platform for visitors to share their experiences openly. This fosters a sense of community among users, where they can share their thoughts and honest reviews. Potential visitors can access this wealth of information to make informed decisions when planning their own trips to Aquaria KLCC. Analyzing data from Aquaria KLCC reviews on TripAdvisor can indeed be a valuable source of information for both potential visitors and the management of the aquarium. By studying the reviews and ratings, the management can gain insights into the strengths and weaknesses of the aquarium's offerings and services. This data-driven approach can help identify areas that need improvement, understand visitor preferences and expectations, and ultimately enhance the overall visitor experience.

By closely monitoring the reviews and ratings, the management can identify common complaints, and positive feedback. For instance, the management might learn if certain exhibits are particularly popular or if there are aspects of the aquarium experience that need attention, such as crowd management, cleanliness, or customer service. Additionally, tracking sentiment analysis from reviews can help the overall visitor satisfaction and identify areas that are consistently well-received or consistently underperforming. This data-driven approach enables the management to prioritize improvements that will have the most significant impact on the visitor experience. Furthermore, by comparing Aquaria KLCC's reviews and ratings with those of other tourist attractions, both locally and globally, especially with the same attraction such as the aquarium, the management can benchmark their performance.

Overall, leveraging the data from Aquaria KLCC's reviews on TripAdvisor can be a powerful tool for continuous improvement and attracting more visitors to the aquarium. By actively listening to the feedback of their visitors, the management can adapt to meet the expectations of tourists,

ensuring a memorable and enjoyable experience for everyone especially after the COVID era.

### B. Project Understanding

The data understanding phase of this study involves leveraging the powerful web scraping software tool, WebHarvy, to efficiently extract data from the TripAdvisor webpage, specifically focusing on the reviews of Aquaria KLCC. As a good tool for data extraction, WebHarvy offers users a versatile platform capable of navigating through diverse types of websites and collecting data in various formats, such as CSV, Excel, and JSON.

The target URL for this data extraction task is the TripAdvisor page dedicated to the Aquaria KLCC attraction. TripAdvisor is a widely acclaimed platform where travellers and visitors share their experiences, impressions, and opinions about various destinations and attractions, making it an ideal source for gathering valuable insights about Aquaria KLCC.

With WebHarvy's user-friendly interface, the data extraction process becomes remarkably straightforward, enabling efficient retrieval of the information. The software's intuitive point-and-click feature empowers users to effortlessly select specific items. The data set collected from the WebHarvy comprises an impressive total of 1100 reviews including with the review titles and the corresponding usernames from the TripAdvisor page dedicated to Aquaria KLCC.

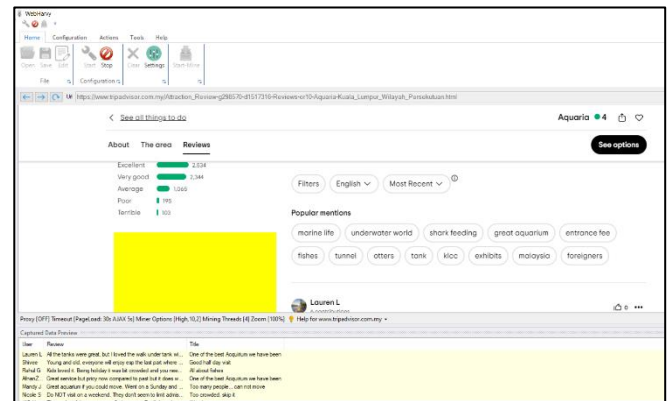


Figure 2 Data Mining Process using WebHarvy

### C. Data Preparation

Data pre-processing is a critical step in the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology to prepare the data for analysis and modeling. In the context of sentiment analysis for reviews in this study, the following steps were taken using Python and Orange data mining software:

- **Sentence Separation:** The initial data contained reviews with multiple sentences. To analyze each sentence independently, the reviews were separated into one sentence per row. This step resulted in a huge number of sentences, hence only 230 rows, with each row representing a single sentence, were used in this study.
- **Sentiment Labelling:** For sentiment analysis, the sentences need to be labeled as either positive or negative. This labeling was done manually, where the sentences are classified according to their sentiment.

**Tokenization:** Tokenization is the process of splitting sentences into individual words or tokens. In this case, word punctuation was used as a basis for tokenization. Punctuation marks were treated as separate tokens rather than being part of the words. For example, "Hello!" would be tokenized as ["Hello", "!", "."].

- **Target Variable Selection:** To train a sentiment analysis model, a target variable needs to be set. In this case, the sentiment labels "POSITIVE" and "NEGATIVE" were used as the target variable. The "Select Column" operator is employed to specify the attribute containing the sentiment labels as the target variable.

absolutely impressed  
weekends  
hokkaido inside distancing should animals  
definitely recommend ventilation never working  
health photos our just well being ticks  
found amount exhibits tanks sharks time activity  
walk way than experience worthy stuffy young  
adults family visitors but this really many other things  
pricey us price life staff with not that about some social  
overall few too have be if we at kids it over crowded  
got ger place of to was they  
without had like an aqua aquarium visit enjoyed  
displays variety fish for in there hot bit or must  
whole city funy all no good and it there old lack  
then kicc had hour see were very aquarium very specialy  
highes people youn and it there as up from here world gift  
creatures enjoy enjoyed great Visita aquaria can crowded need  
packed feel large quite friendly would make something  
species huge etc them under informative interesting situation  
available unfortunately underwater aquariums  
impressive disappointing

description uncomfortable  
 spectacular impossible little disappointment conditioning  
 environment different outside educational centre  
 knowledge true disappointing  
 notice attraction Venetian had price sweat frustrate  
 interesting experience tunnel/sight walk stop stroller  
 job feed animal variety exhibit creature beautiful  
 sense kick distance ticket life hour friendly money better waste  
 matter touch photo staff experience situation function  
 habitat whole swim ray fish star kind lack pretty/porn  
 guess cheap big small visit place holiday pose  
 empty know family shark visit place child adult helpful  
 nothing water world tank aquarium kid/lo move art help  
 glass/heat fun family tank aquarium kid/lo move art help  
 learn/tear give many fish sea nice kill crowd love amazing reward  
 cool shop limit pick city great good price absolutely  
 best special impress old like great good price absolutely  
 control impressive old like great good price absolutely  
 upstairs specific informative/buy enjoy people marine activity busy appear  
 maintain gift spend/click recommend large health heart display return love pretty  
 anyone water weekend plan huge underwater especially encounter scary  
 please friend working young expensive overprice  
 entrance language understand entertainment  
 explore including screaming common  
 interactive various

The wordcloud before the data preprocess generated was observed that the word cloud includes commas and stopwords, which are positioned at the center of the cloud, indicating higher in frequency within the corpus. It can be distracting and may not contribute meaningfully to the overall analysis. Thus, the data pre-process is needed. The wordcloud in Figure 2 is after the data pre-processing.

In the context of sentiment analysis for reviews, the Bag of Words (BoW) model was used to generate a corpus with word counts for each data during the feature extraction phase. By focusing only on word frequency and ignoring grammatical and word order factors, the BoW model displays text data as a set of individual words.

### E. Modelling

Based on the findings from previous studies, it is evident that employing SVM, RF, and LR with bag-of-words as the feature extraction method has been successful in sentiment analysis of movie reviews [6]. Additionally, LR has demonstrated exceptional performance, achieving an impressive accuracy of 93% [7], while another study achieved a notable accuracy of 78.88% using LR [8] and SVM also showed good performance [9].

Considering the success of LR, SVM, and RF in sentiment analysis of restaurant reviews with accuracies surpassing 85% [10], and the recognition of LR's outperformance over other models [11], it would be good to use these models in this study of sentiment analysis. By employing a combination of LR, SVM, and RF along with bag-of-words as the feature extraction method, this study can benefit from the well-established studies of these models in effectively classifying sentiment in textual data.

In the modeling phase, three major different machine learning algorithms were employed for sentiment analysis with different parameter which were Support Vector Machine (SVM) with linear, radial basis function (RBF), and polynomial kernels, followed by Logistic Regression with none regularization, Lasso regularization, and Ridge regularization, and finally, Random Forest with 10, 20, and 30 trees. These algorithms were trained and evaluated on the pre-processed data to build sentiment analysis models. The choice of diverse algorithms and hyperparameters allows for a comprehensive exploration of different modeling approaches, aiming to select the best-performing model for accurately classifying the sentiment of reviews.

- **Support Vector Machine (SVM):** SVM is a powerful supervised machine learning algorithm used for classification tasks. Its main objective is to discover the optimal hyperplane in a high-dimensional feature space that effectively separates data points of different classes. In sentiment analysis, SVM is applied with various kernels, including linear, RBF, and polynomial, to classify reviews into positive or negative sentiments based on the extracted features from the Bag of Words model. The linear kernel is suitable for linearly separable data, while the RBF and polynomial kernels are more capable of handling non-linearly separable data.
- **Logistic Regression (LR):** Logistic Regression is a popular linear classification method utilized to model the probability of binary outcomes, such as positive or negative sentiment in this case. By estimating feature

coefficients, the algorithm calculates the likelihood of a review belonging to a particular class. This study incorporates three variations of Logistic Regression models: the default one, along with L1 regularization (Lasso), and L2 regularization (Ridge). Regularization is employed to mitigate overfitting and to identify crucial features, with Lasso having a higher probability

- of generating sparse feature weights where some features may have zero weights.
- **Random Forest (RF):** Random Forest is an ensemble learning method that builds multiple decision tree classifiers and combines their predictions to make a final decision. Each decision tree is trained on a random subset of the data and features, reducing overfitting and improving model generalization. In this study, three different Random Forest models are built with varying numbers of trees (10, 20, and 30) to evaluate how the model's performance varies with the number of trees used in the ensemble.

TABLE I. MODELS & PARAMETER TESTED

Model	Parameter
Logistic Regression	Regularization = None
Logistic Regression (Ridge)	Regularization = Ridge
Logistic Regression (Lasso)	Regularization = Lasso
SVM (Linear)	Kernel = Linear
SVM (RBF)	Kernel = RBF
SVM (Polynomial)	Kernel = Polynomial
Random Forest (10)	Number of trees = 10
Random Forest (20)	Number of trees = 20
Random Forest (30)	Number of trees = 30

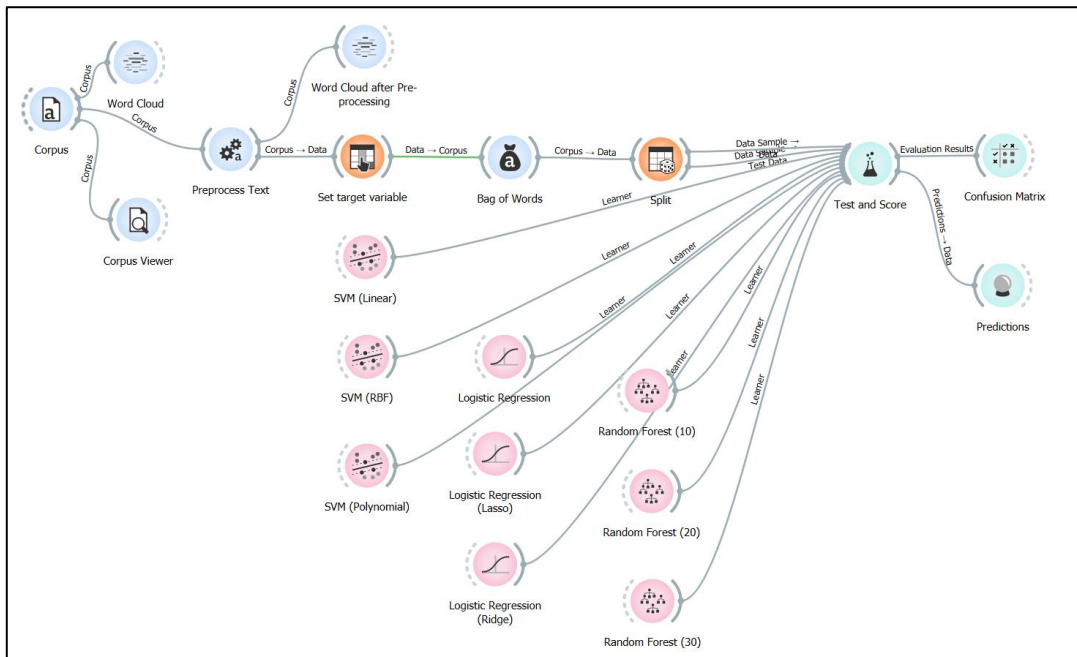


Figure 5 The Workflow of this Study



## F. Evaluation

In this study, various evaluation metrics were utilized to assess the performance of the predictive models. These metrics included Accuracy, Recall, Precision, AUC, F1-Score and Matthews correlation coefficient (MCC). Each metric served a specific purpose in measuring different aspects of the model's predictive capabilities.

$$Accuracy = \frac{TP+TN}{TP+FP+TN+FN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

$$F1-Score = \frac{2 \times Precision \times Recall}{Precision+Recall} \quad (4)$$

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (5)$$

## V. RESULTS & DISCUSSIONS

### A. Performance Evaluation

TABLE II. PERFORMANCE EVALUATION

Model	AUC	CA	F1	Prec	Recall	MCC
Logistic Regression	0.81	0.696	0.694	0.699	0.696	0.394
Logistic Regression (Ridge)	0.805	0.71	0.71	0.71	0.71	0.42
Logistic Regression (Lasso)	0.584	0.609	0.553	0.732	0.609	0.323
SVM (Linear)	0.773	0.725	0.725	0.725	0.725	0.45
SVM (RBF)	0.495	0.565	0.545	0.575	0.565	0.137
SVM (Polynomial)	0.494	0.551	0.533	0.556	0.551	0.104
Random Forest (10)	0.653	0.58	0.579	0.581	0.58	0.161
Random Forest (20)	0.681	0.609	0.609	0.609	0.609	0.217
Random Forest (30)	0.722	0.681	0.676	0.69	0.681	0.37

In comparing the performance of the models (Logistic Regression, SVM, and Random Forest), this study analyzed the results based on various evaluation metrics such as Area Under the Curve (AUC), F1-score, Precision, Recall, and Matthews Correlation Coefficient (MCC) but most importantly the Classification Accuracy. The higher the value of these metrics, the better the model's performance.

Among the nine evaluated models, SVM (Linear) emerges as the most favorable model based on various performance metrics. With a Classification Accuracy (CA) of

0.725, SVM (Linear) achieves the highest percentage of correctly classified instances. Additionally, its balanced F1-score of 0.725 reflects a good trade-off between precision and recall, ensuring both true positives and false positives are minimized effectively. SVM (Linear) also excels in precision and recall, both being 0.725, indicating its capability to make accurate positive predictions while avoiding false predictions. Furthermore, SVM (Linear) exhibits an impressive Area Under the Curve (AUC) value of 0.773, signifying its superior ability to distinguish between positive and negative samples. Lastly, SVM (Linear) achieves the highest Matthews Correlation Coefficient (MCC) of 0.45, emphasizing its overall performance in capturing true positive and true negative rates.

However, SVM (Polynomial) is the worst-performed model. It has the lowest AUC of 0.494 and Classification Accuracy (0.551), indicating lower correct classifications compared to the other models. Its F1-score is 0.551 while both precision (0.533) and recall (0.556) are relatively low, indicating challenges in correctly predicting positive predictions. The Matthews Correlation Coefficient (MCC) is also low at 0.104, indicating overall poor performance. Based on these metrics, SVM (Polynomial) is the worst model in this study.

Taking all these metrics into account, SVM (Linear) stands out as the best model for the classification prediction of sentiment analysis. It demonstrates consistent and good performance across various evaluation criteria, providing robust and reliable predictions.

### B. Sentiment Predictions

The output computed in the confusion matrix are tabulated for each machine learning model, where the total of true positives and true negatives are included. Table below shows the results of sentiment prediction for all models.

TABLE III. SENTIMENT ANALYSIS RESULT

Data	Number of Sentences	Total	
		Positive	Negative
Manual Data	Total: 230 Train: 161 Test: 69	Total: 117 Train: 82 Test: 35	Total: 113 Train: 79 Test: 34
Random Forest			
1st attempt (10)	Test: 69	32	37
2nd attempt (20)	Test: 69	36	33
3rd attempt (30)	Test: 69	43	26
Logistic Regression			
1st attempt (Normal)	Test: 69	40	29
2nd attempt (Lasso)	Test: 69	10	59
3rd attempt (Ridge)	Test: 69	35	34
Support Vector Machine			

Data	Number of Sentences	Total	
		Positive	Negative
1st attempt (Linear)	Test: 69	34	32
2nd attempt (RBF)	Test: 69	49	20

Data	Number of Sentences	Total	
		Positive	Negative
3 <sup>rd</sup> attempt (Polynomial)	Test: 69	48	21

Below is the snippet of predictions table for each model:

Class	Review	SVM (Linear)	SVM (RBF)	SVM (Polynomial)	andom Forest (1C	andom Forest (2C	andom Forest (3C	ogistic Regression	stic Regression (Lr	stic Regression (Ri
POSITIVE	Very informativ...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE
POSITIVE	We were both r...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE
NEGATIVE	We had difficult...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	POSITIVE	NEGATIVE	NEGATIVE
NEGATIVE	Ventilation was ...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	I would like to L...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	Today we went ...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE
POSITIVE	Experience the ...	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE
POSITIVE	The place is ver...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
POSITIVE	Impressed with ...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE	NEGATIVE	POSITIVE
NEGATIVE	Do NOT visit on...	NEGATIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	I'm reviewing t...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
POSITIVE	Worth the price.	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	Very bad on th...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	The amusemen...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
POSITIVE	The line for tick...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	We don't need ...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	They don't see...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE
POSITIVE	Tickets pretty c...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE
NEGATIVE	Unfortunately t...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE
POSITIVE	It was the first t...	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE	NEGATIVE	POSITIVE	POSITIVE	POSITIVE
NEGATIVE	The whole aqua...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE
NEGATIVE	The place was o...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	POSITIVE	NEGATIVE	NEGATIVE
POSITIVE	Good place to ...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE
POSITIVE	Really broad ra...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE
NEGATIVE	Very crowded o...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
NEGATIVE	What I get for p...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE	NEGATIVE	POSITIVE
POSITIVE	Very interesting...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE
NEGATIVE	Photos a bit pri...	NEGATIVE	POSITIVE	POSITIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE	NEGATIVE
POSITIVE	Good place for ...	POSITIVE	NEGATIVE	NEGATIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE
POSITIVE	You will be able...	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	POSITIVE	NEGATIVE	POSITIVE

Figure 6 Prediction of the Polarity (Positive & Negative)

## VI. CONCLUSION

All in all, this study focused on sentiment analysis of Aquaria KLCC based on reviews collected from TripAdvisor. The study followed the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology, which involved several stages, including project understanding, data understanding, data pre-processing modeling, and evaluation.

The dataset of the reviews was self-extracted from TripAdvisor using WebHarvy, and the necessary pre-processing steps were performed using Python and Orange. Nine different models were employed to predict the sentiment of the reviews, and their performance was compared to identify the most accurate one. Interestingly, the study found that the SVM (linear) model is the best model, outperformed the other eight models in all performance metrics.

As for the limitation, the manual labeling process of sentiment data into positive and negative can introduce biases due to the subjectivity of the researcher. The researcher may have different interpretations and criteria for labeling sentiments, resulting in inconsistencies and potential misclassifications as it might be influenced by their personal beliefs, experiences, or emotions, leading to unintended biases in the labeled dataset.

In future research, to address this limitation, hiring professional data labelers who are experts in sentiment analysis or natural language processing can help ensure a more objective and consistent sentiment labeling process. Additionally, it is highly recommended to delve into the exploration of diverse machine learning models while simultaneously harnessing the power of deep learning techniques such as Long short-term memory networks (LSTM) and Artificial Neural Network (ANN).

## ACKNOWLEDGMENT

I would like to express my sincere gratitude to all those who have contributed to the successful completion of this study, especially to my lecturer, PM Dr. Noraini Seman for her invaluable guidance and support throughout the duration of this study.

## REFERENCES

- [1] X. Guo and J. Pesonen, "The role of online travel reviews in evolving tourists' perceived destination image," *Scandinavian Journal of Hospitality and Tourism*, vol. 22, no. 4–5, pp. 372–392, Aug. 2022, doi: 10.1080/15022250.2022.2112414.
- [2] Clerk Maxwell, *A Treatise on Electricity and Magnetism*, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [2] "Why online reviews matter for travel & tourism brands," *Union*. <https://union.co/articles/importance-of-reviews-for-travel-and-tourism-brands>

- [3] N. Raj, "Starters Guide to Sentiment Analysis using Natural Language Processing," *Analytics Vidhya*, Jun. 2023, [Online]. Available: <https://www.analyticsvidhya.com/blog/2021/06/nlp-sentiment-analysis/>
- [4] "What is Sentiment Analysis? - Sentiment Analysis Explained - AWS," *Amazon Web Services, Inc.* <https://aws.amazon.com/what-is/sentiment-analysis/#:~:text=Sentiment%20analysis%20is%20the%20process,so%20cial%20media%20comments%2C%20and%20reviews.>
- [5] N. Hotz, "What is CRISP DM? - Data Science Process Alliance," *Data Science Process Alliance*, Jan. 19, 2023. <https://www.datascience-pm.com/crisp-dm-2/>
- [6] B. Sangeetha, S. Sangeetha, D. T. Goutham, and V. R. N., "Sentiment Analysis on Movie Reviews: A Comparative Analysis," 2023 International Conference on Intelligent Systems for Communication, IoT and Security (ICISCoIS), Feb. 2023, doi: 10.1109/iciscois56541.2023.10100367.
- [7] R. F. Ramadhan, P. H. Gunawan, and N. Aquarini, "Web-Based Sentiment Analysis Application of Hotel Reviews in Indonesia," 2022 2nd International Conference on Intelligent Cybernetics Technology & Applications (ICICyTA), Dec. 2022, doi: 10.1109/icicyta57421.2022.10037946.
- [8] S. Hemalatha and R. Ramathmika, "Sentiment Analysis of Yelp Reviews by Machine Learning," 2019 International Conference on Intelligent Computing and Control Systems (ICCS), May 2019, doi: 10.1109/iccs45141.2019.9065812
- [9] R. Uma, A. S. H, P. Jawahar, and B. V. Rishitha, "Support Vector machine and convolutional neural network approach to customer review sentiment analysis," 2022 1st International Conference on Computational Science and Technology (ICCST), Nov. 2022, doi: 10.1109/iccst55948.2022.10040381.
- [10] K. Zahoor, N. Z. Bawany, and S. Hamid, "Sentiment Analysis and Classification of Restaurant Reviews using Machine Learning," 2020 21st International Arab Conference on Information Technology (ACIT), Nov. 2020, doi: 10.1109/acit50332.2020.9300098.
- [11] N. U. Saaqib, N. Gunika, and H. K. Verma, "Analysis of Sentiment on Amazon Product Reviews," 2023 Third International Conference on Secure Cyber Computing and Communication (ICSCCC), May 2023, doi: 10.1109/icseccc58608.2023.10176787