
Machine Learning-Based Social Media Review Analysis for Recommending Tourist Spots

Prakash Lahagun, Bidur Devkota, Sakchham Giri, Parbat Budha

Gandaki College of Engineering and Science, Pokhara University, Nepal

(Manuscript Received 10/09/2023; Revised 11/11/ 2023; Accepted 01/12/ 2023)

Abstract

In recent years, the tourism industry has witnessed significant growth, resulting in an increased demand for effective and personalized tourist place recommendation systems. In this study, a tourist spot recommendation system is proposed which is built by developing a machine learning model based on a Support Vector Machine (SVM), Decision Tree (DT), and k-Nearest Neighbors(k-NN). Public experiences and opinions regarding the various spots available in popular social media sites such as TripAdvisor, Google, Instagram, and TikTok are utilized to train the model. The system matches the probability of the user query with the predicted probability of reviews for a particular spot. The SVM algorithm, known for its robustness in handling high-dimensional data, is adapted to model the complex relationships between users' reviews, spots, and their attributes. Real-world data is used to evaluate the system's performance, demonstrating its ability to significantly improve the user experience and contribute to the sustainable growth of the tourism sector. The system's capability was demonstrated as it achieved a notable F1-Score of 0.78 when SVM was implemented. Additionally, a promising accuracy rate of 93.023% was observed when random queries were used for tourism spot prediction, emphasizing that SVM outperformed DT and k-NN.

Keywords: Machine Learning; Recommendation System; Social Media data analysis; Support Vector Machine; Tourism Industry

1. Introduction

The global tourism industry is undergoing a transformative surge, driven by heightened connectivity, increased accessibility, and a growing appetite for unparalleled travel experiences. Social media platforms have played an indispensable role for people to discuss and convey their experiences and opinions regarding various aspects of life. Such platforms cater to a massive amount of data useful for the travel and tourism industry [2],[3],[8],[14]. The United Nations has recognized tourism as a major economic activity that can make a significant contribution to achieving Sustainable Development Goals (SDGs) [15]. This study proposed a machine learning model, e.g. SVM, DT, k-NN, and by utilizing the user experiences scattered across various social media platforms to recommend tourism spots to visitors. The effectiveness of SVM as a learning method, in comparison to its deep learning counterparts, is underscored by its robust performance even when confronted with limited data and fewer resources [1]. DT characterized by its

hierarchical decision-making process involving recursive dataset partitioning based on particular features [17], is deemed to be an influential model. Additionally, the k-NN algorithm, which was initially formulated by Evelyn Fix and Joseph Hodges in 1951 [16], is recognized as a non-parametric supervised learning technique. The proposed model trains the SVM classifier using the public reviews of various tourist spots. The model thus developed matches the user's search query and recommends those tourist spots which has the higher probability. In this study, we demonstrate the working of the proposed methodology by utilizing social media data such as TripAdvisor, Google, Instagram, and TikTok related to the tourist spots in and around Pokhara, a popular tourist destination in Nepal.

1.1 Problem-Solution Approach

Conventional sources such as tourism service providers and hotel websites might present biased information due to their inherent bias in promoting

their services. Furthermore, the content on these platforms can become outdated over time, as it may not receive consistent updates to reflect changing circumstances or customer feedback. The reliability of information from traditional sources like tourism service providers and hotel websites can be supplemented by public reviews from actual visitors and multiple data sources. Public reviews provide firsthand experiences, while diverse data sources contribute to a more comprehensive and unbiased understanding of a place. This approach ensures a more accurate depiction of the destination.

While seeking a spectacular view of "Machhapuchhare" in Pokhara, tourists often encounter a dilemma. While "Lakeside" offers a glimpse, "Sarangkot" is often favored for a superior sight. This choice between accessibility and a better view can be a source of difficulty for travelers. The challenge of itinerary planning in Pokhara, especially when seeking a balance between accessibility and a better view of "Machhapuchhare," can be addressed through our method. We recommend the nearest location, such as "Lakeside," for those who desire both lake and mountain views. By offering recommendations for the top 5 and top 10 spots, we aim to simplify the decision-making process, ensuring travelers can enjoy their preferred experiences without unnecessary complexity.

In our implementation, Term Frequency-Inverse Document Frequency (TF-IDF) was employed to analyze text and ascertain the significance of terms within a collection of reviews. Additionally, SVM, DT, and k-NN were utilized to analyze user reviews, thereby extracting valuable insights. This approach facilitated the generation of spot recommendations by harnessing user-generated content and feedback, subsequently contributing to the overall enhancement of the recommendation system's accuracy.

2. Related Works

The approach we introduced can construct spot probability descriptions by utilizing text from manually collected sources such as TripAdvisor¹, Google², TikTok³, and Instagram⁴. Many researchers have conducted studies in the past to understand the meaning of a place by using different sources of information.

As an illustration, Jiang and their research team [3]

undertook a project titled "Travel recommendation through author-topic model-based collaborative filtering." They utilized text data from social media to discern people's preferences in recommending tourist destinations. They came up with a way to suggest interesting places to social media users, calling it the Author Topic Collaborative Filtering (ATCF) method. This study introduces a new way to plan travel routes. It looks at the main themes of places to visit and what makes them special. They collected travel information from the internet, organized it, and used location details to create travel paths. By considering both themes and special features, their method helped pick out the best travel routes from lots of travel notes. They tested it, and it worked well [4]. They use Point of Interest (POI), topic modeling-based collaborative filtering. It uses social media data like Flickr geotags text. It will provide output in the user experience. The limitation exists in the case where no POI is discovered, as geotagged Flickr photos within the region are reviewed, and the photo tag with the highest TF-IDF is used to designate the spot, though it may not always be representative.

Another type is the approach presented is designed to generate descriptions for popular tourist areas of interest (TAOIs) by utilizing text from Tweets and Flickr. One challenge of it is that sometimes there isn't enough information available to describe these spots. In places with limited social media presence, the approach of encouraging users to keep their posts short may not work as effectively on platforms like Flickr and Twitter [2].

Veronika Arefieva and her team conducted research on destination Instagram images [5]. They gathered a substantial collection of Instagram images and applied techniques such as k-NN, topic modeling, and correlation analysis. These methods improved topic modeling and k-NN optimization but had limitations in correlation metrics. Additionally, the study solely relied on Instagram images, potentially lacking comprehensive place-related data for future research on authentic tourist experiences.

Our proposed model matches the user's search query and recommends tourist spots based on their higher probability of relevance. Our data collection spans multiple sources, significantly boosting prediction accuracy. The inclusion of diverse social media user data minimizes the risk of bias in our analysis, contributing to a more robust and reliable model.

¹<https://www.tripadvisor.com/>

²<https://www.google.com/>

³<https://www.tiktok.com/>

⁴<https://www.instagram.com/>

3. Methodology

Figure 1 illustrates the overall methodology employed in this research. We collected data manually from various platforms and emphasized data consistency. During preprocessing, we identified and filtered out non-essential languages like Japanese, Nepali, and Indonesian. Additionally, we removed emojis, stop words, and special characters, and converted text to lowercase. To further enhance text quality, we employed lemmatization and tokenization techniques. In our utilization of machine learning algorithms, we adopt a multi-class classification methodology in which spot names are designated as distinct classes, leveraging a dataset encompassing 37 unique spots. The SVM employs hyperplanes for text classification, while DT employs decision boundaries, constructing a hierarchical tree with internal nodes representing decisions among multiple classes. In parallel, the k-NN algorithm discerns class separation through a majority vote mechanism among the closest neighbors. This integrated approach forms a comprehensive framework for text classification within our system. When a user submits a query, our model classifies it to generate tourist spot recommendations by assessing query relevance and processing data. This methodical procedure guarantees better spot suggestions, thereby augmenting the user experience and facilitating informed decision-making throughout the travel planning process. This meticulous approach entails the utilization of algorithms and data analysis, resulting in a seamless and user-centric interface. Through this mechanism, users benefit from tailored and reliable travel recommendations, ultimately enhancing their overall travel planning experience and ensuring optimal destination choices.

Table 1 provides a breakdown of review counts for different spots in the selected social media platforms. Specifically, a total of 8,810 reviews from 11 spots were collected from Google.

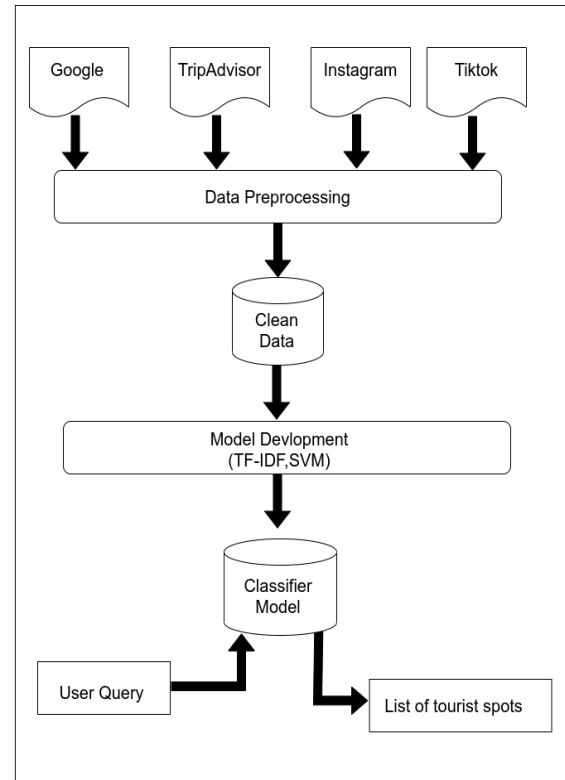


Figure 1: Overall method to recommend tourist spots

Furthermore, 10,000 reviews from 17 spots were gathered from TripAdvisor, and a combined count of 8,341 reviews from 9 spots were collected from Instagram and TikTok. The "Total" row consolidates all the review counts from these platforms, indicating a sum of 27,151 reviews in this dataset. These numbers are valuable as they represent user-generated feedback or reviews, which can be instrumental in developing and enhancing recommendation systems for various tourist spots. Analyzing this data can help improve personalized recommendations for users across these platforms.

Table 1. Data was collected for Pokhara and nearby areas.

Sites	Review Count	Number of Spots
Google	8,810	11
TripAdvisor	10,000	17
Instagram and TikTok	8,341	9
Total	27,151	37

3.1 Term Frequency (TF) and Inverse Document Frequency (IDF)

TF-IDF was introduced by Karen Sparck Jones in 1972 [13], is a crucial tool in natural language processing. TF-IDF is a quantitative metric used in natural language processing to assess the importance

of a word within a particular document relative to its relevance across a corpus of documents. It's particularly valuable for analyzing user reviews extracted from multiple platforms. It gauges word importance by considering both frequency within a document and rarity across all documents. This yields TF-IDF scores, where high scores signify words that are both frequent and unique, carrying significant meaning or context [9]. The TF-IDF scheme is employed as the weighting scheme.

$$TF = N_d / T_d \quad (1)$$

In the given Eq. (1), 'N_d' represents the frequency of occurrence of the word 'w' within a specific document, while 'T_d' denotes the total word count in that document.

IDF is a statistical measure that assesses the importance of a word within a corpus of text. IDF is calculated by taking the logarithm of the total number of documents in the corpus divided by the number of documents that contain the word in question. This means that words that appear in more documents will have a lower IDF score, while words that appear in fewer documents will have a higher IDF score. The rationale behind IDF is to downplay the significance of common words. Common words like "of", "in", "the", "is", "that", and "it" appear frequently in many documents, but they carry little or no meaningful information. By giving lower IDF scores to common words, we can focus our attention on the words that are more likely to be important to the document in which they appear.

$$IDF = \log(T_c / D_w) \quad (2)$$

From Eq. (2) the variables 'T_c' represents the total number of documents in the corpus, and 'D_w' represents the number of documents containing the word 'w'. The TF-IDF weight is determined using the following formula

$$TF-IDF \text{ Weight} = TF \times IDF \quad (3)$$

3.2 Support Vector Machine

SVM is a machine learning algorithm utilized for both classification and regression tasks. It constructs a hyperplane that optimally separates data points into distinct classes, maximizing the margin between them. Data points closest to the hyperplane, known as support vectors, influence the model's decision boundary. The parameter 'C' for regularization plays a pivotal role in balancing the trade-off between the

training error and the margin. Previous research often applied SVMs to text classification but overlooked their probability estimation functionality. This paper investigates SVMs in text classification, with a focus on leveraging their probabilistic outputs to improve result interpretation [10]. The concept of multi-class classification in SVM is employed [11], with classes designated based on spot names, resulting in 37 distinct categories such as Lakeside, AnnapurnaMountainRange, Rupalake, MardiHimal, BarahiMandir, and more. Class separation is essential in SVM for multi-class classification as it allows the algorithm to create clear boundaries between different classes, enabling accurate categorization of data points into their respective categories. For example, consider the review 'Amazing Nice Lake View and Temple,' with the class label 'Lakeside.' In this case, the SVM algorithm learns to find the optimal hyperplane that separates feature vectors into distinct classes, including 'Lakeside.' The positioning of this hyperplane by the SVM model is such that it effectively distinguishes reviews related to 'Lakeside' from other user queries. As a result, when a new user enters their query, the SVM model determines which side of the hyperplane the query falls on. An SVM model has been trained, and decision function scores are extracted. To convert the decision scores into probabilities, a logistic regression model is applied. After that, the probability of the text is calculated using the logistic regression formula [12].

$$P(y = 1/x) = 1 / (1 + \exp(a \times dfv + b)) \quad (4)$$

Here, 'dfv' represents the value of the decision function, while 'a' and 'b' denote parameters that are acquired through the training process of the logistic regression model.

3.3 Decision Tree

In machine learning, DT are constructed to facilitate data analysis, and data is recursively split into subsets based on feature values. Each split, guided by a chosen attribute, partitions the data into more homogenous groups. These splits continue until a predetermined stopping criterion is met. DT is extensively used for classification and regression tasks due to its interpretability and ability to handle both categorical and numerical data. They are constructed, pruned, and evaluated based on specific algorithms. By traversing the tree from a root to a leaf node we finally made a prediction, following the decision path [17]. Like in SVM, as described in

section 3.2, DT also employs multi-class classification by designating classes based on spot names, leading to the creation of 37 distinct categories. For example, a review designated as 'Lakeside' is effectively segregated from other categories. When a new user query is introduced, the classification in the DT model is determined based on which side of the decision boundary the query falls onto. Entropy and information gain are employed for the selection of the next attribute. Entropy is a measure of impurity or disorder in a set of data, and within the context of DT, it is utilized for quantifying the uncertainty or randomness connected with the class labels of the data. Information gain, on the other hand, is a measure of how much the entropy of a dataset is diminished subsequent to the utilization of a specific attribute for splitting the data. Its purpose is to gauge the effectiveness of an attribute in enhancing the purity of subsets. In the following Eq. (5), 'H(s)' represents entropy, and 'IG(s)' signifies information gain. Information gain quantifies the disparity in entropy between parent and child nodes, guiding the selection of the attribute with the highest information gain for the subsequent internal node.

$$H(s) = - \sum P_c * \log(P_c) \quad (5)$$

where, H(s) represents the entropy of a set or a probability distribution. 'P_c' represents the probability of a particular category 'c' occurring in the data.

$$IG(s) = H(s) - \sum_t P_t * H(t) \quad (6)$$

where H(s) represents the entropy of the original dataset or node 's' before the split. 'P_t' represents the probability of a particular outcome or subset 't' occurring as a result of the split. H(t) represents the entropy of each outcome or subset 't' after the split.

3.4 k-Nearest Neighbors

The k-NN algorithm is a versatile machine-learning technique utilized for both classification and regression purposes. In k-NN, data points are assigned to a specific class based on the majority class among their k-nearest neighbors in the feature space. This method involves measuring distances between data points, commonly using Euclidean distance. k-NN is categorized as a non-parametric, instance-based learning method because it doesn't rely on strong assumptions regarding the underlying data distribution. Instead, it relies on the local

characteristics of the data. k-NN's adaptability and simplicity make it valuable for various applications, although it can be sensitive to the choice of k and requires efficient algorithms for large datasets [16]. Much like in SVM, as outlined in section 3.2, multi-class classification is also implemented in k-NN. Classes are designated based on spot names, resulting in the creation of 37 distinct categories. For instance, a review labeled as 'Lakeside' is effectively isolated from other categories. When a new user query is introduced, the classification in the k-NN model is determined based on the majority vote among its nearest neighbors. In the k-NN algorithm, the Euclidean distance is a commonly used metric to measure the similarity or dissimilarity between data points. It helps determine which data points are the nearest neighbors to a given data point.

$$d(p,q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2} \quad (7)$$

Eq. (7) shows 'd (p, q)' denotes the Euclidean distance, which measures the distance between points 'p' and 'q' in 'n' dimensions, where 'p_i' and 'q_i' refer to the ith components of the vector's 'p' and 'q,' respectively.

4. Results and Discussion

Utilizing the proposed technique, we elucidate the tourist spots by analyzing social media data encompassing reviews and descriptions. This section is dedicated to showcasing the outcomes generated from the application of our method within our study area. We have judiciously selected several renowned tourist attractions not only within Pokhara but also extending to regions beyond the valley, encompassing Baglungkalika, Sikles, the Annapurna Mountain Range, Machhapuchhare, Rupakot, Mardi Himal, and PoonHill for in-depth exploration. In Table 2, a comprehensive depiction of the respective performances of diverse algorithms employed is presented for a thorough understanding.

Table 2. Evaluations of Different Algorithms

Algorithms Used	F1-Score
SVM	0.78
DT	0.63
k-NN	0.59

Table 2 highlights SVM's versatility in classifying both linear and non-linear data, emphasizing its selection for its lightweight nature and proficiency in categorizing reviews and identifying spot names. Despite imbalanced datasets, it delivered a respectable 0.78 F1-Score. In the experiment, a

regularization parameter (C) of 3 was employed, along with default settings for kernel, degree, and gamma. The dataset was split into training and testing sets with an 80:20 ratio, ensuring reproducibility by using a random seed of 100. Various C values (1, 2, 3, 5, and 10) were tested, revealing that moderate regularization (C=3) delivered optimal performance. This underscores that, for our dataset and task, a balance between underfitting and overfitting was achieved with C=3, surpassing other tested values.

In decision-making, we use a hierarchical decision support model known as a DT. This model utilizes a tree-like structure to represent decisions and their possible outcomes, including chance events, resource allocations, and utility considerations. In our analysis, a DT with default parameters was utilized, resulting in an F1-Score of 0.63. The k-NN algorithm operates on the assumption of similarity between the new data or case and the existing cases, placing the new case into the category most closely resembling the available categories. In our analysis, k-NN with default parameters was employed, resulting in an F1-Score of 0.59. When compared to SVM and Decision Tree, it is notable that k-NN exhibits the lowest F1-Score.

In the comparative analysis among SVM, Decision Trees, and k-NN, it is evident that SVM achieves the highest F1-Score, securely reaching 0.78, thereby yielding superior results. So that we employ the formidable SVM algorithm in our research.

Table 3. Performance of SVM model in Each Platform

Platforms Utilized	F1-Score
Google	0.86
TripAdvisor	0.87
Instagram and TikTok	0.53
Total	0.78

In Table 3, we meticulously gather diverse datasets from various platforms and clean the text. Applying SVM to different platforms yields varying F1-Score, with Google at 0.86, TripAdvisor at 0.87, and Instagram and TikTok at 0.53. To overcome the challenge of insufficient spots for recommendations, we aggregate data from all sources, achieving an overall 0.78 F1-Score. With only 17 classes on TripAdvisor, the classification task becomes less complex, as it necessitates distinguishing between a smaller number of categories. This circumstance potentially facilitates the attainment of a high F1-Score. However, for the recommendation system's comprehensive coverage, it is imperative to encompass all 37 classes, necessitating the amalgamation of data from TripAdvisor, Google, Instagram, and TikTok. Despite yielding a lower F1-Score than TripAdvisor individually, the integrated

dataset ensures the inclusion of all recommended spots, thus enhancing its overall utility and completeness. Our model adeptly predicts the top 5 recommended spots with a high level of reliability. However, it occasionally predicts irrelevant spots due to associations with unrelated keywords or user queries. In the future, to improve this, we plan to refine our approach by using SVM's probability method at a word level. Analyzing individual words enhances prediction precision, ensuring accurate alignment with user queries and spot name retrieval through query probability calculation.

Table 4. Samples of User Queries and their Evaluation

User Query	Probability	Expected Outcome	Desired Outcome
Lakeside	0.60	Lakeside	Yes
Sarangkot	0.97	Sarangkot	Yes
DikiDada	0.0014	Dikidada	No

Total: 37 spots

Correctly predicted (highest probability):33

Incorrect predictions: 4

Accuracy = $33/37 * 100$

Accuracy = 89.19%

Table 4 displays sample test results for our evaluation. In our dataset of 37 tourist spots, 33 spots were accurately predicted, while 4 places were predicted incorrectly. When specific queries were input by users, the system's ability to provide accurate recommendations could be hindered by the limited number of user reviews, often fewer than ten, for certain spots. These reviews are crucial for matching user query keywords with spot content. With few reviews, relevant keywords may not be prioritized by the system, potentially causing relevant spots to be omitted or ranked lower. This highlights the significance of review quantity in accurately assessing spot appeal and attributes. Insufficient reviews can lead to discrepancies in user recommendations, as spot value is challenging to gauge accurately. From Table 4, it can be observed that the probabilities of "DikiDada" are lower due to insufficient reviews. In contrast, "Lakeside" and "Sarangkot" achieved prominence due to abundant reviews and a higher probability of success. During evaluation, our model was subjected to testing by users, a process that yielded positive feedback. The mobile app, which was installed on a select group of users, was developed for evaluation purposes. Positive outcomes were observed in user engagement with the model. Favorable comments and reactions from users affirmed its effectiveness and user-friendliness. This feedback validated our development efforts and aligned the model with user expectations and needs.

Table 5. Samples of random User Queries and their Evaluation

User Query	Result	Desired Outcome
lake	BegnasTal Lakeside Rupalake Kaskikot Sikles	Yes
zipline	Lekhnath Sikles Kaskikot Sarangkot TibetanRefugeeCamp	No
Total query collected: 43 Correctly predicted: 40 Incorrect predictions: 3		
Accuracy = $40/43 * 100$ Accuracy = 93.023%		

We evaluated our model's performance with input from diverse users who generally found it effective. The mobile application was developed and installed on a randomly selected group of users locally, ensuring controlled access and evaluation. Table 5 displays outcomes when users input queries, generating recommendations. Satisfactory results were yielded by most queries, but some fell short due to limited user reviews for certain spots, thereby affecting the quality of recommendations, and causing the "cold start" problem. This caused the expected outcomes to either rank last or not appear among the top five recommendations. We collected 43 keywords from users during the test. To illustrate this, consider the results obtained for queries associated with keywords like "lake" and "zipline". For instance, when the keyword "lake" was used, the recommendations for spots like "BegnasTal" and "Lakeside" received scores of 0.49 and 0.47, respectively, indicating a relatively good match. It's worth noting that specific spots like "Rupalake," "Kaskikot," and "Sikles" received relatively lower scores, such as 0.023, 0.0016, and 0.0016, respectively. This occurred because when users entered the query "lake," the presence of the keyword "lake" within the reviews of these spots, contributed to their high ranking among the top 5 recommended spots. Similarly, this pattern continued with other keywords "zipline," where "Lekhnath" and "Sikles" scored 0.12 and 0.11, while "Kaskikot," "Sarangkot," and "TibetanRefugeeCamp" received scores of 0.10, 0.094, and 0.054, respectively. These variations underscored the varying degrees of relevance

associated with the presence or absence of reviews. The system failed to produce the desired outcome primarily because of the limited number of reviews, resulting in a lower weighting of specific keywords in the recommendation algorithm.

Table 6. Result Evaluation of Rupalake

User Query	Evaluation in the Presence of Limited Single-Source Reviews	Evaluation Using Multiple Source Reviews
Rupalake	BarahiMandir (0.47) ShantiStupa (0.18) Rupalake (0.11) GupteshworCave (0.08) Machhapuchhre (0.04)	Rupalake (0.81) Lekhnath (0.14) BegnasTal (0.009) Kaskikot (0.003) Pumdikot (0.003)

Table 6 illustrates that when low review quantities are available from a single source, the resultant accuracy diminishes. Notably, it portrays Barahi Mandir in the top position, whereas our intended outcome is to place Rupa Lake at the forefront. Barahi Mandir is assigned a probability of 0.47, followed by Shanti Stupa (0.18), Rupa Lake (0.11), Gupteshwor Cave (0.08), and Machhapuchhre (0.04). This discrepancy in results can be attributed to the limited number of reviews. However, when multiple sources are considered, the accuracy of the rankings significantly improves. Rupa Lake emerges as the top choice with a substantial probability of 0.81, followed by Lekhnath (0.14), Begnas Tal (0.009), Kaskikot (0.003), and Pumdikot (0.003). This approach, which incorporates reviews from various sources, not only yields promising results but also effectively addresses the challenge posed by a scarcity of reviews.

5. Conclusions

In conclusion, a Machine Learning-based Tourist spot recommendation System was introduced in this study, which harnessed real user experience and employed SVM to generate tailored and relevant recommendations. The likelihood of different tourist attractions at the location was assessed, and they were categorized based on their captured probabilities through the utilization of a support vector machine. The system's potential to enhance user satisfaction and engagement is highlighted by the empirical findings. During our study, it was observed that among the testing algorithms, a high

F1-Score of 0.78 was achieved by SVM. When the algorithm was tested on single-source data, a high F1-Score of 0.87 was obtained (TripAdvisor). However, when multiple-source data were aggregated, a F1-Score of 0.78 was achieved. Additionally, a promising accuracy of 93.023% was obtained when random queries collected from users were tested. Throughout this study, limitations were encountered, including a lack of data and a limited number of collected random queries (43).

In future research, a deeper exploration into the realm of deep learning may be pursued, with an emphasis on the collection of a more extensive dataset, facilitated through the accumulation of data derived from image location details. This endeavor is intended to address the challenges associated with mitigating the "cold start" problem and addressing data scarcity. Furthermore, the consideration of gathering an increased volume of random queries, potentially reaching 100 queries, from diverse sources of random users, is also incorporated into the research plan. A performance improvement is anticipated as a result of these measures.

References

- [1] Cortes C. and Vapnik V. Support-vector networks. *Machine learning*, 20 (1995) 273-97.
- [2] Devkota B., Miyazaki H. and Pahari N. Utilizing user generated contents to describe tourism areas of interest. *2019 First International Conference on Smart Technology & Urban Development (STUD)*, (2019) 1-6.
- [3] Jiang S., Qian X., Shen J. and Mei T. Travel recommendation via author topic model based collaborative filtering. *In MultiMedia Modeling: 21st International Conference, MMM 2015, Sydney, NSW, Australia, January 5-7, 2015, Proceedings, Part II 21*, (2015) 392-402.
- [4] Du S., Zhang H., Xu H., Yang J. and Tu O. To make the travel healthier: a new tourism personalized route recommendation algorithm. *Journal of ambient intelligence and humanized computing*, 10 (2019) 3551-3562.
- [5] Arefieva V., Egger R. and Yu J. A machine learning approach to cluster destination image on Instagram. *Tourism Management*, 85 (2021) 104318.
- [6] Upreti B. R., Upadhyaya P. K. and Sapkota T. Pokhara Tourism Council NC. Tourism in Pokhara: Issues, trends and future prospects for peace and prosperity, (2013).
- [7] Adhikari S. Prospects of tourism in Nepal: A study of Pokhara city. (2019).
- [8] Wenan T., Shrestha D., Gaudel B., Rajkarnikar N. and Jeong S. R. Analysis and Evaluation of TripAdvisor Data: A Case of Pokhara, Nepal. *In Intelligent Computing & Optimization: Proceedings of the 4th International Conference on Intelligent Computing and Optimization 2021 (ICO2021) 3*, (2022) 738-750.
- [9] Shrestha K. Comparative Analysis of TF-IDF and Word2vec Algorithm for Content-based Job Recommendation System. (2020).
- [10] Ilangovan P. Support Vector Machine based a New Recommendation System for Selecting Movies and Music. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, 12(10) (2021) 1425-1429.
- [11] Weston J. and Watkins C. Multi-class support vector machines, Technical Report CSD-TR-98-04, Department of Computer Science, Royal Holloway, University of London. (1998).
- [12] Shrestha K. Comparative Analysis of TF-IDF and Word2vec Algorithm for Content-based Job Recommendation System. (2020).
- [13] Platt J. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3) (1999) 61-74.
- [14] Sparck Jones K. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1) (1972) 11-21.
- [15] Devkota B., Miyazaki H., Witayangkurn A. and Kim SM. Using volunteered geographic information and nighttime light remote sensing data to identify tourism areas of interest. *Sustainability*, (2019) 11(17) 4718.
- [16] Hall C. M. Constructing sustainable tourism development: The 2030 agenda and the managerial ecology of sustainable tourism. *Journal of Sustainable Tourism*, 27(7) (2019) 1044-1060.
- [17] Fix E. and Hodges J. L. Discriminatory analysis. Nonparametric discrimination: Consistency properties. *International Statistical Review/Revue Internationale de Statistique*, 57(3) (1989) 238-247.
- [18] Breiman L., Friedman J. H., Olshen R. A. and Stone C. J. Classification and regression trees. Statistics/probability series. (1984).