

Evaluating the Effect of Weather on Tourist Revisit Intention using Natural Language Processing and Classification Techniques

Evripides Christodoulou¹, Andreas Gregoriades¹, Maria Pampaka², Herodotos Herodotou¹

¹Cyprus University of Technology, Limassol, Cyprus

²The University of Manchester, Manchester, UK

{ep.xristodoulou}@edu.cut.ac.cy, {andreas.gregoriades,herodotos.herodotou}@cut.ac.cy,
{maria.pampaka}@manchester.ac.uk

Abstract—Tourists’ revisit has significant monetary benefits to destinations because the cost of retaining existing visitors is less than attracting new visitors. Re-visit intention is often based on tourists experience and satisfaction at a destination. An important aspect that influences the relationship between satisfaction and intention to revisit is the weather conditions at a destination given the increased frequency of heatwaves that strike summer holiday destinations over the summer months. This work applies natural language processing and classification techniques to evaluate the impact of weather information on revisit intention utilizing reviews from TripAdvisor and online weather data. Information retrieval techniques (Doc2Vec) are applied on online reviews collected during the summer months between 2010-2019 from tourists that visited Cyprus. Reviews are labeled as “revisits” or “neutral” based on their textual content. The labelled reviews dataset is enhanced with weather information based on the reviews’ timestamp, such as temperature and humidity of tourists’ country of origin and Cyprus at the time of the visit to the hotel/destination. To account for the influence of hotel infrastructure and available services to deal with heatwaves (i.e., climate-controlled), the training dataset included hotel star rating as an additional parameter. An ensemble gradient boosting tree classifier is trained utilizing the compiled dataset to predict revisit intention. The classifier is evaluated against the area under the curve. To interpret the classifier’s inherent patterns, a popular machine learning interpretation technique is used, namely Shapley Additive Explanation (SHAP). Visualizations of the model using SHAP indicate that the heat index and weather difference between destination and country of origin influence revisit intention. Such preliminary insights are encouraging for further investigations with an end goal to develop a decision support system to assist destination managers during their target marketing campaigns.

Keywords— XGBoost, Doc2vec, Heat Index, Revisit Intention, Data Mining, eWOM.

I. INTRODUCTION

In tourism, there are two classes of tourist-consumers: the first-time consumer and the returning consumer [1]. Destination marketers are keen to understand what drives tourists’ intention to revisit because the cost of retaining revisitors is less than that of attracting new tourists [2], hence

the subject of tourist revisit has significant monetary benefits. Tourists intention to visit a destination is the result of marketing strategies, promoting among others the destination’s weather [3]. Hall [5] highlights that climate is a key factor influencing travel motivation and destination choice. For instance, in countries with colder summers, one degree increase on average temperature over the summer period results in an increase in tourism expenditure at the destination. The literature on the effect of weather and specifically high temperatures on tourists satisfaction highlights that tourists expect specific climate and weather conditions based on what is advertised about the country of visit [3] and when this gap is large the level of satisfaction drops [6]. The recent climate change is affecting negatively summer destinations such as Cyprus with tourists experiencing higher temperatures from what they were expecting [6], with high star rating hotels seem to be dealing better in satisfying tourists expectations than lower rating hotels, possibly due to better climate control services [6].

With the introduction of Web 2.0, users are able to express their opinions globally in terms of micro blogs or electronic word of mouth (eWOM). Harnessing such data enables the identification of patterns that can be utilised to improve service and touristic experience. TripAdvisor is a popular platform for tourist related eWOM with numerous studies utilizing such data to understand and measure customer satisfaction [7]. Most studies on tourists’ revisits use surveys or interviews to examine this effect [1]. However, such methods are expensive and time consuming. This research proposes the analysis of eWOM through natural language processing to assess revisit intention and hence has a comparative monetary advantage to traditional approaches. Weather related information obtained from historical records are used to enhance reviews with weather conditions during the time of visit. Preprocessed dataset is used to predict revisit intention using extreme gradient boosting classification (XGBoost). Insight from this analysis will be embedded in a decision support system that will assist destination manager in their effort to satisfy tourists during periods of heatwaves. The research questions addressed in this work are:

1. Does the heat index affect tourist intention to revisit and if so, which is the heat index threshold that decreases that intention?
2. Does the difference between the expected and actual heat index of a destination affect the intention to revisit?
3. Does the difference between the heat index of the tourists' home country and the destination affect their intention to revisit?
4. How does hotel star rating affect the above relationships?

The paper is organized as follows. The next section addresses the literature on revisit intention, the effect of weather and climate on tourism, and natural language processing techniques for analyzing eWOM. Subsequent sections elaborate on the method followed and the results obtained. The paper concludes with the discussion and future directions.

II. LITERATURE REVIEW

The literature related to tourists' revisit intention and how this is linked to weather conditions is presented next.

A. Revisit Intention and weather comfort

Revisit intention is mainly the intention of a visitor to visit the same tourist destination more than one time [8], with much research investigating the patterns that lead to increased repeated purchase [9]. Tourists choose their destination based on the possibility that this will best meet their needs. This is influenced by variables such as economical state of the destination, its image, environment and weather [3]. During decision making, travelers consider the satisfaction levels from past destination visits. This is influenced by the destinations' ability to satisfy their expectations in terms of weather comfort, attractiveness, quality, and risk [10]. Climate factors, such as temperature, wind, sunshine, and humidity, influence positively summer destinations performance. Climate and weather conditions can strongly influence the tourism industry [11] and are considered as key factors in the success of a touristic destination. Climate is considered an expected property of a summer destination that affect prospective tourists' intention to visit them for the first time. However, the overwhelming consensus on global climate change and its effects in the Mediterranean region started to have a negative effect on tourists' experience at these destinations that seem to affect revisit intention.

High temperature weather is an important threat to human health and a strong inhibitor to comfort and pleasure and thus has a strong effect on the intention to revisit a destination that is stroked by heatwaves regularly. This occurs because higher heat creates the feeling of unfamiliarity and discomfort in undertaking physical activities [6]. Thus, with the increase of heat waves duration, frequency, and intensity and with predictions that the global average surface temperature would rise significantly by the end of the 21st century [12], the future of summer destinations is threatened.

Climate conditions are used in marketing promotions to create a positive destination image [6] with Cyprus promoted as a sunny destination with nice beaches [13]. Tourists' past

experiences with weather conditions at a destination and the activities they are able to do during their holidays act as influencing factors to revisit intention. Weather is a factor for which tourists have easy access to information prior to their visit (from online destination descriptions), and hence their decision can be based on their climate preference. In the case of potential revisit, this decision is affected by previous years' weather conditions and experience[14].

Heat Index (HI) is a popular method of measuring air temperature combined with atmospheric humidity and it can determine a person's degree of comfort and satisfaction when exposed to it [14] is also known as the "real feel". The HI has been used as a measure of heat exposure in several studies throughout the world. During high temperature, the human body sweats but this ability can be reduced when the temperature is accompanied by high humidity. This combination is an indicator of a high HI and creates discomfort, tiredness, and stress in those who extend it. HI can either be obtained from Steadman's HI tables or calculated using one of many algorithms that reproduce the values in these tables. The algorithm more suitable to our case is the method used by [15] that employs air temperature and relative humidity since these were readily available from historical weather databases. The following equation emulates heat and moisture transfer to calculate HI.

$$\begin{aligned}
 HI = & -8.7847 + 1.6114T - 0.012308T^2 \\
 & + H[2.3385 - 0.14612T \\
 & + (2.2117 \times 10^{-3})T^2] \\
 & + H^2[-0.016425 + (7.2546 \times 10^{-4})T \\
 & + (-3.582 \times 10^{-6})T^2]
 \end{aligned}$$

T = Air Temperature in Celsius
 H = Relative Humidity

Equation 1: Steadman's HI

Similarly, the impact of summertime overheating on buildings is a growing issue with accommodation in dense urban areas with limited open spaces that increase heat discomfort. Hotels of different ratings have different indoor discomfort indexes (depending on quality of climate control and building materials) with different capabilities in supporting indoor activities in different sizes of open spaces in airconditioned or climate-controlled areas (bigger lounge, swimming pool areas, breakfast rooms etc.) that affect tourists' satisfaction during heat wave periods.

III. TECHNICAL BACKGROUND

The main natural language processing techniques relevant to this study with emphasis on word embedding that is used to label reviews according to revisit intention are presented next. The section also introduces the classification technique used to predict the intention of tourists to revisit the destination.

A. Natural Language Processing

Natural Language Processing (NLP) refers to the field of research that focuses on the analysis/processing of human language using algorithms. It falls under the umbrella of artificial intelligence and for its operation, it uses machine learning, with the aim of finding patterns and information

through raw data [16]. Regression and classification are the main methodologies used to find patterns. Regression is used to predict values ranges, such as finding the change in the stock market from the verbal content or news articles. Classification is used to find the category to which input data falls, for example finding the genre of a film through the script or summary .

Natural language processing has many functions [16], the most popular and relevant to this study are:

1. Language Modelling: This function focuses on modelling the input language so that it can be used to understand and predict the connections between words.
2. Text Classification: This is the most common function in NLP and focuses on the classification of text in categories using certain criteria, such as revisit intention and non-revisit intention.
3. Information Retrieval: Retrieval is about finding relevant information linked to a keyword or keywords, for example “revisit intention”.

In this work, we utilize information retrieval to label reviews according to their authors’ intention to revisit the hotel or touristic destination.

B. Word Embedding

To process unstructured information in text, we often need to transform this data into numerical format that can be digitally processed. This transformation of raw text into a vector format is known as word representation and has been successfully applied in recommendation systems by generating embeddings for millions of data, helping users to retrieve relevant information [17]. Word embeddings represent words as numerical vectors based on the context in which they appear and is a popular method of analyzing text. This numerical word representation allows the mapping of words in a vocabulary to a point in a vector space. This technique is popular in capturing semantic relations among words in large corpuses and are usually learned from words co-occurrence information. This technique is based on the distributional hypothesis stating that words occurring in similar linguistic contexts have similar semantic meaning [18]. Word2Vec is one of the most popular word embedding algorithm that is based on artificial neural networks, bag of words, and skip gram models.

Word2vec is a prediction-based method that can be implemented in two ways: as a continuous bag-of-words and as a skip-gram. Skip-gram is an unsupervised learning technique used to find the most related words to a given word while the former attempts to predict a focus word given its context. This technique also employs neural networks to learn the mapping of words to a point in a multi-dimensional vector space. The two key parameters used during word2vec model training are the number of the embedding dimensions and the number of words before and after the target word, which is considered as its context. The advantages of word2vec compared to other word embeddings algorithms is its ability to predict the probability of a word’s similarity to a given context and the ability to predict context similarity in a given word in addition to word similarity prediction [19].

Doc2Vec is a document embedding-representation algorithm used in natural language processing extending the word2vec with the option of converting a sequence of words into a vector. This model has the ability not only to convert a series of words into a vector but can understand the position of the word, its surrounding text, the order of the words, the semantic distances between them to create a vector that represents the general meaning of the document . Similar with word2vec, doc2vec can use similarity metrics to assess the similarity between two documents using a pretrained model [19].

C. Extreme Gradient Boosting(XGBoost)

XGBoost is a gradient boosting decision tree model [20] for classification and regression problems and has been extensively used in academia and industry due to its excellent performance in machine learning tasks. XGBoost offers improvement over traditional boosting due to the optimization of the loss function, it prevents overfitting with regularization hence models produced can be generalized; and it is computationally very efficient [20]. XGBoost combines multiple decision trees with each tree built based on the result of the previous tree[20]. It is an ensemble method since it combines multiple classification and regression trees (CARTs), each composed of a number of nodes. Because XGBoost generates high performance on predictions, it has been used in research related to weather predictions, traffic accidents with the combination of Shapley Additive exPlanation (SHAP) for plotting models’ insights [21]. XGBoost has been proven effective in studies aiming to find patterns among eWOM data for opinion mining.

IV. METHODOLOGY

The methodology employed to address our research questions include information retrieval techniques using doc2vec and word2vec as well as binary classification using XGBoost. The workflow employed is depicted in Figure 1. The process starts with the data extraction from trip advisor using a python scrapper for tourists that visited Cyprus between 2010 and 2019. The scraper saved in comma delimited format the travelers username, rating of hotel, user helpful votes and contributions, dates of stay and date of feedback, city of stay, hotel stars, country of origin, and the review text. The total number of reviews collected were 65000 all in the English language, for tourists coming from 27 countries and stayed at 2- to 5-star hotels.

The next step in the process is the pre-processing of the data to remove irrelevant ascii codes and the elimination of reviews with incomplete information (i.e., country of origin). Next step is the elimination of reviews from local visitors and the focus of the reviews produced during the summer months (June-September).

A subsequent step is the calculation of the HI for each review (country of origin and destination for specific dates with different HI for each one), utilizing the date of visit from reviews’ metadata, the tourists’ country of origin and the holiday destination in Cyprus. Based on these spatiotemporal information, the corresponding weather data (temperature and humidity) is obtained from the Cyprus Department of

based on game theory. SHAP assigns each feature of the model with an importance value, also referred to as shapely value, used to estimate the importance and effect of each feature on the model's output [24] and the interaction between variables. The final step of the method utilizes the SHAP techniques to interpret the patterns of the learned XGBoost model.

V. RESULTS

Initial descriptive statistics are depicted in Figure 3, showing the percentage of revisit intentions across the past 10 years over all reviews. The level of revisit intention is marginally decreasing throughout the years, and this might be attributed to climate changes.

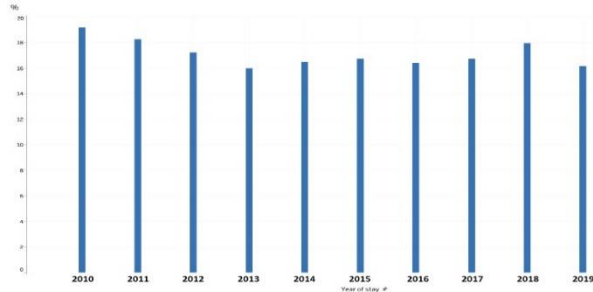


Fig. 3. Percentage of reviews classified with revisit intention across the years based on the similarity threshold of 90% of the Word2vec analysis

After data preprocessing and classifier training the resulting XGBoost model yielded a prediction accuracy of 70% in the test data, 77% F1 score as an indication of precision and recall measurement, and 71% effective class separation measured using the Area Under the Curve (AUC). Using the SHAP techniques, the model's inherent patterns have been visualized using firstly the summary plot of Figure 4 that combines each model's feature importance with their effects on the target variable in term of log-odds. Each point on the summary plot is a Shapley value for an instance of a feature. The features are ordered according to their importance on the y-axis. The color of the points represents the value of the feature from low to high. The points on the graph form a distribution of the Shapley values per feature.

Based on the summary plot of Figure 4, it can be served that the hotel star rating is the strongest feature in the model with higher star rating hotels related with increasing revisit intention (log-odds) and lower stars with neutral hotel revisit. The next feature in terms of importance is the expected weather and actual weather difference. This feature shows that when the expected and the actual weather conditions are similar the revisit intention is increasing. The third feature refers to the difference in climate conditions between the home country and the destination. This shows that when the difference is low to medium, then the intention to revisit is increasing. Finally, the heat index is an influencing factor to revisit intention when is of moderate value, while when the heat index is high the intention to revisit is negative.

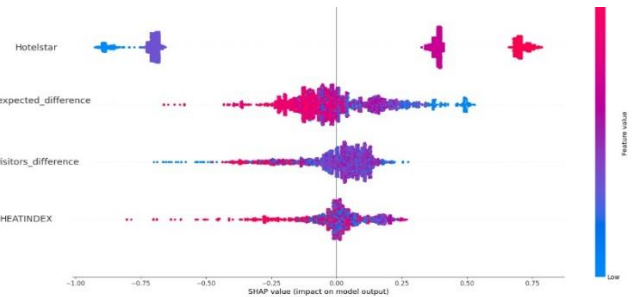


Fig. 4: XGBoost Classifier SHAP Summary with features importance in hierarchical order and intention to revisit in terms of log-odds on the x-axis

Drill down analysis of each of the weather variables and how it interacts with hotel star rating (expressed as ordinal scale 2-5 stars) using the SHAP dependence plot reveal that the negative difference from expected and actual weather conditions tend to increase revisit intention for higher star rating hotel guests, while when the actual temperature is higher than the expected one, the revisit intention drops drastically for higher star rating hotel guests in contrast to lower star rating hotel guests (Figure 5). This could be attributed to the fact that high star rating hotel guests tend to be tourists of higher income which tend to relate with higher age. From the graph we also notice that the high star rating hotel guests demonstrate the highest variability in intention to revisit given the differences from expected weather conditions, while hotel guests of low star-rating have higher return intention when the expected and actual weather conditions are similar.

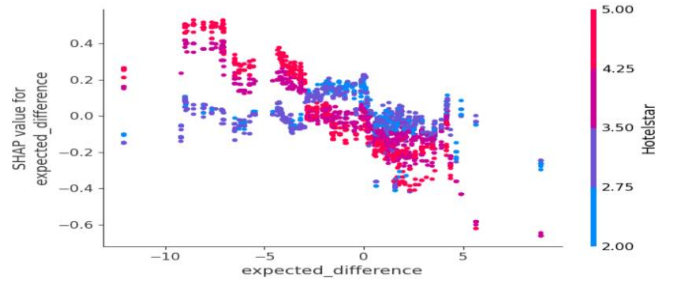


Fig. 5: Difference from Expected and actual condition in degree Celsius and how it interacts with Hotel Star rating ranging from 2-5 stars

The impact of the difference between the tourists' country of origin and the holiday destination at the date of visit is shown in Figure 6. Most tourists prefer the heat index not to exceed 15 to 20 celciusm from their home country. In the case when the heat index difference is gretaer than 20, the revisit intention drops for all hotel star rating visitors.

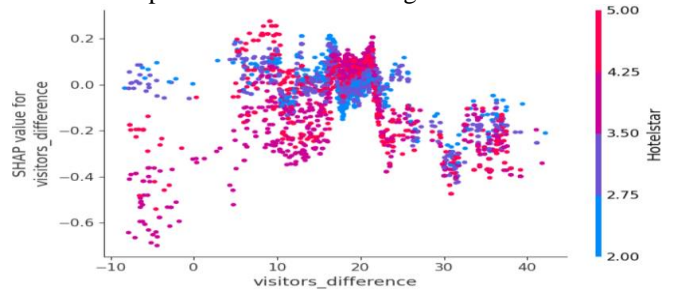


Fig. 6: Difference in heat index between home country and destination and the interaction with hotel star rating

The heat index feature and its interaction with hotel star rating is depicted in Figure 7 with the heat index of 38 and greater being the value with negative revisit intention for high star rating hotels, while the low rating hotels guests lose their interest in revisiting the destination at the heat index of 36 Celsius. It is interesting also to note that the low star rating guests are keener to revisit the destination when the heat index is below 36 and greater than 32 compared to the high hotel star rating guest possibly due to difference in age range and hotel star rating.

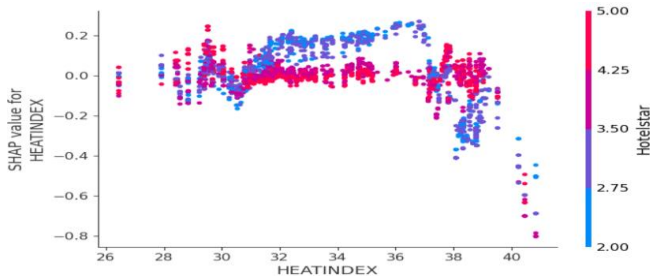


Fig. 7: XGBoost Classifier SHAP Heat Index/Hotel Star and revisit intention

VI. CONCLUSIONS

This research addressed the problem of identifying how weather and hotel star ratings influence tourist revisit intention using natural language processing and classification techniques. Insight from this analysis can be utilised during the design of decision support systems to assist destination managers decision making and recommending them with ways to maintain tourists' satisfaction in extreme weather conditions. In addition, such systems could assist destination management target tourists from countries that are more likely to be satisfied with the country's weather conditions and hence minimise spending for retaining them as customers in subsequent years.

The method utilizes online reviews from trip advisor using information retrieval techniques to label them according to the intention of the author to revisit the hotel or destination. The labelled reviews are augmented with weather related data such as heat index of the touristic destination, difference between the destination and home country, and difference between expected and actual weather conditions. An XGBoost classification model is developed and validated using the augmented dataset. The model's inherent patterns are visualized using the SHAP interpretation technique.

The results show that the ideal range of heat index for travelers to return is 27-37. Regarding the expected heat index value that tourists expect before their arrival and the actual one during their visit, it seems that the comeback intention differs according to the hotel of choice. Guests of 2-3 stars hotels prefer lower heat index while people who choose accommodation with 4-5 stars are tolerant of higher values up to 2 degrees. This might be due to the fact that people who have more luxurious accommodation during their vacation can be satisfied with the amenities, while people with lower hotel-stars rely on outdoor activities to satisfy their needs and expectations [25]. Another important insight is the difference between the heat index of the country of

destination and the country of origin. Tourists show higher intention to return when the heat index difference is between 10- 20 degrees.

Limitation of this work lie in availability of local weather data from tourists' countries of origin. The second limitation of this work regards the labelling of reviews based on their authors' revisit intention. The technique used in this study might miss out reviews that refer to revisit and this need to be examined further. To address this limitation, additional labelling techniques will be evaluated in the future to compare their results against this approach.

VII. REFERENCES

- [1] S. Huang and C. H. C. Hsu, "Effects of Travel Motivation, Past Experience, Perceived Constraint, and Attitude on Revisit Intention," *J. Travel Res.*, vol. 48, no. 1, pp. 29–44, 2009
- [3] N. Stylos, V. Bellou, A. Andronikidis, and C. A. Vassiliadis, "Linking the dots among destination images, place attachment, and revisit intentions: A study among British and Russian tourists," *Tour. Manag.*, vol. 60, pp. 15–29, 2017
- [4] T. Mazzarol and G. N. Soutar, "'Push-pull' factors influencing international student destination choice," *Int. J. Educ. Manag.*, vol. 16, no. 2, pp. 82–90, 2002
- [5] S. Nicholls, "Tourism, Recreation, and Climate Change," *Ann. Tour. Res.*, vol. 33, no. 1, pp. 275–276, 2006
- [6] J. H. G. Jeuring, "Weather perceptions, holiday satisfaction and perceived attractiveness of domestic vacationing in The Netherlands," *Tour. Manag.*, vol. 61, pp. 70–81, 2017
- [7] M. D. Sotiriadis and C. van Zyl, "Electronic word-of-mouth and online reviews in tourism services: the use of twitter by tourists," *Electron. Commer. Res.*, vol. 13, no. 1, pp. 103–124, 2013
- [9] S. (Shawn) Jang and R. Feng, "Temporal destination revisit intention: The effects of novelty seeking and satisfaction," *Tour. Manag.*, vol. 28, no. 2, pp. 580–590, 2007
- [11] D. Scott and G. Mcboyle, "Using a 'tourism climate index' to examine the implications of climate change for climate as a tourism resource," *FIWCTR* pp. 69–88, 2001 .
- [14] J. Xu, Q. Wei, X. Huang, X. Zhu, and G. Li, "Evaluation of human thermal comfort near urban waterbody during summer," *Build. Environ.*, vol. 45, no. 4, pp. 1072–1080, 2010
- [15] E. M. Fischer and C. Schär, "Consistent geographical patterns of changes in high-impact European heatwaves," *Nat. Geosci.*, vol. 3, no. 6, pp. 398–403, 2010
- [16] A. E. Yilmaz, "Natural Language Processing," *Int. J. Syst. Serv. Eng.*, vol. 4, no. 1, pp. 68–83, 2014
- [17] H. Caselles-Dupré, F. Lesaint, and J. Royo-Letelier, "Word2Vec applied to Recommendation: Hyperparameters Matter," *RecSys 2018 - 12th ACM Conf. Recomm. Syst.*, 2018
- [18] S. K. Sien, "Adapting word2vec to Named Entity Recognition," *NODALIDA 2015*
- [19] K. W. CHURCH, "Word2Vec," *Nat. Lang. Eng.*, vol. 23, no. 1, pp. 155–162, 2017
- [20] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proceedings of the ACM SIGKDD-2016*, pp. 785–794.
- [21] A. B. Parsa, A. Movahedi, H. Taghipour, S. Derrible, and A. (Kouros) Mohammadian, "Toward safer highways, application of XGBoost and SHAP for real-time accident detection and feature analysis," *Accid. Anal. Prev.*, vol. 136, p. 105405, 2020
- [22] A. G. Barnett, S. Tong, and A. C. A. Clements, "What measure of temperature is the best predictor of mortality?," *Environ. Res.*, vol. 110, no. 6, pp. 604–611, 2010
- [23] H. Lakkaraju, E. Kamar, R. Caruana, and J. Leskovec, "Interpretable & Explorable Approximations of Black Box Models," *arXiv*, 2017 <http://arxiv.org/abs/1707.01154>
- [24] D. Lubo-Robles, D. Devegowda, V. Jayaram, H. Bedle, K. J. Marfurt, and M. J. Pranter, "Machine learning model interpretability using SHAP values: Application to a seismic facies classification task," in *SEG 2020*, Sep. 2020, pp. 1460–1464.
- [25] W. G. Kim, J. J. Li, J. S. Han, and Y. Kim, "The influence of recent hotel amenities and green practices on guests' price premium and revisit intention," *Tour. Econ.*, vol. 23, no. 3, pp. 577–593, 2017