



ISSN: 2454-9940



**INTERNATIONAL JOURNAL OF APPLIED
SCIENCE ENGINEERING AND MANAGEMENT**

E-Mail :
editor.ijasem@gmail.com
editor@ijasem.org

www.ijasem.org

A Deductive Decision Tree-Based Traveler Recommendation System

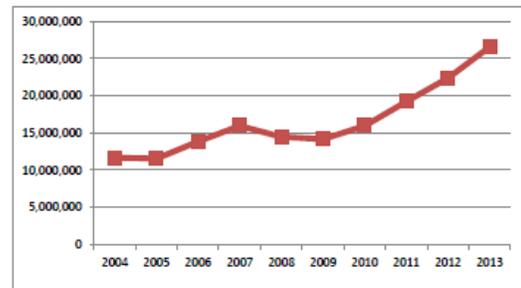
DR BIJAYA KUMAR NANDA

Abstract—

One of the most difficult things for travelers is deciding where to go on vacation from the plethora of information that can be found online and elsewhere. Whilst making preparations for a trip and while on the road. In the past, many Travel Recommendation Systems (TRSs) have sought to address this issue. But some of the technical factors, like system correctness, and the practical ones, like usability and pleasure, have been ignored. Novel models for the information-seeking behavior of visitors are necessary to solve this problem. In this study, we suggest an innovative, human-centric TRS to help travelers get around a strange city. We use a dataset we acquired from the actual world to balance the technical and practical considerations. Recommendations are generated using decision tree C4.5, and the system is constructed utilizing a two-step feature selection strategy to limit the amount of inputs. The testing findings demonstrate the effectiveness of the suggested TRS in making tailored recommendations for vacation spots that are sure to please.

I. INTRODUCTION

In 2013, tourism generated 9.5% of the world's GDP, making it a highly significant economic sector. Predictions for the travel industry are optimistic. Provide a GDP-boosting impact of almost 10.3 percent in 2023. When it comes to the economic impact of the travel and tourism industry, South East Asia is predicted to have the most expansion. Specific nations with the most appealing tourist attributes in 2013 were selected as Thailand, Indonesia, Singapore,



ore, and Myanmar [1].

PROFESSOR, Mtech, Ph.D
Department of CSE
Gandhi Institute for Technology, Bhubaneswar.

Figure 1: The Annualized Number of Overseas Visitors to Thailand, 2004–2013 [1]

The number of foreign visitors to Thailand has increased by a factor of 2. Within the previous decade and a half (See Fig 1). Thailand is the world's tenth most visited country in 2013[1]. The number of

brought in 1.79 trillion BHT (\$55.49 billion) in 2013[2]. Tourists today rely on the Internet more than any other source for learning about local businesses and their offerings [3]. Searching for places, often known as trip planning, may be overwhelming for visitors due to the vast amount of disparate information that is readily accessible online. The quality of attractions, travel routes, accommodations, numbers of travelers, leisure activities, weather, etc. are just a few examples of the numerous variables that must be considered while arranging a vacation. Technology, notably the Internet, has recently provided significant benefits to the tourist industry [5]. Decision-making technologies, also known as Recommendation Systems (RS), have made it easier

processes. It is also important to remove any doubts that may arise during the research phase of a tourist's decision-making process. Model complexity might be reduced by lowering the number of system parameters. As a result, the system's suggestion performance and user happiness may both improve.

II. BACKGROUND

A. Method for Making Suggestions

Recommendation systems (RSs) are a kind of decision support system (DSS) that may provide advice on what to do next. Product depending on the user's choices as a whole [6]. It helps people out by giving them resources to utilize in making choices that are meaningful to them and address their issues [7]. Many well-known online retailers, like Amazon, Netflix, Pandora, etc., use RS extensively. The e-commerce RSs will recommend things according on the user's interests in news, publications, individuals, URLs, and so on [8].

B. Trip-Planning Software

The decision-making processes involved in tourism are difficult because of the wide variety of locations, attractions, activities, and services available to tourists. This is why TRS are of interest to scientists in both academia and the private sector. Many

foreign visitors to the country, now at 26.5 million, is up 18.76% from 2012 [2]. The government of Thailand has made boosting the number of visitors (both foreign and local) and the revenue generated by tourism a top priority. The tourist industry in Thailand

than ever before for travelers and service providers to find exactly what they're looking for, narrow down their options, make informed comparisons, and settle on a course of action. Most prior TRSs have concentrated on cost estimates for planning a trip's location, its activities, and its attractions and services (e.g. restaurants, hotels, and transportation) personalized for each individual user. In terms of technology, these TRSs only provide simple procedures for filtering, sorting, and matching things to the user's strict requirements. However, they are deficient in both theoretical and practical features (such as sparsely, scalability, transparency, system accuracy, theories to enhance personalization, etc.). Improving the traveler's ability to make decisions is a major difficulty in the creation of a TRS that offers individualized suggestions of tourism sites. To do this, it creates unique models of the information-seeking behavior of visitors and needs a thorough knowledge of their decision-making

different sorts of platforms have seen the development or deployment of different TRS (e.g. desktop, browser, mobile). User interest may be estimated, preferred Points of Interest (POIs) can be selected, services and routes can be identified, ranked in order of preference, and a whole trip can be planned with the help of a TRS. While some TRSs also aim to aid travel agencies, the vast majority are designed with the individual traveler in mind [9]. They're both based on analogous frameworks but use different technologies, theories for enhancing recommendation methods, data inputs, interaction styles, and personalization Fig. 2 provides an overview of the current TRSs' overall design.

The repository serves as a central storage facility for data collected from a wide variety of sources (sensors, GPS coordinates, surveys, reviews, etc.). The recommendation engine can be broken down into various modules, such as an optimization module, a statistical module, an intelligent module, and so on. The goal is to provide recommendations, rankings, or forecasts for the user based on their needs and preferences, as well as any relevant hard and soft constraints (such as user demographics, the length of their trip, their available funds, the type of trip they're taking, etc.). Typically, a tourist's inputs (implicit, explicit, or both) are collected by the TRS prior to or during their trip so that a user profile can be created

and the recommended result can be calculated and sent back to the tourist. The system's output can be interpreted in a variety of ways by tourists, including as icons representing destinations on a map interface along with a point-to-point route, an agenda, or an itinerary. When displaying the outcome, most TRSs employ spatial web services like the Google Maps API service. Recently, TRSs have been developed that can tailor their output based on factors like as the user's current location and the current weather. A few TRSs provide a capacity for the user to make adjustments to the produced result and for the result to change in response to user ratings [10, 11].

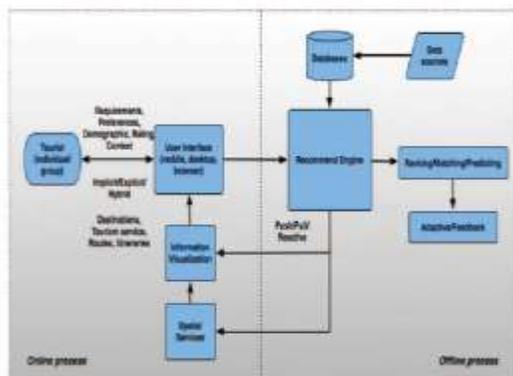
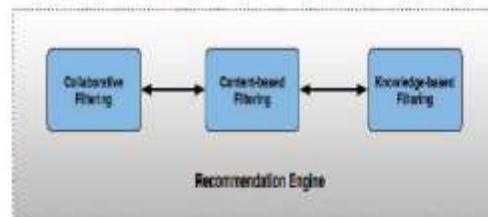


Figure 2. Overall structure of vacation advice apps

Procedures for making suggestions as stated in [12], RS may be categorized in terms of its severity. Of customization, such as how helpful and precise the suggestions are. No personalization, temporary personalization, and long-term personalization are all ways in which the level of customization may be described (long term). A no personalized RS is a basic system that makes suggestions without considering the user's tastes. For example, the RS only compiles a list of the most sought-after products (i.e., editors' picks or best-sellers) based on the total number of reviews and/or sales. Thus, the suggested outcomes would be helpful for other general users of the system. Non-personalized systems have not been a focus of RS research because of their lack of autonomy in making decisions [7]. An ephemeral and personalized RS is superior than a no personalized RS in terms of the inclusion of information connected to the system's users (i.e., user preferences, sociodemographic information, etc.). As a result, each user would be presented with a unique set of suggestions. As an example, Trip-advisor1 would suggest a place to visit based on the user's profile data, including their age, gender, and marital status. Previous studies have examined a wide variety of tailored RSs, and researchers have classified them based on the information-filtering mechanisms they

use [7, [13]-[15]. Following this, we'll take a quick look at the recommendation engine (Fig. 3), which is built from a variety of recommendation methods based on research from [14]. We will talk about the merits and demerits of each, as well as the hybrid filtering strategy used (i.e., the interconnection of many RSs).



a) Collaborative filtering: This strategy is generally utilized by the most deployed recommendation engines (see Figure 3). Systems. Users with similar characteristics are taken into account when making recommendations, and well-liked products are also suggested. However, this method still has a cold-start issue since it requires an initial rating of the new item or user before making a suggestion. The second kind of recommendation strategy is called "content-based filtering," and it makes suggestions to the user based on the user's prior searches and queries. The user has to start from scratch and submit a lot of information before the algorithm can provide a suggestion, which is the biggest negative. Unless a sufficient amount of previous data has been stored, the system will not be able to provide reliable findings [13]. Overspecialization is another prevalent issue [7] due to the system's tendency to propose the item that the user loved the best. Expertise-based filtering (c): Recommends products to the user based on prior subject knowledge. That is to say, the system understands how the item pertains to the person in question. Case-based reasoning and ontological approaches are particularly useful for this purpose. Both [9] and [16] use systems that use the prior experiences of travel firms and groups of experts to provide recommendations.

All the above-mentioned recommendation methods have their advantages and disadvantages, which is why d) hybrid filtering was developed. Reasons for Advising a Hybrid Approach the goal of combining techniques is to maximize performance while eliminating any drawbacks to one approach. In addition, there are a plethora of hybridization approaches, such as the mixing of several recommendation systems (weight, switching, mixed, feature combinations, cascades, feature augmentations, and met levels) [13]. Researchers now have the tools they need to create a TRS that is

intelligent, interactive, and adaptive; that can be automated; that can support a higher level of user satisfaction than ever before thanks to advancements in ICT like Artificial Intelligence (AI), the Semantic Web, communication networks, and so on. To that end, we're working on a system design to meet those goals.

III. METHODOLOGY

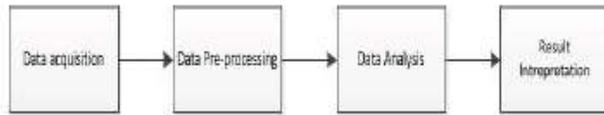


Figure 4. Data Mining Framework

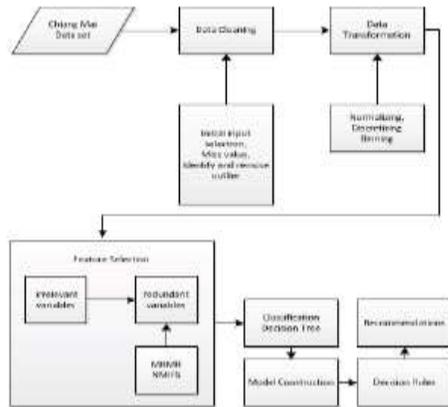


Figure 5. Theoretical framework for the targeted radio spectrum

Figure 4 depicts the proposed DM architecture, which includes four stages: data collecting, data preparation, information dissection and hypothesizing about the findings. (1) A four-part questionnaire was developed for data collection and was sent out and collected in Chiang Mai, Thailand. (2) A variety of data pre-processing procedures, including data cleaning, data transformation, and feature selection approaches, are applied to the gathered data before it is used. In the third and final stage, a decision tree C4.5 is used as a classifier to analyze the data. Down the last stage, we seek to zero in on the most useful aspects and discover the best models. The fourth and last step is to analyze the resulting optimum decision trees and derived rules of thumb. Fig. 5 depicts the general process flow.

The Gathering of Information We choose a questionnaire as a data collection approach due to its efficacy as a tool for gathering information from tourists, allowing us to better comprehend their search behavior when evaluating travel information and their decision-making procedure when selecting

a location. In order to further tailor the questionnaire to potential respondents, researchers conducted preliminary analyses to identify the myriad of variables that affect vacationers' top choices of places.

Each of the four sections of the questionnaire is dedicated to a different collection of variables relevant to travelers' most-preferred vacation spots. Following:

1) Characteristics of the trip itself: these factors are the most influential [17]. The duration of the journey, its intended use, its content, and so on are all relevant factors. Tourists' psychological, philosophical, and economic circumstances are all factors in their final decision [17]. Third, why people travel: literature evaluations have shown that a traveler's motive is a major element in their choice of vacation spot. Indicator of why vacationers prefer a certain location [18]. Socio-demographic details about vacationers: individuals' demographics may affect their information-seeking behavior [19]. It was determined that the five most popular tourist spots in Chiang Mai, Thailand should each distribute and collect 4,000 questionnaires. The most popular places were compiled from user reviews on the travel website Trip Advisor. The poll was sent to both foreign (60%) and local (40%) vacationers. Art in Paradise (27.7%), Mae SA Waterfall (22.06%), and Hay Tung Tao Lake (19.18%), the Museum of World Insects and Natural Wonders (16.97%), and Boa Thong Waterfall (14.09%) were among the popular attractions. On average, it took the respondents 15-30 minutes to finish the survey. The data pre-processing step was entered for 3,695 valid questionnaires including 145 variables; 35 samples were discarded due to insufficient data.

This suggested framework utilizes questionnaire-derived variables to categories the tourist's most desired location based on factors such as travel preferences, budget, and personality. This article describes the demographics of tourists as well as their behavior, spending habits, reasons for travelling, and other relevant facts.

Data Cleaning and Preparation

Incomplete, noisy, and inconsistent data is typical in the real world. In the case of surveys like the ones we conduct, for instance, mistakes in data entry or the deliberate submission of false information by respondents who want to protect their anonymity are both possibilities. Data of high quality is essential for accurate categorization. We integrated the data, cleaned the data, transformed the data, and selected the variables to analyze utilizing feature selection techniques to get the job done. Selecting subsets of

important characteristics that characterize the output classes is known as feature selection or variable selection. This procedure is crucial not just for increased efficiency and usefulness, but also for enhanced precision. In this study, we focus on minimizing the number of variables while retaining as much information as possible. So, the goal is to improve classification model performance while simultaneously decreasing the quantity of required user inputs. In this research, we offer a Mutual Information (MI)-based two-stage filtering approach to priorities features and eliminate duplicates. In the process of feature selection, MI is employed as a metric to describe the variables' importance and redundancy. Assuming the variables are unrelated to one another, the MI value would be 0. In general, the bigger the MI value, the more important the dependent variable. The marginal probability distribution functions of X and Y is p(x) and p(y), whereas the joint probability distribution function of X and Y is p(x, y). In this case, the MI is expressed as:

$$MI(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right)$$

1.) The Basic Filtering Technique

The first stage of filtering is meant to order the variables and get rid of the ones that aren't independent. Separated from the dependent variable. To filter out superfluous information, we used the Max- Relevance feature selection method [20], using MI as the metric of choice. For every set of explanatory and criterion variables, we calculated the MI score. To exclude characteristics that contributed less or were unrelated to the predictive capacity, we sorted them in decreasing order and applied a threshold value (the threshold value is determined manually).

2) An alternate filtering strategy

As a second round of filtering, we employed the mutual information-based feature-selection algorithms Minimum Redundancy Maximum Relevance (MRMR) [20] and Normalized Mutual Information Feature Selection (NMIFS) [21] to get rid of superfluous information. Taking the highest MI G value into account, we determined that this was the best feature space. When G 0, further feature selection will be halted.

a) Algorithm for MRMR

The MRMR method [20] is based on the concept of utilizing the MI value to order features according to minimum redundancy and maximum relevance. Redundancy between features and their relevance to a class are both calculated by MRMR. It may be stated as (1).

$$MRMR = \max_{i \in \Omega_s} [I(i, h) - \frac{1}{|S|} \sum_{j \in S} MI(i, j)] \quad (1)$$

$$MI2(i; j) = \frac{MI(i; j)}{\min(H(i), H(j))} \quad (2)$$

$$NMIFS = \max_{i \in \Omega_s} [I(i, h) - \frac{1}{|S|} \sum_{j \in S} MI2(i, j)] \quad (3)$$

3. Analyzing the Information

The suggested TRS uses a decision tree as its classifier/model because of the many advantages it offers. A decision maker, like ease of use and clarity. The decision-making process is modeled as a flowchart, making it simple to grasp. In terms of technological considerations, it solves the sparsely and scalability problems plaguing the TRS. There are nodes and leaves in the decision tree. Test set instances begin their journey to a leaf node at the root node. Internal nodes entail checking a specific property, which results in a binary or multi-way split. A class label (the result of the classification) or the instance's ultimate verdict from the test data is represented by the leaf nodes. [22]. we must follow the decision tree from its trunk all the way out to its leaves before we can confidently advise tourists on where to go. There are several decision tree algorithms available, such as Hunt's algorithm, Top-down Induction of Decision Tree (TDIDT), ID3, CHAID, CART, and C4.5. The criteria for splitting, the extent of pruning, the types of characteristics, etc., are all different.

IV. EXPERIMENT DESIGN

1. Data set representation

Details about the data set utilized for this analysis are provided in Table 1. There are five travelers' records in the dataset. Locations of choice. In spite of include all five locations in the decision tree model; the classification accuracy was just 36.1%. The decision tree model was also overly complicated, with a high tree size and a significant number of leaves, both of which made the model opaque to the decision-maker. This multi-classes classification problem is broken down into manageable chunks by first learning which types of tourists visit which cities; then using that

data along with insights from Chiang Mai's tourism experts and Trip Advisor to determine which cities are most popular with each of those tourists. This led to the development of the two groups, which are shown in Table 2. Decision tree models were built using these taxonomies as inputs. The data from the Museum offers a binary classification challenge, whereas the data from Nature presents a multi-classification problem. Because both types of museum in the Museum data set serve distinct purposes, we divide them into distinct categories. There are a total of three categories in the Nature dataset, with two of them representing the waterfall and one representing the lake.

TABLE 1. CHARACTERISTICS OF THE DATA SET USED IN THIS STUDY

Data set	# Features	# Classes	# Sample
Tourist destination choice	145	5	1,632

B. Preparing the Data

The data cleansing procedure starts with the first selection. At this stage, we use what we've learned from the tourist industry to filter out characteristics that aren't directly linked to the final product. Next, we subject both sets of data to missing value analysis. The binning technique was used to discredited continuous variables. With a bin size of 10, we can divide the data well. In order to standardize some of the discrete variables, we tapped the expertise of professionals in the tourist industry. The suggested two-stage filtering procedure was implemented once the data set had been cleansed and modified. This was completed in an effort to cleanse the data collection of superfluous or duplicative elements. In the first filtering stage, we employed between 17 and 18 criteria to determine which attributes were most important, depending on the data set. The features in the subset were then run through the MRMR and NMIFS feature selection algorithms to get rid of the ones that weren't necessary.

C. Grouping and Model Building

We used a decision tree to build a classifier after removing superfluous characteristics and honing in on the intended ones. The effectiveness of the two feature selection methods in C4.5 is analyzed. This study used a technique known as K repeat holdout. Sixty percent of each data set was chosen at random for training, whereas 20 percent was stratified (i.e., the percentage of each class in the training, validation, and testing sets is the same). Training and

validation sets' prediction accuracy over iterations was averaged. Finding the best models for each batch of data requires a variety of confidence level settings for decision tree pruning. With a step size of two, the confidence levels may be anywhere from 0.1 to 0.5. In terms of validation sets, the best-case scenario is when their mean accuracy is highest. Second, the validation set's mean accuracy must be lower than the training set's mean accuracy.

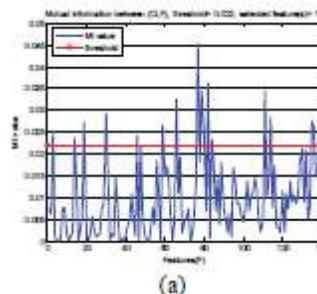
V. RESULTS AND SYSTEM EVALUATION

C4.5's categorization rate results are shown in Table 2. The Museum dataset was successfully classified at an 80% rate using the single best learner. An overall classification success percentage of 49.72% was found in the Nature dataset. When comparing the two feature selection techniques, NMIFS is regarded as more effective than MRMR for both datasets.

TABLE 2. ACCURACY RATE FOR EACH DATA SET

Data set	# of classes	#Sample	Confidence level	Single best learner accuracy rate
Museum	2	729	0.39	80%
Nature	3	903	0.24	49.72%

As may be seen in Fig. 6, the results of the data pre-processing on the Museum dataset are shown. The MI value from the first filter technique is shown in Fig 6(a), where the threshold was 0.022, 128. The data set was cleaned up by removing variables. The MI G values for both feature selection methods are shown in Fig. 6(b). Whenever a negative number was reached, feature selection ceased.



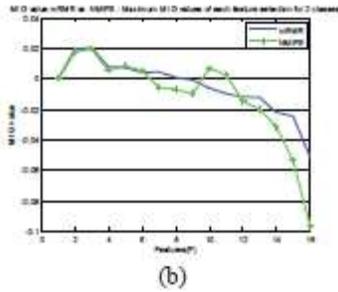


Figure 6: The MI value (a) and the MI G value (b) from the Museum data set's two-step feature selection approach

Feature selections from each are shown in Table 3. Methods for selecting features from the Museum's data collection. Features that are part of the "optimal subset" are denoted by bold variables. After using the second filtering strategy, the MRMR algorithm chose eight ideal features from the Museum data set, whereas the NMIFS chose six. Feature a stands out as the most crucial. Assuming that one of the museums specializes in insects provides an explanation for this observation. Three features (c, d, and b) were used in tandem to aid in the data set's classification. Using a combination of four characteristics taken from the NMFIS, the best decision tree for the Museum dataset is determined, and decision rules are constructed (See Fig 7 and 8). Due to its small size (17 nodes) and plenty of leaves (10 in total); the resulting decision tree is considered very easy to comprehend. An analysis of the Nature dataset revealed that b2 (travel goal) is the most crucial variable to consider.

TABLE 3. FEATURE RANKING BASED ON THE MRMR AND NMIFS ALGORITHMS (MUSEUM DATA SET)

Algorithm	Selected feature
MRMR	a c d b e f g h i j n k l m o p
NMIFS	a c d b g h f j i o k e n l m p

a: deepest impression is wildlife b: to visit place I have never been before c: The people who are companying are friend d: books and guides influences your decision to visit Chiang Mai

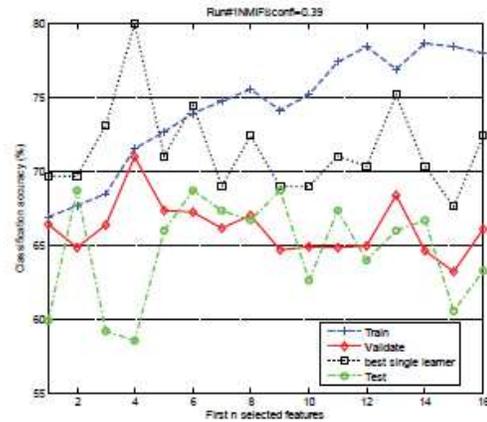


Figure 7. Accuracy rate for the Museum data set

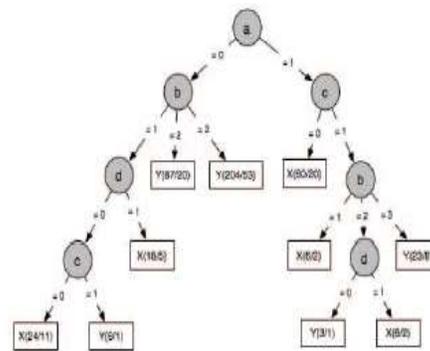


Figure 8: The validation data-driven best-case scenario decision tree for the Museum dataset. (Art in Paradise, Chiang Mai 3D Art Museum and Museum of World Insects and Natural Wonders, X and Y, respectively) The confusion matrix, which includes both the true and anticipated classifications made by the best decision tree, is used with the accuracy rate to assess the model's efficacy. There were a greater number of false positives (samples from the Museum of World Insects that were wrongly labeled as samples from the 3D Arts Museum) in the Museum of World Insects, as seen in Table 4.

TABLE 4. CONFUSION MATRIX OF THE MUSEUM DATA SET

		Predict	
		Museum of world insect	3D arts Museum
Actual	Museum of world insect	26	21
	3D arts Museum	8	90

To aid in decision making, the derived optimum decision tree is used to produce decision rules of the Museum data set, which are then presented as Table 5 shows. For the Museum dataset, eight rules are developed.

TABLE 5. THE DECISION RULES OF THE MUSEUM DATA SET

```

if a == 0, then
  if b == 1 then
    if d == 0
      if c == 0 then , class = X;
      elseif c == 1 then, class = Y;
      end
    elseif d == 1 then, class = X
    end
  elseif b == 2, then class = Y;
  elseif b == 3, then class = Y;
  end
elseif a == 1
  if c == 0 then, class = X;
  elseif c == 1
    if b == 1, then class = X;
    elseif b == 2
      if d == 0 then, class = Y;
      elseif d == 1, then class = X;
      end
    elseif b == 3, then class = Y;
    end
  end
end
end

```

VI. CONCLUSION

In this research, we introduce a decision tree-based tourist recommendation system to address the problem with existing TRS for specific destinations. The employing expertise in the tourist industry, the data set was partitioned into two sub-sets. This was executed to lessen the decision tree's complexity and improve its categorization accuracy rate. As a result of analyzing the data from NMIFS, the most accurate and straightforward (i.e., smallest in terms of number of leave and overall size) decision trees possible have been crafted for final destination selection. Rules for making selections were mined from decision trees. It is clear that NMIFS is the superior strategy since it employs a smaller feature set than MRMR does while dealing with both data sets. In conclusion, testing findings support the practicality of the suggested TRS. The needs of visitors coming to or already in Chiang Mai are met by the projected TRS. To further improve the data sets' categorization accuracy in future study, several classifiers might be evaluated. As an added bonus, a front-end web application with an interactive and adaptable UI will be developed.

REFERENCES

[1] "Economic Impact of Travel & Tourism 2014 Annual Update: Summary." *World travel & tourism council*.
 [2] "Thailand Annual Report 2013."

[3] E. Pentane and L. D. Petro, "From e-tourism to f-tourism: emerging issues from negative tourists' online reviews," *J. Hosp. Tour. Technol.*, vol. 4, no. 3, pp. 211–227, 2013.
 [4] B. Pan and D. R. Fesenmaier, "Semantics of Online Tourism and Travel Information Search on the Internet: A Preliminary Study," *Inf. Commun. Technol. Tour. 2002 Proc. Int. Conf.Innsbr. Austria 2002*, pp. 320–328, Jan. 2002.
 [5] E. Pitoska, "E-Tourism: The Use of Internet and Information and Communication Technologies in Tourism: The Case of Hotel Units in Peripheral Areas," *Tour. South East Eur.*, vol. 2, pp. 335–344, Dec. 2013.
 [6] G. Häubl and V. Thrifts, "Consumer Decision Making in Online Shopping Environments: The Effects of Interactive Decision Aids," *Mark. Sci.*, vol. 19, no. 1, p. 4, Winter 2000.
 [7] F. Ricci, L. Roach, and B. Shapiro, "Introduction to Recommender Systems Handbook," in *Recommender Systems Handbook*, F. Ricci, L. Roach, B. Shapiro, and P. B. Kantor, Eds. Springer US, 2011, pp. 1–35.
 [8] P. Redneck and H. R. Varian, "Recommender Systems," *Common ACM*, vol. 40, no. 3, pp. 56–58, Mar. 1997.
 [9] G. I. Alptekin and G. Buyukozkan, "An integrated case-based reasoning and MCDM system for Web based tourism destination planning," *EXPERT Syst. Appl.*, vol. 38, pp. 2125– 2132, 2011.
 [10] R. Analects, L. Figueiredo, A. Almeida, and P. Novas, "Mobile application to provide personalized sightseeing tours," *J. Newt.Compute. Appl.*, vol. 41, pp. 56–64, May 2014.
 [11] L. Sebastian, I. Garcia, E. Oneida, and C. Guzman, "e- TOURISM: A TOURIST RECOMMENDATION AND PLANNING APPLICATION," *Int. J. Art if. Intel. Tools*, vol. 18, no. 5, pp. 717–738, Oct. 2009.