

Exploring China's 5A global geoparks through online tourism reviews: A mining model based on machine learning approach

Yuyan Luo^{a,b}, Jinjie He^a, Yu Mou^a, Jun Wang^{c,*}, Tao Liu^a

^a College of Management Science, Chengdu University of Technology, Chengdu 610059, China

^b Post-doctorate R&D Base of Management Science and Engineering, Chengdu University of Technology, Chengdu 610059, China

^c Business School, Sichuan Normal University, No.1819, Section 2, Chenglonglu, Longquan Yi, Chengdu 610101, China



ARTICLE INFO

Keywords:

Geotourism
Online review
Text mining
Topic model
Sentiment analysis

ABSTRACT

This study aims to advance the existing research on geoparks by incorporating machine learning models in the analysis of online reviews to provide valuable suggestions for managers in increasing their understanding of the psychological cognition of tourists and evaluating the status of geoparks. A total of 120,532 online reviews of 24 AAAA UNESCO Global Geoparks in China were linguistically analyzed through multiple machine learning models of the support vector machine, the improved latent Dirichlet allocation, and the importance–performance analysis (IPA). Results uncovered 10 attributes and 80 elements of attributes. From the identified attributes, four negative attributes were specified, namely travel cost, tour services, well-known degree, and transportation and accommodation. As the performance of these geoparks in terms of experience products and geological knowledge education was insufficient, several actions for geoparks and geotourism were suggested on the basis of IPA.

1. Introduction

In July 2020, the UNESCO's Executive Board approved the designation of 15 new sites demonstrating the diversity of the planet's geology as new UNESCO Global Geoparks (UGGp). With this year's additions, the number of UGGp is brought to 161 in 44 countries (UNESCO, 2020). Among them, China has 41 UGGp, and 5 of which were new ones added in 2019. (Fig. 1 shows the Zhangye UGGp, which was newly added in 2020). China has the largest number and fastest growth rate in the UGGp network. The construction of the geopark has become an important means to alleviate poverty in China, and the effect of promoting local residents' income and employment is viewed as positive (Xinhua News, 2019). Geotourism not only reduces immigration and unemployment (Zouros, 2010) but also is becoming an important source of income in several rural areas (Brocx, Brown, & Semeniuk, 2019). However, under such an important background of geotourism, a series of problems emerges. For example, in, 2013, UNESCO gave yellow cards to three famous UGGp in China, and urged them to make rectifications to popularize earth science knowledge to the public (Xinhua News, 2013). On April 1st, 2020, the Huangshan UGGp was seriously congested with tourists, and the scenic area management department immediately stopped ticket sales to relieve the congestion

(Xinhua News, 2020). These problems that the UGGp and geotourism are facing call for a better evaluation of the current development status of geoparks that will also allow for easier and faster problem identification.

Geotourism, which is a relatively new concept, has rapidly developed over the recent decades. It is defined as a form of natural regional tourism that focuses exclusively on landscapes and geology. In the traditional quantitative analysis of geoparks, questionnaire analysis is commonly used, which involves a small sample size. However, small samples are prone to selection bias and estimation bias, leading to incorrect analysis results. Large-scale sample research is needed in the future. With the rise of online emergence of travel and online reviews platforms, such as TripAdvisor (<https://www.tripadvisor.cn/>), Qunar (<https://www.qunar.com/>), and Ctrip (<https://www.ctrip.com/>), big data is promptly applied to entering the field of tourism research (Fuchs, Höpken, & Lexhagen, 2014). The online reviews of tourists are widely accepted because of their brevity, timeliness, and sample size. Online reviews influence consumer perception and attitudes toward products (Vermeulen & Seegers, 2009), purchasing decisions (Sparks & Browning, 2011), and product sales (Ye, Law, & Gu, 2008). When consumers purchase travel products online, they often browse through online reviews that other consumers have submitted, which serve as an important

* Corresponding author.

E-mail addresses: luoyuyan13@mail.cduto.edu.cn (Y. Luo), hejinjie@stu.cduto.edu.cn (J. He), my@stu.cduto.edu.cn (Y. Mou), wangjun@sicnu.edu.cn, wangjun.sicnu@qq.com (J. Wang), liutao@stu.cduto.edu.cn (T. Liu).

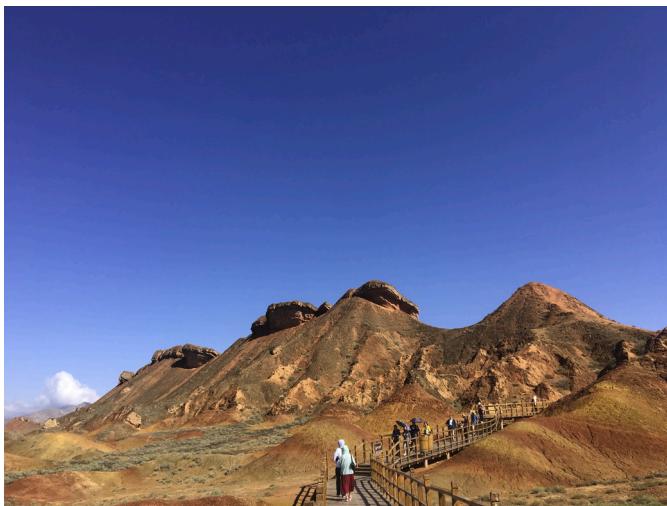


Fig. 1. Zhangye UGGp.

reference for their own purchase decisions. Travelers can access online platforms to provide feedback and advice for other travelers (Neidhardt, Rümmel, & Werthner, 2017). Consequently, research on online reviews has surged and has begun to emphasize the use of new analytical techniques (Cheng, Fu, Sun, Bilgihan, & Okumus, 2019). However, enhancing the practical application of this method in UGGp has become a problem. From the perspective of tourists, utilizing online reviews to continuously analyze the development trend of the UGGp, assess the quality of service management of the tourism destination, and determine the needs of tourists is worthy of research for the development of UGGp in China.

Tourist information websites, such as Qunar and Ctrip, provide a classification of “good” and “bad” reviews for travel products or destinations. These simple classifications often possess tremendous errors, though. Some tourists choose “good” for unknown reasons, but they are express dissatisfaction in their reviews. An enhanced scientific and effective method is needed to analyze the online reviews directly and reveal the actual emotional tendency of tourists. In recent years, scholars have studied sentiment analysis to understand the views and characteristics of the mass or market groups and determine the credibility of the content and motivation for writing comments (Ehsan, Leman, & Begum, 2013). Different sentiment analysis methods have been launched in various fields (Beretić, Đukanović, & Cecchini, 2019). Among the sentiment analysis methods based on supervised machine learning, the super vector machine (SVM) model is the most widely used (Moraes, Valiati, & Gaviao Neto, 2013). The SVM is a classifier that uses annotated data for training to obtain the best separated hyperplane/line to classify new sample data accurately into different categories. This study selects the SVM to classify sentiment. Owing to the multi-dimensional connotation requirement of geotourism, mining the multi-dimensional features of the online reviews is also important. Online reviews about geoparks are extensive and often accompanied by various evaluation noises. Therefore, filtering feature words from the online reviews to collect useful information is difficult for managers. In addition, intuitively discovering problems in the management process from the connotation and requirements of the geotourism is impossible for managers. Topic recognition means to integrate the discrete and isolated online reviews on the same topic and extract the key words of the topic (Wei, 2006). Among the numerous topic models, the latent Dirichlet allocation (LDA) model is the most representative topic model. The LDA is a probability generation model based on Bayesian model, which can extract the implicit topic of texts. By introducing Dirichlet prior distribution to the multinomial distribution, it becomes a complete probability model. The probability model not only maintains the comprehensiveness of the original data, but also reduces the dimensions

of the characteristic vocabulary matrix (Blei, Ng, & Jordan, 2003).

This study aims to propose a machine learning method based on the multi-dimensional connotation and requirements of geoparks and geotourism to evaluate the development status of UGGp through online reviews from the perspective of tourists and discover the advantages and disadvantages during the development. This study first analyzes the emotional sentiment through the SVM, which understands the trend of tourist satisfaction. Subsequently the study uses the LDA to explore the highlights of reviews, understand the perception dimension of visitors, and strengthen the identification of negative subjects. Finally, the attributes and elements of attributes are analyzed by two dimensions, namely, importance and performance using the importance–performance analysis (IPA), which is combined with the SVM and LDA models. Suggestions on geotourism development are then put forward to help UGGp managers further understand the views of tourists. On the one hand, this study offers decision support for the management and planning of UGGp. On the other hand, it provides references for other geoparks to be successfully declared.

2. Literature review

This study uses Citospace, a bibliometric tool, to select information from the Web of Science core database and locate keywords with the theme, “geopark.” Toward the end of October 2019, 662 data samples were obtained. Keywords are important signs of reaction research hot-spots. This study selects the keyword node, displays the high-frequency words in the keyword, and manually adjusts it to obtain Fig. 2. A large node reflects a high number of occurrences of the keyword. Fig. 2 shows that the research on geoparks has been carried out in two main aspects. One is the research on the tourism direction of geoparks represented by the word “geotourism” (Guo & Chung, 2019; Newsome, Dowling, & Leung, 2012; Strba, Krsak, & Sidor, 2018). The other aspect focuses on “geoconservation” (Benado, Herve, Schilling, & Brilha, 2019; Brocx, Semeniuk, & Meney, 2019; Sallam et al., 2020) and “conservation” (Cetin, Zeren, Sevik, Cakir, & Akpinar, 2018; Shekhar, Kumar, Chauhan, & Thakkar, 2019). The words are representative of the study on the protection of geoparks. The main research methods used include qualitative analysis (Esfehani & Albrecht, 2019; Nikolova & Sinnovskiy, 2019; Perotti, Carraro, Giardino, De Luca, & Lasagna, 2019; Ruban, 2015; Shahhoseini, Modabberi, & Shahabi, 2017) and quantitative analysis (Liang et al., 2019; Nobre da Silva, Leite do Nascimento, & Mansur, 2019; Perotti, Carraro, Giardino, De Luca, & Lasagna, 2019). And the further analysis will be made with the key nodes through the relevant literature.

2.1. Concept and development of geopark

A geopark holds an extensive amount of geological paleontology with special scientific importance, rarity, and beauty. It has not only geological and paleontological significance but also archaeological, ecological, historical, and cultural value (UNESCO, 2014). The term geotourism was first defined by Hose (1995, 2006), a pioneering scientist in this field. He described geotourism as the provision of interpretive and service facilities that enable tourists to acquire knowledge and understanding of the geology and geomorphology of a site beyond the level of mere aesthetic appreciation. Dowling (2011) defined geotourism as a form of sustainable tourism, which focused on enhancing the environmental appreciation and cultural understanding of geological features to promote and protect geological heritage. The National Geographic Society defines geotourism as “tourism that sustains or enhances the geographical character of a place—its environment, culture, aesthetics, heritage, and the well-being of its residents” (Stokes, Cook, & Drew, 2003). Geotourism contains all elements of a geoparks geographical characteristics working together to create an experience and is richer than the sum of its parts, appealing to visitors with diverse interests. Based on geoparks, geotourism builds on a destination’s “sense

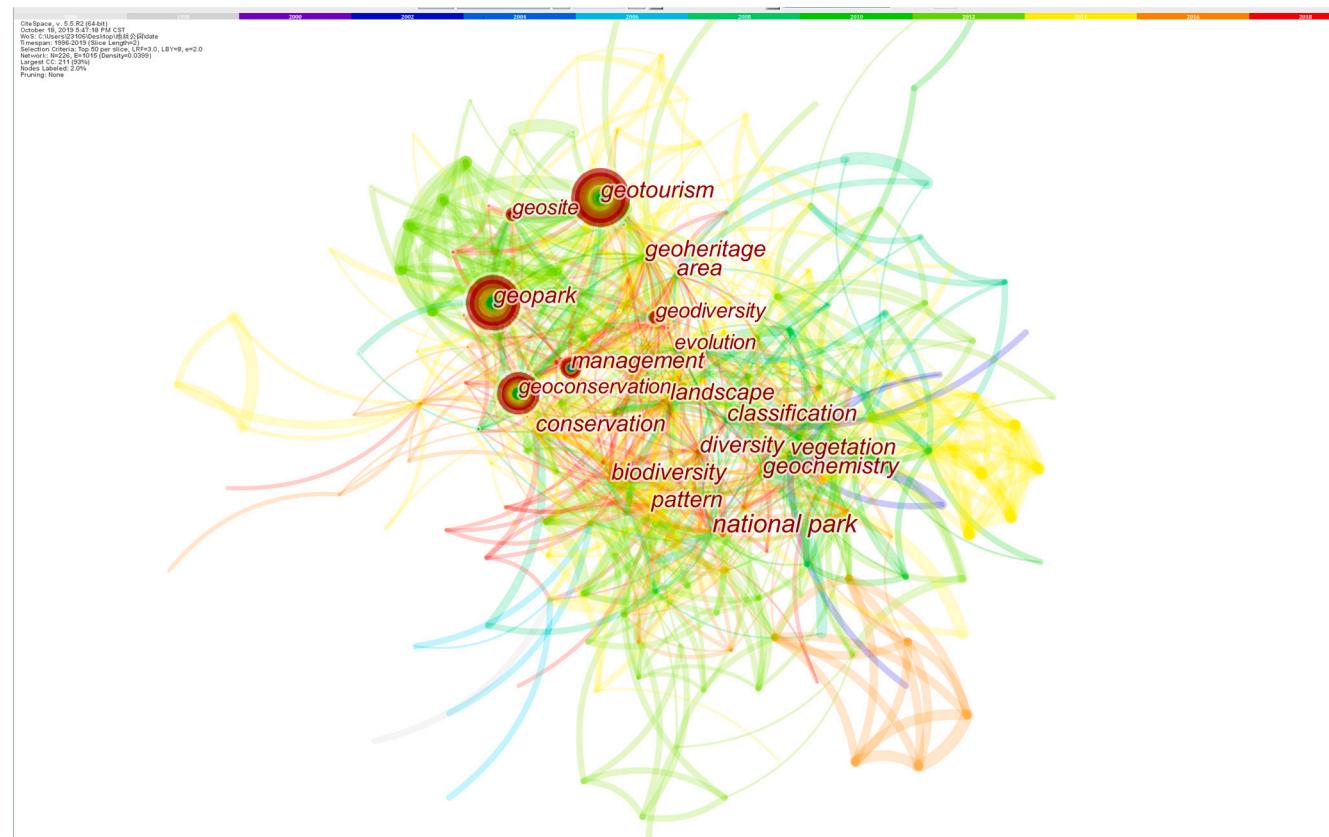


Fig. 2. Key words co-occurrence map of geopark research.

of place" to emphasize the distinctiveness of its locale and benefit visitors and residents alike.

Geoparks directly contribute to the spread of geological knowledge, scientific research, and entertainment activities (Alexandrowicz, 2006). Apart from geological heritage, a geopark should have areas with relevant biodiversity and archaeological heritage and areas where historical and cultural sites have a connection with local geodiversity (Carvalhido, Brilha, & Pereira, 2016). These definitions promote geotourism to encourage geoconservation and sustainable management and for people to appreciate geographic diversity (Han, Wu, Tian, & Li, 2017). Geoparks are also favored by domestic and foreign tourists. Pazari and Dollma (2019) used a questionnaire survey to determine the geographic location of tourist resorts from four aspects, namely, accessibility, state of preservation, scientific value, and educational value. It shows that the geotourism industry should not only promote community economic development but also consider geological protection and the general education of geosciences, it has become increasingly concerned about the harmonious coexistence between man and nature, that is, the sustainable development of the geopark system, including multiple stakeholders.

Most of the existing literature on the development of geoparks conducts conceptual discussions on a qualitative level. Fanwei (2014) evaluated the perspectives of residents on the Mt. Huaying Grand Canyon, Sichuan through questionnaires. Wang, Tian, and Wang (2015) analyzed the geological protection and geotourism of HongKong UGGP from the characteristics of coastal geological heritage and formation of landforms. Wang, Wu, Li, and Chen (2019) selected Dunhuang Geopark as research object. They explained its geological protection and geotourism through qualitative research and shared successful experiences. Purdie, Hutton, Stewart, & Espiner (2020) combined geophysical surveys with visitor surveys ($n = 400$) and semistructured interviews with major information providers ($n = 12$) to explore the impact of environmental changes on tourist experience.

On the basis of the analysis results of the literature measurement tools of Citespace, we find that the existing research from the macro level qualitatively explores the development of geoparks through the multi-attribute issues of the economic development, geological protection, and the geosciences popularization. Scant research analyzes and evaluates geoparks from a micro perspective and multi-dimensionality in the practical management application. Most traditional quantitative analysis of geoparks use small sample questionnaires. However, small samples and questionnaire compilation methods are likely to cause attribute omission or deviation of results. In the era of big data, the generation of online reviews provides data support for large-scale research. Determining the perspective of tourists is an important method to promote geotourism development. Therefore, this study aims to evaluate the development status of the geopark from a micro perspective through online reviews, and discuss how to achieve the sustainable development of the geopark system, including multiple stakeholders.

2.2. Method and application of online tourism reviews

In the digital age, the vigorous development of the Internet and social media has increasingly attracted attention. It not only provides tourists with a wide platform to contribute user generated content data but also generates a large amount of online comment data. Online review data, blog data, and other relevant text format data constitute a special type of big data in tourism research—online text data. Online reviews show destination image in tourist perception (Wong & Qi, 2017), and the affective and cognitive attributes are interrelated in online tourism communities (Garay, 2019). The existence of potential value information in online text data provides support for tourism research. Managing online reviews can promote business (Nguyen & Coudounaris, 2015). Existing research uses such reviews to measure the satisfaction of tourists (Tsai, Chen, Hu, & Chen, 2020; Zhao, Xu, &

Wang, 2019), explore the related attributes of the perception of tourists (Ahani et al., 2019; Mellinas, Nicolau, & Park, 2019), evaluate and improve the electronic word of mouth of a hotel (Garay, 2019; Khorshand, Rafiee, & Kayvanfar, 2020), measure tourism satisfaction (Jia, 2020; Kim, Park, & Lee, 2020), and improve attraction management.

To extract and utilize useful information hidden in online text data, research on the tourism industry has widely adopted a variety of text mining techniques, which include two typical stages: data collection and data mining. The first step in analyzing online reviews is to collect online text data (including travel-related reviews and blogs) from relevant social media sites through the Internet (Xiang, Du, Ma, & Fan, 2017). One widely used approach is a web crawler. Web crawler is defined as a program or software which traverses the web and downloads web documents in a methodical, automated manner (Abukausar, Dhaka, & Singh, 2013).

Data mining is the other key stage aimed at exploring valuable information in text information after data collection is completed. Typical technologies in current tourism research include subject extraction and sentiment analysis. The topic model is a modeling method for the hidden topics in the text, that can be used to mine the latent topics and semantic structure of the text, and it is widely cited in many fields. The existing theme models are the LDA (Chen, Zhang, Liu, Ye, & Lin, 2019), LSA (Kim, Park, & Lee, 2020), and NMF (Chen et al., 2019). The LDA has been extensively used in prior studies to discover key dimensions from online reviews (Guo, Barnes, & Jia, 2017), identify influential subjects from text messages (Pournarakis, Sotiropoulos, & Giaglis, 2017), and derive a set of variables to predict the success of crowdfunding projects (Yuan, Lau, & Xu, 2016).

Sentiment analysis, which originated in the late 1990s, is an analysis method that can be used to help decision makers obtain emotional information by mining and analyzing the emotional content expressed in the text (Pang, Lee, & Vaithyanathan, 2002). This type of analysis is mainly divided into two categories: dictionary-based sentiment analysis (Hu & Liu, 2004) and machine learning-based sentiment analysis (Mowlaei, Abadeh, & Keshavarz, 2020). Dictionary-based sentiment analysis is limited by the richness of context and semantic expression, and the accuracy rate may be low. By contrast, sentiment analysis methods based on machine learning are becoming mainstream. In the existing studies on binary sentiment classification, different machine learning algorithms are used to develop methods for binary sentiment classification (Medhat, Hassan, & Korashy, 2014). The SVM is regarded as one of the most effective machine learning algorithms for sentiment classification (Balazs & Velásquez, 2016). The SVM is a supervised machine learning algorithm based on the principle of structural risk minimization and Vapnik–Chervonenkis (VC) dimension theory (Cortes & Vapnik, 1995). Compared with other machine learning algorithms, machine learning algorithms based on SVM can reduce the VC dimension and improve the classification efficiency (Liu, Bi, & Fan, 2017; Zhang, Xu, Su, & Xu, 2015).

The above research shows that the studies on geoparks need to pay attention to many aspects to promote its development—not only from the macro perspective but also from the micro perspective. Therefore, this study selects the important participants in geotourism, based on the perception of tourists to explore the sustainable development strategy of geoparks. Most of the existing geopark research starts from a small sample, and the research objects are few, which may lead to deviations in the conclusions. In the era of big data, online reviews provide a stable data foundation for this research. However, studies on online reviews from a certain type of theme scenic area are few, and the attributes of existing research are relatively simple. Existing research methods, such as the LDA and the SVM, provide good technical support for research in feature extraction and sentiment classification. They cannot obtain the sentiment tendency of specific attributes though. In this study, the results of topic extraction and emotional orientation are merged to obtain the emotional attitudes of the corresponding attributes, thus providing a fine-grained reference for the research. This

research not only tackles a computer-related problem but also a management problem. To enhance the relevance of the research results to management decision-making, this study integrates the classic management model of the IPA, an ideal model of evaluating services (Sever, 2015), and further proposes a new model based on the decision-making paradigm of geopark management for mining and forming promotion strategies.

3. Methodology

As a resource utilization method, the establishment of UGGp has shown comprehensive benefits in terms of the protection of geological relics and ecological environment and the development of local economy. From the multi-dimensional connotation and requirements of UGGp and geotourism, evaluating the status of UGGp in real time through the online reviews and discovering advantages and disadvantages in its development are vital to the development of UGGp in China. National AAAA (5A) scenic spots are scenic spots classified by the quality of tourist attractions in China. There are five levels, and 5A is the highest level of China's tourist attractions, representing the country's world-class boutique tourist attractions (China National Tourism Administration, 2003). As of April 2019, China has 39 UGGp, and the distribution of which is shown in Fig. 3. Given that 5A scenic spots have numerous tourists and are quite representative, this study has selected 24 5A scenic spots among the UGGp in China as research object, which are presented in red in Fig. 3. This study intends to construct an online review analysis framework for UGGp based on the machine learning method. Through the SVM, we analyze the emotional sentiment and understand the satisfaction of visitors. Simultaneously, this study uses the LDA model to explore the attributes of the reviews to understand tourists' perception. Fundamentally, the attributes and elements of attributes are analyzed from two dimensions, namely, performance and importance by using the IPA model combined with the SVM and LDA models. Under the high-quality and connotative requirements of the geotourism industry, the large number of tourists in 5A UGGp also represents higher requirements and deserves attention. At the same time, 5A UGGp can be used as benchmarks, which are worth learning from. Some of the problems in 5A UGGp are unique to 5A scenic spots, but other parts are also common to other scenic spots. The research paradigm can be used as a reference for other types. The overall framework of the model is presented in Fig. 4.

3.1. Data collection and pre-processing

The data collection and related preparations used in this study are as follows:

The initial step is to collect the data. The online reviews on the various dimensions and aspects of scenic locales are highly suitable for the evaluation and analysis of tourist attractions. The sample data used in this study come from several major China tourism communities, including Baidu Travel (www.lvyou.baidu.com), Ctrip Travel (<https://www.ctrip.com/>), Tongcheng Travel (<https://www.ly.com/>), and Qunar (<https://www.qunar.com/>). As of April 2019, 121,443 online reviews of 24 5A UGGp in China from 2000 to 2019 were collected through the Python crawler program. The reviews obtained are presented in Table 1.

The second step is to complete the preprocessing of the text data. After removing empty words, stop words (default stop words and self-built vocabulary), punctuation, expressions, and other unnecessary texts, 120,532 online reviews remain.

3.2. SVM model

The SVM model is a classifier that uses annotated data for training to obtain the best separated hyperplane/line to classify new sample data accurately into different categories. This algorithm has an excellent



Fig. 3. Distribution of UGGPs in China.

effect on solving two-classification problems. Hence, the SVM model is highly suitable for emotional classification and it is a key machine learning method commonly used in sentiment analysis (Markopoulos, Mikros, Iliadi, & Liotatos, 2015).

The SVM is mainly solved by mapping the vector to the high-dimensional space, establishing the maximum interval hyperplane in the high-dimensional space, and converting the plane segmentation problem into a convex quadratic programming problem through the interval linear classifier. With this method, the two classification problem can be converted into a minimum problem with constraints:

$$\min_{w,b} \frac{1}{2} \|w\|^2 \quad (1)$$

Subject to:

$$y_i [w^* x_i + b] - 1 \geq 0, \quad (i = 1, 2, \dots, n) \quad (2)$$

where w is a normal vector used to determine the direction of the hyperplane, and b is the displacement term, showing the distance between the hyperplane and the origin. The sample training set $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$, $y_i \in \{-1, +1\}$ is provided to separate the different categories in D . A hyperplane should be divided in the sample space represented by D .

This study uses Python program to construct the sentiment classification model. The specific process is as follows:

(1) Text preprocessing

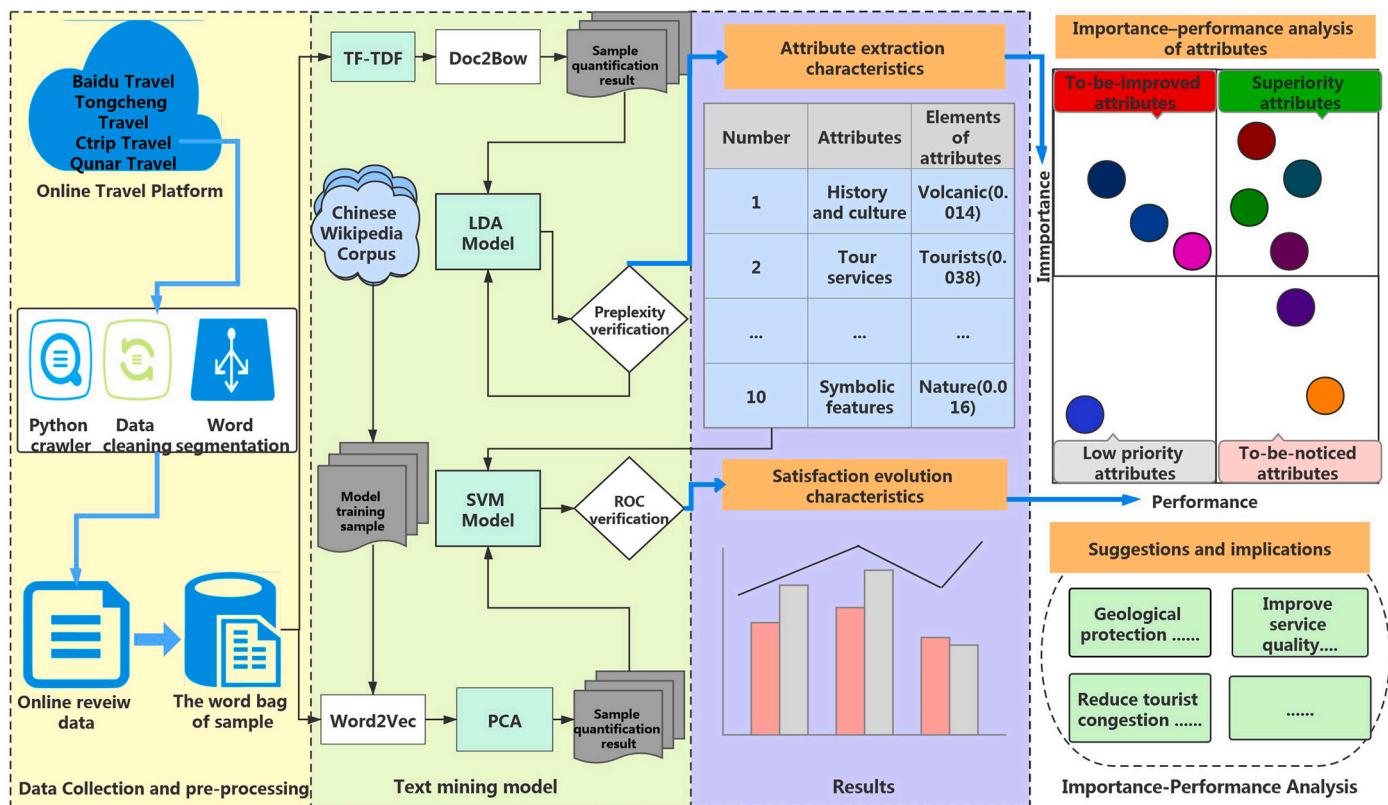


Fig. 4. Online reviews mining framework.

Table 1
24 5A UGGp of China reviews data.

Number	Geopark	Number of reviews
1	Alxa Desert UGGp	2235
2	Arxan UGGp	4150
3	Danxiashan UGGp	4723
4	Dunhuang UGGp	6972
5	Funiushan UGGp	911
6	Huangshan UGGp	15,005
7	Jingpooh UGGp	1315
8	Jiuhuashan UGGp	4790
9	Keketuohai UGGp	882
10	Longhushan UGGp	2699
11	Lushan UGGp	6141
12	Ningde UGGp	7194
13	Sanqingshan UGGp	5314
14	Shennongjia UGGp	2093
15	Shilin UGGp	7973
16	Songshan UGGp	6589
17	Taining UGGp	241
18	Taishan UGGp	10,347
19	Tianzhushan UGGp	5031
20	Wudalianchi UGGp	843
21	Yangdangshan UGGp	2938
22	Yimengshan UGGp	2709
23	Yuntaishan UGGp	4417
24	Zhangjiajie UGGp	15,931
Total		121,443

This study uses Python and Jieba Chinese word segmentation library to complete the processing of the sample text-based data. During word segmentation, this study initially removes all punctuations and various symbols by regularity, obtains pure text, and then loads the stop word library and self-built vocabulary for word segmentation. This study selects the stop words used by the latest version of the NLPPIR participle of the Chinese Academy of Sciences. Concurrently, this study constructs a tourism-specific lexicon on the basis of the existing thesaurus to prevent the division of names of scenic spots, place names, and other words with special characteristics in the tourism field.

(2) Word2Vec model training

The SVM model training only supports numerical samples. The sample text data should be quantized, that is, converted into numerical data. In 2013, Google's Word2Vec open source project triggered a wave of research and application of word vectors (Turian, Ratinov, & Bengio, 2010). Word2Vec is widely used in natural language processing tasks (e.g., text sentiment analysis) as the basic technology for deep learning in the field of natural language processing. This study uses the Wiki Chinese corpus as the original sample. In addition, the study selects the Word2Vec model of the Gensim library for training and uses the skip-gram model. We then input the word bag and set the output vector to 400 dimensions. The window size is set to 5. Afterward, the final Word2Vec model is obtained, serialized to the local destination, and saved as a binary file.

(3) Feature extraction

Representing text by vectors typically faces sparse vector space and high feature dimension problems. In view of this situation, feature dimension reduction processing is needed. Dimensionality reduction can reduce the feature dimension of the script, decrease the number of iterations during model training, and eliminate the features of similar semantics. This reduction improves the accuracy and recall rate of sentiment classification and efficiency. Thus, this study selects the principle components analysis (PCA) method. The PCA, also known as K-L transform, is a linear data analysis method based on statistical properties (Zhou, Lan, Qiang, & Wei, 2013). The PCA method is used to consider the correlation among feature items completely, and the

original feature document matrix is transformed into a lower-dimensional orthogonal feature matrix. This matrix comprises the principal components of the original feature document matrix, retaining most of the feature information from the original feature matrix. Moreover, the matrix ensures that the new features are irrelevant.

Using the PCA algorithm provided in the scikit-learn library, the 400-dimension samples are trained and estimated, and the relationship between the dimension and the variance is represented by matplotlib. The result is shown in Fig. 5.

The graph exhibits a significant transition in the near dimension of approximately 50–80. Thus, this study selects 80 dimensions as the effective dimension of each comment, which is used to express the information of the vector of the comment.

(4) SVM model training

This study manually marks 10,000 positive online reviews and 5000 negative online reviews as the training set. To obtain the ideal model, this study uses Python's scikit-learn library and continuously debugs the parameters of the model. First, the kernel type used in the specified algorithm is a radial kernel function. Second, the penalty parameter C of the error term is set, and the parameter affects the accuracy and generalization ability of the model. Finally, the parameter decides whether to enable probability estimation.

(5) Model test

To verify the actual classification effect of the model, we manually mark 15,000 online tourist reviews. 10,000 of which are positive reviews, and 5000 are negative ones. The model with annotated comments is validated. For the two-class problem, the receiver operating characteristic (ROC) curve shows the correct positive rate and the false positive rate obtained by the classifier prediction. The vertical axis corresponds to the true positive rate (TPR), and the horizontal axis corresponds to the false positive rate (FPR). A graph is drawn from different (FPR, TPR) data pairs, which is called a ROC plot (Zou, O'Malley, & Mauri, 2007). The calculation process is as follows:

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

$$FPR = \frac{FP}{TN + FP} \quad (4)$$

The ROC test curve of the final SVM model trained in this study is shown in Fig. 6.

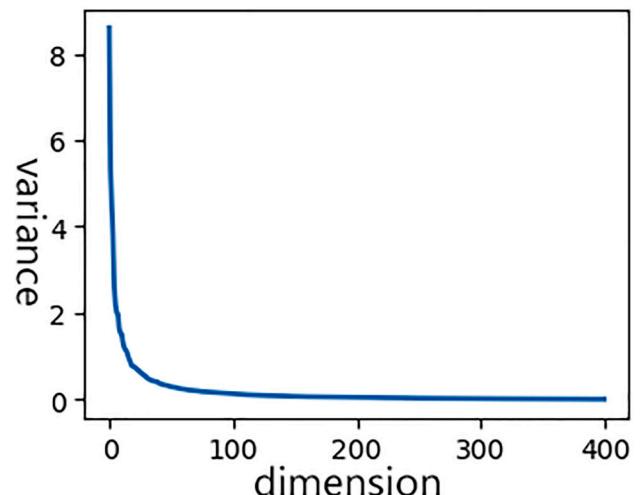


Fig. 5. PCA output result.

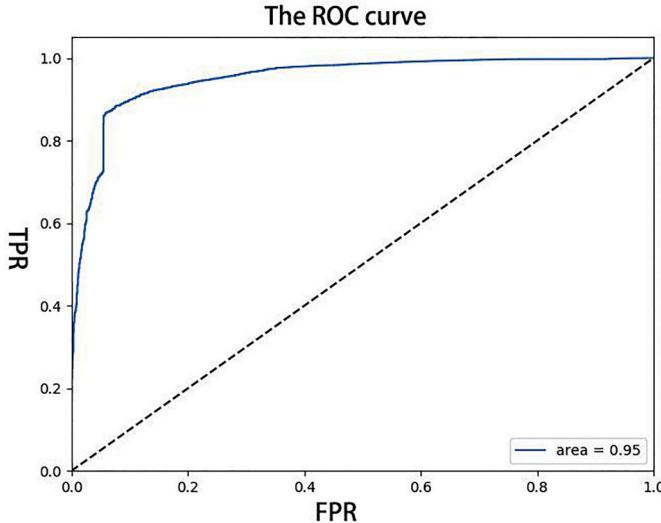


Fig. 6. ROC model test curve.

The output ROC curve combined with the AUC output value of the training model is 0.95. According to the characteristics and theoretical knowledge of the AUC, the model is considered to have high accuracy and is in line with the expectations of this study.

3.3. Improved LDA model

The LDA is a text topic representation method that introduces a full probability model. Its core purpose is to calculate the text posterior theme on the basis of the Dirichlet prior hypothesis of text topic distribution and topic word distribution using Bayesian estimation (Blei, Ng, & Jordan, 2003). The specific structure is depicted in Fig. 7.

The model diagram can be broken down into the following physical processes.

The first step is to select a document d_i according to probability $P(d_i)$. In the second step, the topic distribution θ_m of the document d_i is sampled from the Dirichlet distribution α . The third step is to extract the topic $z_{m,n}$ of the n word of document d_i from the topic distribution θ_m . The fourth step is to sample the generated word distribution β from the LDA φ_k . In the fifth step, the word $w_{m,n}$ is generated by the sampling from the word distribution φ_k . The design formula is as follows:

(1) The following formula is used to generate the probability of the

subject number of all the words:

$$P(\vec{z}|\vec{\alpha}) = \prod_m^M \frac{\Delta(\vec{n}_m + \vec{\alpha})}{\Delta \vec{\beta}} \quad (5)$$

(2) Given that the choice of topic number does not change the word distribution of the topic, the formula can be further expressed as:

$$P(\vec{w}|\vec{z}, \vec{\beta}) = P(\vec{w}|\vec{\beta}) = \prod_k^V \frac{\Delta(\vec{v}_k + \vec{\beta})}{\Delta \vec{\beta}} \quad (6)$$

(3) Finally, the corpus of the LDA model generates probability expressions.

$$P(\vec{w}, \vec{z}|\vec{\alpha}, \vec{\beta}) = \prod_k^V \frac{\Delta(\vec{v}_k + \vec{\beta})}{\Delta \vec{\beta}} \prod_m^M \frac{\Delta(\vec{n}_m + \vec{\alpha})}{\Delta \vec{\alpha}} \quad (7)$$

For evaluating the LDA, Blei et al. (2003) proposed the use of the degree of confusion (perplexity value) as the criterion. The degree of perplexity between the distribution of the probability measure or the probability model and the sample is as follows. For document d , the degree of the trained model depends on which subject the document d belongs to. A small perplexity enhances the effect of the model. The calculation formula for the perplexity is as follows:

$$\text{Perplexity } (w|a_d) = e \left[-\frac{\sum_1^m \ln(p(w_m|a_d))}{\sum_1^m N_d} \right] \quad (8)$$

where w is the test set and w_m is the observable word in the test document m . $p(w_m|a_d)$ is the probability that the text w_m is generated for the model. N_d is the number of words of the document w_m .

On the basis of the LDA model, the text is clustered through an unsupervised topic clustering model. However, this clustering requires the text to have a certain normativeness that is suitable for news and other corpus that have been processed for the first time. The online tourism reviews are short and colloquial, the noise is large, and the semantic information is sparse. Therefore, this study uses the term frequency-inverse document frequency (TF-IDF) method to optimize the method to enhance the clarity of the semantic distribution. TF-IDF (Meng, Lin, & Li, 2011) is a classical statistical method for assessing the importance of entries for texts. The method tends to filter out high-frequency words

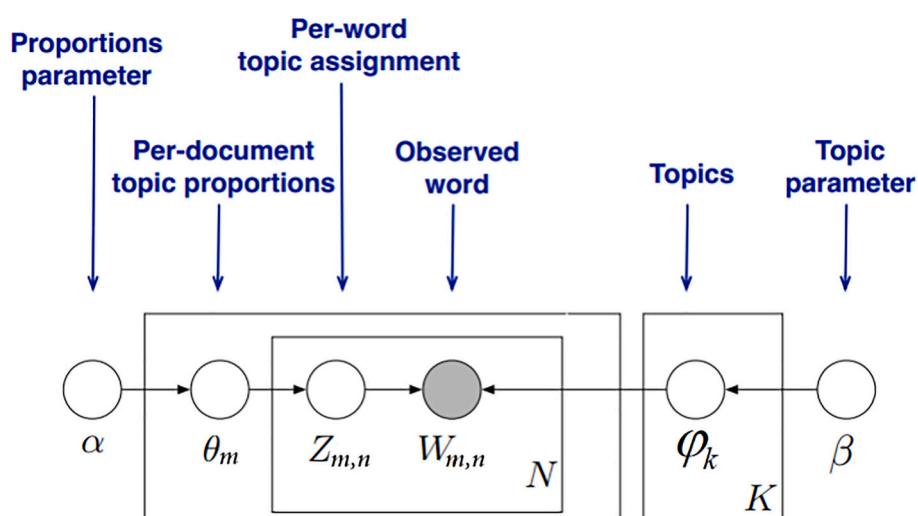


Fig. 7. The LDA model.

with low discrimination and retains low-frequency words with high discrimination. The TF is the frequency of words that appear in a document. Given the different lengths of various documents, these frequency gaps are large. They should be normalized to enable comparisons of frequencies in the same environment, which can be expressed as:

$$tf_{ij} = \frac{n_{ij}}{\sum_k n_{kj}} \quad (9)$$

where tf_{ij} represents the word frequency of the word n_i in the document j . $n_{i,j}$ represents the frequency of the word n_i , and $\sum_k n_{kj}$ is the total number of words contained in the document j .

The IDF represents the weight of a given word. A large weight indicates that the word is important. On the basis of word frequency, light weight is given to common words. Words that have great weight can be expressed as:

$$idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}| + 1} \quad (10)$$

where $|D|$ represents the total number of files in the corpus. $|\{j : t_i \in d_j\}|$ represents the number of files containing words t_i .

Finally, the formula for calculating TF-IDF is:

$$TF - IDF = TF * IDF \quad (11)$$

This study uses the TF-IDF-based improved LDA subject extraction model. The steps are as follows:

- (1) Preprocess the review sample to obtain the sample word bag.
- (2) Use the TF-IDF model method to calculate the degree of influence of each word in the word bag on the text topic. Then, set the threshold to filter to determine the expressive strength of the final theme.
- (3) Use the Gensim NLP handler method doc2bow to quantize the text into a numeric matrix. Then, apply the Gensim NLP handler function to record the mapping between text and identity.
- (4) Using the output of (3) as input, use Gensim NLP to perform model training.
- (5) Determine the output model, model confusion, and topic details.
- (6) Evaluate whether the model passes the output; otherwise, record the parameters, adjust the parameters, and return to (1).

Upon completing the above steps, the model parameter model is modified by evaluating the result theme and confusion. The relationship between the output theme and confusion degree is shown in Fig. 8. This

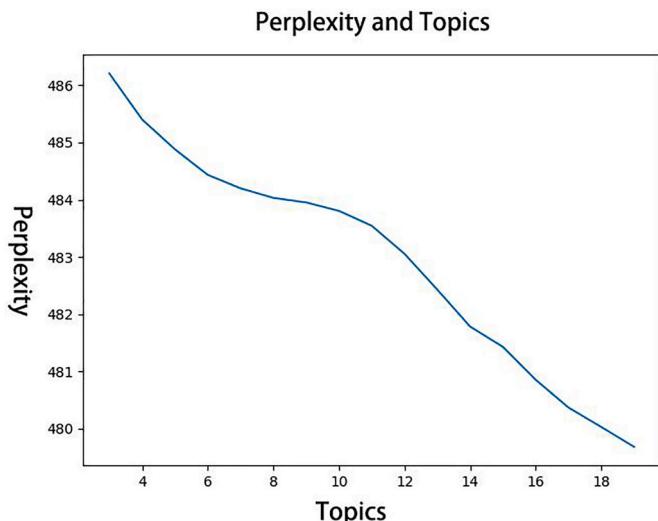


Fig. 8. Relation curve between number of topics and perplexity.

study concludes that the effect is best when the number of topics is 10.

3.4. IPA model

Through the above the LDA and SVM models, we identify the sensory dimension and overall sentiment orientation of tourists. The objective is to understand the perception of tourists toward various attributes and help geoparks in China improve their management and monitor warnings with the perspective of visitors as basis. This study introduces the IPA model. It was proposed by Martilla (Martilla & James, 1977) and was soon applied to the tourism and catering industries. The managers of geoparks can observe the satisfaction of tourists, as the results obtained by the IPA method are close to the needs of the tourism market. The basic idea is to determine the significance of each factor improvement by comparing the importance and actual performance of each factor. Thus, managers can use limited resources on the “blade,” a method to optimize resources effectively. In this study, the themes extracted by the LDA model and their thematic elements are used as the objects. The methods are as follows:

(1) Importance: Using the LDA model completed by the training, the subject reviews are separately discriminated, and the topic tendencies of each comment are refined. Moreover, the distribution of topics in all samples is counted. The attribute distribution frequency is then used as the importance indicator.

(2) Performance: People unconsciously become positive on social networks (Qiu, Lin, Leung, & Tov, 2012). The positive reviews in the emotional classification results are far greater than the negative reviews. If the number of times a word appears in a positive comment is high enough, then negative topics may also be identified as positive topics. Such a circumstance is not conducive to the mining of negative topics. Interestingly, negative reviews have a strong impact (El-Said, 2020) and are more persuasive than positive reviews (Park, Kim, & Ryu, 2019). Therefore, we need to strengthen the recognition of negative attributes. We calculate the frequency of occurrence of positive and negative comments in each topic through strict probability statistics based on the LDA model. Topics that appear more frequently in negative comments than in positive comments are categorized as negative topics. At the same time, we randomly selects the same number of negative and positive comments to alleviate the analysis error caused by sample imbalance.

4. Results

4.1. Satisfaction evolution characteristics

This study uses the SVM model to classify the sentiment orientation for mining the emotional sentiment of visitors directly from the text review. In this step, the SVM model directly takes an online review as a whole and judges the emotional tendency of each review. The model reflects the comprehensive satisfaction of tourists with a geopark and does not involve the analysis of specific attributes. In addition, this study divides the online reviews by year from 2010 to 2019. Among them, the amount of data from 2010 to 2012 is small, so such data are omitted. The final positive reviews are 97,277, and the negative reviews are 23,255. The satisfaction rate is approximately 80.70%. Fig. 9 shows that the overall satisfaction of tourists has been volatile in recent years. The number of tourists in the domestic tourism market has been steadily rising with the increasing of population. The number of tourists is increasing, on the one hand, so higher requirements are put forward on the geoparks. On the other hand, it brings certain pressure on the management of the geoparks. At the same time, with the development of the mobile Internet, growing number of people choose to express their feelings through online platforms. In 2017, the satisfaction of geoparks gradually declined. A significant downward trend in satisfaction was noted in 2018, when geological disasters in this year were serious (People's Daily Online, 2018). Geoparks are more susceptible to natural

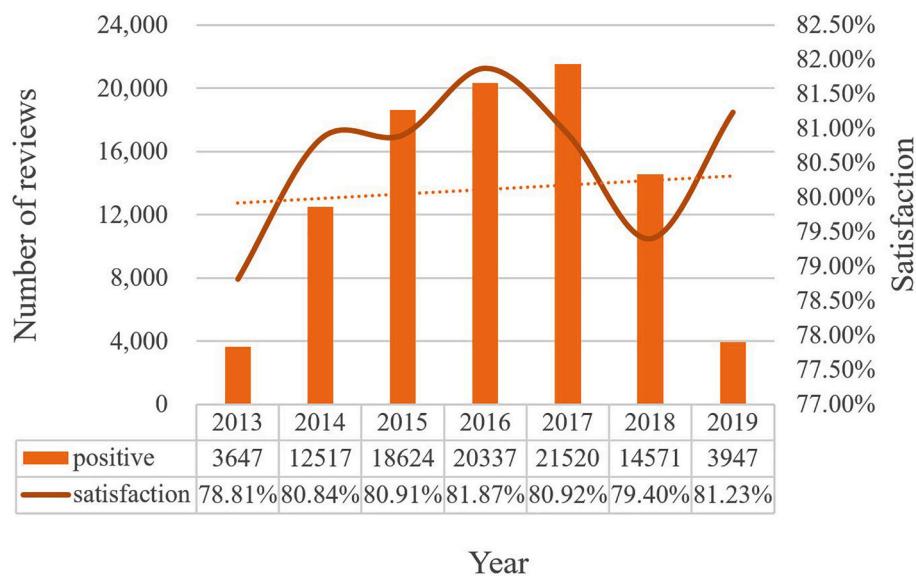


Fig. 9. Satisfaction evolution from 2013 to 2019.

disasters due to their geomorphic features, which can be inferred to be related to the decline in satisfaction. Moreover, the managers of geoparks were actively responding to the changing needs of tourists and striving to improve the competitiveness of these spots. During this period, tourists were less satisfied with the 5A UGGp in China. Nevertheless, an upward trend was observed in 2019. The scenic spots should seize development opportunities by promoting geotourism further.

4.2. Attribute extraction characteristics

This study uses the LDA model to obtain the final ten attributes found in tourist reviews, as shown in Table 2. We extract 10 attributes and 100 elements of attributes. We artificially remove meaningless elements of attributes and finally retained 80 elements of attributes. The first column briefly summarizes the meaning of each topic tag. The second column shows the most frequent and exclusive words related to the attributes, that is, the words that appear most frequently in the attribute but the lowest in other attributes. The numbers in the parentheses are the weights of the attributes, indicating the importance of each word in the attribute.

The attributes mentioned in Table 2 comprehensively reflect the attributes and elements perceived by tourists during tourism. Such attributes show certain emotional attitudes, indicating that the model proposed in this study works effectively. The following section uses the IPA direction to analyze the results further and to clarify the emotional orientation of visitors and the importance of each attributes.

4.3. Importance–performance analysis

In this study, the IPA method is used in combination with the SVM and LDA models to obtain the saliency and expressiveness of each attribute and element of attributes. The average of the performance and importance mean is divided into IPA quadrants. Figs. 10 and 11 show the results. In Fig. 11, different colors denote different attributes, and the size of the graph points the importance of the attribute. The larger the graph dots, the higher the importance of the attribute is. We retain the same elements under different attributes, because they may have different importance or performance. For example, the elements “tickets” and “by tickets” may represent a similar meaning in Chinese semantics, but because they are under different attributes, it means that tourists have varying views on the same thing under distinct circumstances. In the “travel cost,” it may express the views of tourists on fares,

Table 2
Attributes extraction.

Number	Attributes	Elements of attributes (Top words)
1	History and culture	Volcanic(0.014), natural(0.014), reputation (0.014), mountain view(0.014), culture (0.014), famous(0.013), altitude(0.011), history(0.010)
2	Tour services	Tourists(0.038), management(0.029), staff member(0.019), queue(0.019), tourist guide (0.019), tickets(0.015), Internet(0.013), buy tickets(0.010)
3	Well-known degree	Country(0.024), landscape(0.020), tourism (0.016), geoparks(0.015), cruise(0.014), national level(0.014), focal point(0.013), nationwide(0.011)
4	Travel cost	Tickets(0.055), hour(0.014), good(0.012), queue(0.011), time(0.010), deserve(0.009), price(0.009), cheap(0.009)
5	Way of tour	Cableway(0.046), hour(0.032), cable car (0.018), mountain(0.013), time(0.012), the next day(0.012), choose(0.012), walk(0.010)
6	Internet service	Check-in(0.047), queue(0.026), service (0.021), online(0.020), booking(0.018), buy tickets(0.014), tickets(0.013), service window(0.012)
7	Natural scenery	Canyon(0.029), landscape(0.022), scenery (0.017), geology(0.016), nature(0.015), amazing pines(0.014), fantastic rock peaks (0.014), landform(0.014)
8	Transportation and accommodation	Parking lot(0.018), driver(0.017), taste (0.013), transfer(0.011), eat(0.010), service area(0.010), accommodation(0.009), hotel (0.009)
9	Emotional experience	Pretty(0.056), landscape(0.048), deserve (0.045), weather(0.026), mountain climbing (0.017), sea of clouds(0.012), regret(0.012), like(0.011)
10	Symbolic features	Nature(0.016), stone(0.014), mountain peak (0.010), clouds(0.009), amazing(0.009), uncanny workmanship(0.009), volcanic rock (0.008), landform(0.008)

and it is the most important, indicating that tourists have the highest perception of ticket expenditure and price; In the “tour services,” it may express tourists’ views on the ticketing service of the scenic spot. In the “Internet service,” it shows tourists’ attitude on online ticket purchase. In addition, regardless of the attributes, the performance of tickets is

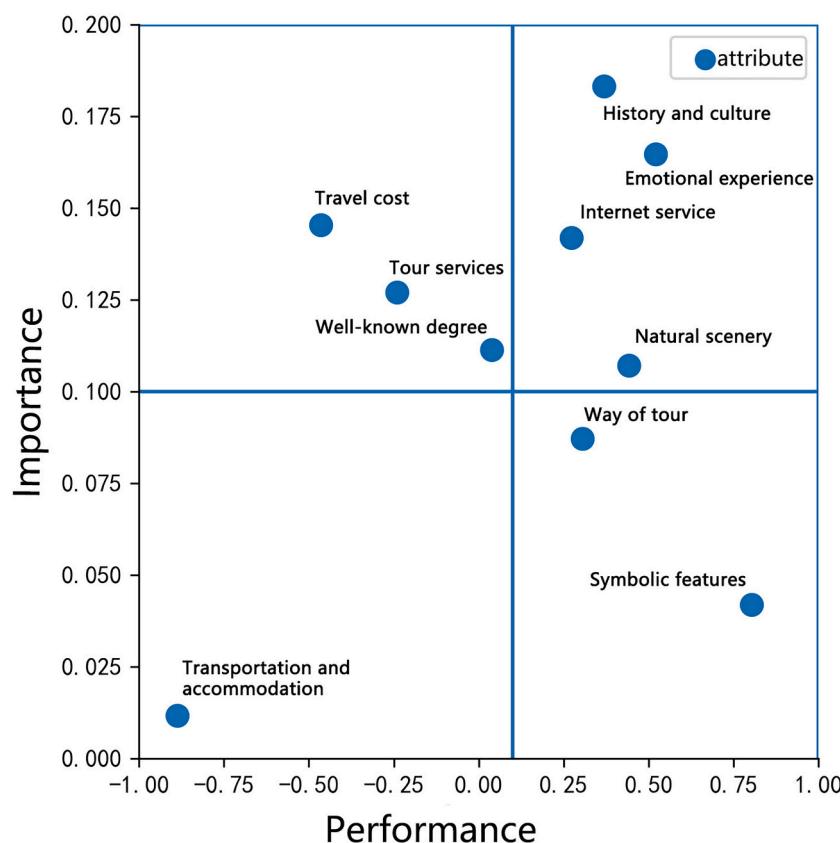


Fig. 10. Importance–performance analysis of attributes.

low, reflecting that tourists are always dissatisfied with tickets.

(1) Superiority attributes of geoparks are presented in the first quadrant. The satisfaction and significance of tourists are high. The attributes in this quadrant include "history and culture," "emotional experience," "Internet service," and "natural scenery." The elements of attributes include "landscape," "volcanic," "history," "deserve," and "reputation." This quadrant plays an important role in prompting the satisfaction of tourists. At the same time, its performance meets the expectations of tourists because the satisfaction is high. The UGGp has a high resource endowment, and tourists recognize its natural scenery, history, and cultural significance. Moreover, tourists enjoy a better experience here. With the development of online travel platforms, the ticket collection service has become convenient, thus improving the experience for tourists. Scenic spot management personnel should continue to maintain good landscape quality, strengthen the popularization and education of geosciences, and maintain the efficiency of ticket collection.

(2) Attributes that geoparks need to improve are revealed in the second quadrant. For tourists, the satisfaction is low, but the significance is relatively high. The attributes in this quadrant are "travel cost," "tour services," and "well-known degree." The elements include "tickets," "queue, time," and "tourist guide." The elements in this quadrant are the key improvements. The main dissatisfaction of tourists lies in the time and cost of travel and the service of scenic spots. As the number of tourists grows, the reception capacity of scenic spots becomes increasingly overwhelmed. Congestion in scenic spots occurs frequently. At the same time, the quality of service provided by scenic spots declines, which greatly reduces the tourist experience. In future developments, scenic spots should reasonably determine the value of tickets, effectively improve the congestion situation, and improve the service level of their personnel to improve the overall satisfaction of these spots. In addition, publicity should be strengthened to reinforce the popularity of geoparks.

(3) Low priority attributes of geoparks are shown in the third

quadrant. For tourists, both the satisfaction and the significance are low. The attribute in this quadrant is "transportation and accommodation." The elements include "transfer," "food," "hotel," "weather," "parking lot," and "driver." The elements in this quadrant are low priority and do not need to be focused on development, but not totally unnecessary. They are not suitable for priority development in the case of limited resources. As the market changes in the future, these elements can be further upgraded after the opportunity is ripe. At the same time, the bus station, railway station, hotel, and other factors are not among the resources of scenic spots but rather the factors involved in the entire travel process. Relying on the scenic spots alone to improve the development is difficult.

(4) Attributes to which geoparks need to pay attention are shown in the fourth quadrant. For tourists, the satisfaction is high, but the significance is relatively low. The attributes in this quadrant are the "way of tour" and "symbolic features." The elements include the "cableway," "cable car," "view of the mountain," and "sea of clouds." The tourists are satisfied with the elements in this quadrant, but do not attach significance to them. Although the symbolic characteristics of tourists in the National Geopark are recognized by tourists, the perception of tourists is low. The above factors can be further excavated to transform them into the competitive advantages of scenic spots. However, investing much time and energy are not necessary to maintain and improve these spots in a timely manner.

4.4. Suggestions and implications

The establishment of a geopark has three main purposes: protecting geological relics, popularizing geological knowledge, and developing tourism and promoting local economic development. This study proposes the following development recommendations to achieve three aforementioned purposes effectively.

(1) Pay attention to geological protection and maintain landscape

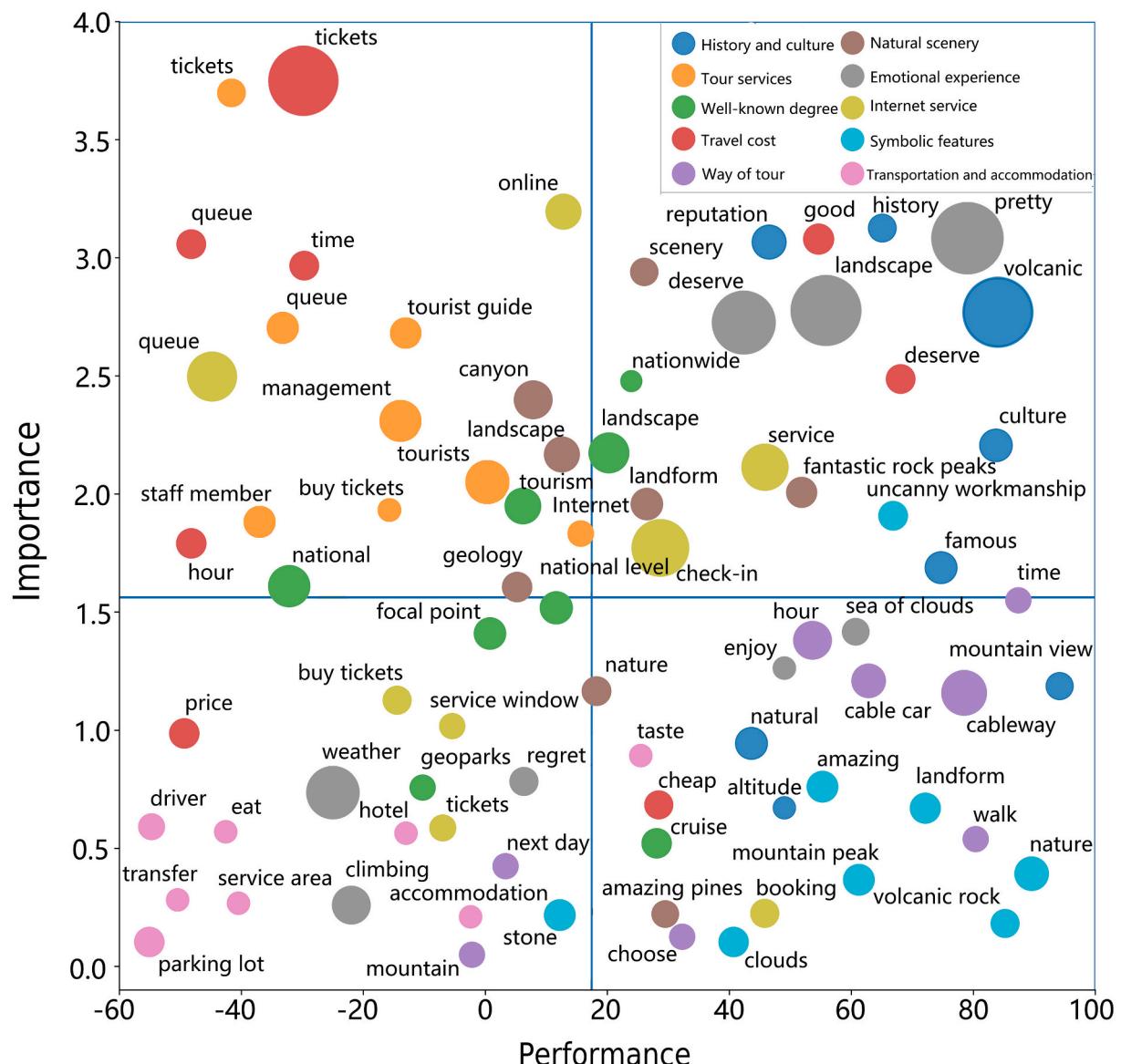


Fig. 11. Importance–performance analysis of attributes elements.

development. The development of geotourism should be based on geological protection. Geological protection should always be prioritized. The IPA results show that attributes, such as natural scenery and emotional experience, are in the first quadrant, indicating that geoparks have rich natural resource endowment and high tourist satisfaction. The development of geoparks should focus on the existing natural scenery and geological protection. In addition, the rational use of resources should be emphasized, and landscape development should be carried out in moderation.

(2) Improve service quality and control ticket prices reasonably. The IPA results show that tourists express dissatisfaction with tour services and travel cost. These attributes have a high degree of perception by tourists, and they should be improved. Geoparks need to strengthen the management of the staff. While improving the service quality, increasing the geological and humanistic knowledge reserves of service personnel and enhancing the level of service are necessary. Reducing the cost of scenic spots and appropriately reducing the price of tickets should also be considered.

(3) Reduce tourist congestion and improve the tourist experience. With the development of geoparks, the number of tourists has increased annually. However, the congestion in geoparks has been criticized. The

above analysis shows that the importance of tourists toward attributes, such as queue and time, is very strong, but the performance of tourists is low. Geoparks can reduce tourist congestion by controlling the flow of scenic spots, rationally planning routes, and promoting free public transportation. Within suitable conditions, scenic spots can provide direct trips to open scenic spots with local bus stations and train stations. Moreover, scenic transportation lines should be open among various scenic spots.

(4) Increase the popularization of science education and experience products of geotourism and highlight the symbolic features of geoparks. One of the main purposes of developing geoparks is to popularize geo-science knowledge. However, popular science education in geoparks has been seriously criticized by the public. These tourist perception dimensions point that tourists have low perception of geosciences and other content, well-known degree attributes, and perception of symbolic features. Thus, promoting information on these aspects must be further strengthened. The construction of geoparks requires strengthening of the marketing promotion of geopark education, increasing the address of tourism experience products, and stimulating the motivation for science tourism.

(5) Strengthen the construction of infrastructure and information

technology. Infrastructure is a necessary element for the development of tourism, and it can promote the comparative and competitive advantages of tourist destinations. According to the IPA results, Internet service highly satisfies tourists. In the context of smart tourism construction, scenic spots can further strengthen information technology construction, such as adding an intelligent voice navigation system and providing free WiFi so that tourists can obtain additional information on scenic spots and understand the general knowledge on geosciences.

5. Conclusions

Given the important strategic position of geoparks and geotourism, and at the same time, certain problems they are facing, such as insufficient geoscience popularization and crowd congestion, this study starts from the connotation and requirements of geoparks and geotourism and proposes a multi-dimensional evaluation of the development status of UGGp from the perspective of tourists using online reviews. This study further proposes a model for geopark management strategies based on “data acquisition and cleaning, data mining and analysis, and strategy formation” to help managers better understand the geoparks, discover the advantages and disadvantages during development, and provide references for other geoparks.

With regard to theoretical implications, the current method the LDA model and the SVM model provide a sound technological support for feature extraction and emotion classification, but there are some shortcomings. Combining the disadvantages and advantages of the SVM, the LDA and the IPA models, and the TF-IDF and PCA algorithms, the objective is to build an online reviews mining model to understand the tourist expression, specially the negative attributes. The negative attributes identified are not extracted only from negative reviews. Such attributes are found more frequently in negative reviews than in positive reviews with statistical significance. First, the SVM and LDA models and algorithm improvement based on the characteristics of geoparks to identify emotional tendencies and contribute attributes from massive online reviews. Second, the IPA model is provided to help managers understand the importance and performance of each attributes. Lastly, this study puts forward suggestions on geotourism and provides new ideas for geopark management and planning. The model evaluation proves that the model achieves good results.

In the managerial implications, this study shows that the overall satisfaction of 120,532 reviews of 24 5A UGGp in China is 80.70%. We identify 10 attributes and 80 elements of attributes. The main negative attributes are “travel cost,” “tour services,” “well-known degree,” and “transportation and accommodation.” At the same time, we reveal the importance of each attributes and attributes elements. Visitors have high perceptions toward “attractions,” “tickets,” “time,” “queue,” and “tourist guides,” the more important elements are consistent with the negative attributes. This proves that the model proposed in this study can effectively identify negative topics. At last, the results of online tourism review analysis show that the geological experience products, local disciplines, and geological protection performance are insufficient, reflecting that the current geoparks are lacking in scope and necessitate further improvement. We suggest that the development of geoparks should increase geoscience popularization and geoscience tourism experience projects based on geological protection. At the same time, from the tourist experience level, reducing ticket prices appropriately, improving service quality, relieving tourist congestion, and enhancing infrastructure construction are necessary.

This study has certain limitations that can be explored in future research. First, the data search volume of the exclusive vocabulary of geotourism is not extensive enough. It can be further expanded to demonstrate the uniqueness of the vocabulary used in various fields. Second, the topics identified in this study are from specific time points and comments, and the subject matter may change as the text data changes. A machine-based online reviews analysis and analysis platform

can be developed in the future to improve the satisfaction of tourists. Moreover, given the dynamics of visitors' needs over time, scenic attributes, and visitor characteristics can be incorporated into the model to monitor trends in the perceived dimensions of visitors easily.

Acknowledgments

This work was supported by the Humanities and Social Sciences Program of the Ministry of Education of the People's Republic of China (Grant Nos. 20YJC630095, 19YJC630119), the National Natural Science Foundation of China (Grant Nos. 71501019, 71971151), Postdoctoral Research Foundation of China (Grant No. 2018M631069), the Philosophy and Social Science Planning Program of Chengdu (Grant No. 2018A09), General Project of Regional Public Management Informationization Research Center of Key Social Science Research Base in Sichuan (Grant No. QGXH20-03), the Funding Program for Middle aged Core Teachers in Chengdu University of Technology (Grant No. 2019KY37-04203), Philosophy and Social Science Research Foundation of Chengdu University of Technology (Grant No. YJ2019-NS004), the Program of Graduate Education Reform Program at Chengdu University of Technology (Grant No.10800-00009824) and Key Project of National Park Research Center of Sichuan Social Science Research Base (Grant No. GJGY2020-ZD001) “Research on the path to improve tourism ecological value of Giant Panda National Park from the perspective of host-guest governance.”

Credit Author Statement

Yuyan Luo: Yuyan contributed this paper in theory development, mining model design, drafting and revising the manuscript.

Jinjie He: Jinjie contributed to this paper on the data processing and analysis, drafting and revising the manuscript.

Yu Mou: Yu contributed to the model development, data analysis, and revision of this manuscript.

Jun Wang: Jun organized this research and made substantial contributions on the research design, manuscript writing and revision.

Tao Liu: Tao contributed to this paper on collecting and organizing research-related data and materials, and revising the manuscript.

Declaration of Competing Interest

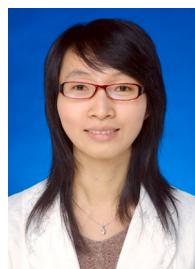
None.

References

- Abukausar, M., Dhaka, V. S., & Singh, S. K. (2013). Web crawler: a review. *International Journal of Computer Applications*, 63, 31–36.
- Ahani, A., Nilashi, M., Yadegaridehkordi, E., Sanzogni, L., Tarik, A. R., Knox, K., ... Ibrahim, O. (2019). Revealing customers' satisfaction and preferences through online review analysis: The case of Canary Islands hotels. *Journal of Retailing and Consumer Services*, 51, 331–343.
- Alexandrowicz, Z. (2006). Geopark—nature protection category aiding the promotion of geotourism (Polish perspectives). *Geoturystyka*, 2, 3–12.
- Balazs, J. A., & Velásquez, J. (2016). Opinion mining and information fusion: A survey. *Information Fusion*, 27, 95–110.
- Benado, J., Herve, F., Schilling, M., & Brilha, J. (2019). Geoconservation in Chile: State of the art and analysis. *Geoheritage*, 11, 793–807.
- Beretić, N., Dukanović, Z., & Cecchini, A. (2019). Geotourism as a development tool of the geo-mining Park in Sardinia. *Geoheritage*, 11(3), 1–16.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Brocx, M., Brown, C., & Semeniuk, V. (2019). Geoheritage importance of stratigraphic type sections, type localities and reference sites—Review, discussion and protocols for geoconservation. *Australian Journal of Earth Sciences*, 66, 1–14.
- Brocx, M., Semeniuk, V., & Meney, K. (2019). Geoheritage and geoconservation in Australia: Introduction. *Australian Journal of Earth Sciences*, 66, 751–752.
- Carvalhido, R. J., Brilha, J. B., & Pereira, D. I. (2016). Designation of natural monuments by the local administration: The example of Viana do Castelo municipality and its engagement with Geoconservation (NW Portugal). *Geoheritage*, 8, 279–290.
- Cetin, M., Zeren, I., Sevik, H., Cakir, C., & Akpinar, H. (2018). A study on the determination of the natural park's sustainable tourism potential. *Environmental Monitoring and Assessment*, 190(3), Article 167.

- Chen, Y., Zhang, H., Liu, R., Ye, Z., & Lin, J. (2019). Experimental explorations on short text topic mining between LDA and NMF based schemes. *Knowledge-Based Systems*, 163, 1–13.
- Cheng, X., Fu, S., Sun, J., Bilgihan, A., & Okumus, F. (2019). An investigation on online reviews in sharing economy driven hospitality platforms: A viewpoint of trust. *Tourism Management*, 71, 366–377.
- China National Tourism Administration. (2003). Standard of rating for quality of tourist attractions. Retrieved from <http://openstd.samr.gov.cn/bzgk/gb/newGblInfo?hcno=511242B712ED75EAFB5FFD34C8154E1B>.
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3), 273–297.
- Dowling, R. K. (2011). Geotourism's global growth. *Geoheritage*, 3, 1–13.
- Ehsan, S., Leman, M. S., & Begum, R. A. (2013). Geotourism: A tool for sustainable development of Geoheritage resources. *Advanced Materials Research*, 622–623, 1711–1715.
- El-Said, O. A. (2020). Impact of online reviews on hotel booking intention: The moderating role of brand image, star category, and price. *Tourism Management Perspectives*, 33, Article 100604.
- Esfehani, M. H., & Albrecht, J. N. (2019). Planning for intangible cultural heritage in tourism: Challenges and implications. *Journal of Hospitality & Tourism Research*, 43, 980–1001.
- Fanwei, Z. (2014). An evaluation of residents' perceptions of the creation of a geopark: A case study on the geopark in Mt. Huaying grand canyon, Sichuan Province, China. *Environmental Earth Sciences*, 71, 1453–1463.
- Fuchs, M., Höpken, W., & Lexhagen, M. (2014). Big data analytics for knowledge generation in tourism destinations – A case from Sweden. *Journal of Destination Marketing & Management*, 3(4), 198–209.
- Garay, L. (2019). #Visitspain: Breaking down affective and cognitive attributes in the social media construction of the tourist destination image. *Tourism Management Perspectives*, 32, Article 100560.
- Guo, W., & Chung, S. (2019). Using tourism carrying capacity to strengthen UNESCO global Geopark management in Hong Kong. *Geoheritage*, 11, 193–205.
- Guo, Y., Barnes, S. J., & Jia, Q. (2017). Mining meaning from online ratings and reviews: Tourist satisfaction analysis using latent dirichlet allocation. *Tourism Management*, 59, 467–483.
- Han, J., Wu, F., Tian, M., & Li, W. (2017). From Geopark to sustainable development: Heritage conservation and Geotourism promotion in the Huangshan UNESCO global Geopark (China). *Geoheritage*, 10, 1–13.
- Hose, T. (1995). Selling the story of Britain's stone. *Environmental Interpretation*, 10, 16–17.
- Hose, T. (2006). Geotourism and interpretation. In R. Dowling, & D. Newsome (Eds.), *Geotourism, sustainability, impacts and opportunities* (pp. 221–241). Oxford: Elsevier. Reino Unido.
- Hu, M., & Liu, B. (2004). Mining and summarizing customer reviews. *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Seattle, Washington, USA, August*, 168–177.
- Jia, S. (2020). Motivation and satisfaction of Chinese and U.S. tourists in restaurants: A cross-cultural text mining of online reviews. *Tourism Management*, 78, Article 104071.
- Khorsand, R., Rafiee, M., & Kayvanfar, V. (2020). Insights into TripAdvisor's online reviews: The case of Tehran's hotels. *Tourism Management Perspectives*, 34, Article 100673.
- Kim, S., Park, H., & Lee, J. (2020). Word2vec-based latent semantic analysis (W2V-LSA) for topic modeling: A study on blockchain technology trend analysis. *Expert Systems with Applications*, 152, Article 113401.
- Liang, X., Niu, Q., Qu, J., Liu, B., Liu, B., Zhai, X., & Niu, B. (2019). Applying end-member modeling to extricate the sedimentary environment of yardang strata in the Dunhuang Yardang National Geopark, northwestern China. *Catena*, 180, 238–251.
- Liu, Y., Bi, J. W., & Fan, Z. P. (2017). A method for multi-class sentiment classification based on an improved one-vs-one (OVO) strategy and the support vector machine (SVM) algorithm. *Information Sciences*, 394–395, 38–52.
- Marcopoulos, G., Mikros, G., Iliadi, A., & Liotatos, M. (2015). Sentiment analysis of hotel reviews in Greek: A comparison of unigram features. In V. Katsoni (Ed.), *Cultural tourism in a digital era* (pp. 373–383). Springer Proceedings in Business and Economics.
- Martilla, J. A., & James, J. C. (1977). Importance-performance analysis. *Journal of Marketing*, 41, 77–79.
- Medhat, W., Hassan, A., & Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5, 1093–1113.
- Mellinas, J. P., Nicolau, J. L., & Park, S. (2019). Inconsistent behavior in online consumer reviews: The effects of hotel attribute ratings on location. *Tourism Management*, 71, 421–427.
- Meng, J., Lin, H., & Li, Y. (2011). Knowledge transfer based on feature representation mapping for text classification. *Expert Systems With Applications*, 38(8), 10562–10567.
- Moraes, R., Valiati, J. F., & Gaviao Neto, W. P. (2013). Document-level sentiment classification: An empirical comparison between SVM and ANN. *Expert Systems with Applications*, 40, 621–633.
- Mowlaei, M. E., Abadeh, M. S., & Keshavarz, H. (2020). Aspect-based sentiment analysis using adaptive aspect-based lexicons. *Expert Systems with Applications*, 148, Article 113234.
- Neidhardt, J., Rümmele, N., & Werthner, H. (2017). Predicting happiness: User interactions and sentiment analysis in an online travel forum. *Information Technology & Tourism*, 17(1), 101–119.
- Newsome, D., Dowling, R., & Leung, Y. F. (2012). The nature and management of geotourism: A case study of two established iconic geotourism destinations. *Tourism Management Perspectives*, 2–3, 19–27.
- Nguyen, K. A., & Coudounaris, D. N. (2015). The mechanism of online review management: A qualitative study. *Tourism Management Perspectives*, 16, 163–175.
- Nikolova, V., & Sinnovskiy, D. (2019). Geoparks in the legal framework of the EU countries. *Tourism Management Perspectives*, 29, 141–147.
- Nobre da Silva, M. L., Leite do Nascimento, M. A., & Mansur, K. L. (2019). Quantitative assessments of geodiversity in the area of the Serido Geopark project, Northeast Brazil: Grid and centroid analysis. *Geoheritage*, 11, 1177–1186.
- Pang, B., Lee, L., & Vaithyanathan, S. (2002). Thumbs up? Sentiment Classification using Machine Learning Techniques. *Empirical Methods in Natural Language Processing*, 10, 79–86.
- Park, O. J., Kim, M. G., & Ryu, J. H. (2019). Interface effects of online media on tourists' attitude changes. *Tourism Management Perspectives*, 30, 262–274.
- Pazari, F., & Dollma, M. (2019). Geotourism potential of Zall Gjoçaj national park and the area nearby. *International Journal of Geoheritage and Parks*, 7, 103–110.
- People's Daily Online. Ministry of Natural Resources: A total of 2,966 geological disasters occurred across the country in 2018, resulting in 105 deaths. Retrieved from <http://sn.people.com.cn/n2/2019/0109/c37827-32510225.html>.
- Perotti, L., Carraro, G., Giardino, M., De Luca, D. A., & Lasagna, M. (2019). Geodiversity evaluation and water resources in the Sesia Val Grande UNESCO Geopark (Italy). *Water*, 11(10), Article 2102.
- Pournarakis, D. E., Sotiropoulos, D. N., & Giaglis, G. M. (2017). A computational model for mining consumer perceptions in social media. *Decision Support Systems*, 93, 98–110.
- Purdie, H., Hutton, J. H., Stewart, E., & Espiner, S. (2020). Implications of a changing alpine environment for geotourism: A case study from Aoraki/Mount Cook, New Zealand. *Journal of Outdoor Recreation and Tourism*, 29, 100235.
- Qiu, L., Lin, H., Leung, A. K., & Tov, W. (2012). Putting their best foot forward: Emotional disclosure on Facebook. *Cyberpsychology, Behavior and Social Networking*, 15, 569–572.
- Ruban, D. A. (2015). Geotourism — A geographical review of the literature. *Tourism Management Perspectives*, 15, 1–15.
- Sallam, E. S., Ruban, D. A., Mostafa, M. T., Elkholery, M. K., Alwili, R. L., Molchanova, T. K., & Zorina, S. O. (2020). Unique desert caves as a valuable geological resource: First detailed geological heritage assessment of the Sannur Cave, Egypt. *Arabian Journal of Geosciences*, 13(3), Article 141.
- Sever, I. (2015). Importance-performance analysis: A valid management tool? *Tourism Management*, 48, 43–53.
- Shahnoosineh, I., Modabberi, S., & Shahabi, M. (2017). Study of factors influencing the attitude of local people toward geotourism development in Qeshm National Geopark, Iran. *Geoheritage*, 9, 35–48.
- Shekhar, S., Kumar, P., Chauhan, G., & Thakkar, M. G. (2019). Conservation and sustainable development of geoheritage, geopark, and geotourism: a case study of cenozoic successions of Western Kutch, India. *Geoheritage*, 11, 1475–1488.
- Sparks, B. A., & Browning, V. (2011). The impact of online reviews on hotel booking intentions and perception of trust. *Tourism Management*, 32, 1310–1323.
- Stokes, A.M., Cook, S.D., & Drew, D. (2003). Geotourism: The new trend in travel. Travel Industry America and National Geographic Traveler.
- Strba, L., Krsak, B., & Sidor, C. (2018). Some comments to geosites assessment, visitors, and geotourism sustainability. *Sustainability*, 10(8), Article 2589.
- Tsai, C. F., Chen, K., Hu, Y. H., & Chen, W. K. (2020). Improving text summarization of online hotel reviews with review helpfulness and sentiment. *Tourism Management*, 80, Article 104122.
- Turian, J. P., Ratinov, L. A., & Bengio, Y. (2010). Word representations: A simple and general method for semi-supervised learning. *ACL 2010, Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 384–394.
- UNESCO. (2014). Retrieved from <http://www.unesco.org/new/en/natural-sciences/environment/earth-sciences/unesco-global-geoparks/list-of-unesco-global-geoparks/>.
- UNESCO. (2020). List of UNESCO global geoparks (UGGp). <http://www.unesco.org/new/en/natural-sciences/environment/earth-sciences/unesco-global-geoparks/list-of-unesco-global-geoparks/>.
- Vermeulen, I. E., & Seegers, D. (2009). Tried and tested: The impact of online hotel reviews on consumer consideration. *Tourism Management*, 30, 123–127.
- Wang, L., Tian, M., & Wang, L. (2015). Geodiversity, geoconservation and geotourism in Hong Kong global Geopark of China. *Proceedings of the Geologists' Association*, 126, 426–437.
- Wang, Y., Wu, F., Li, X., & Chen, L. (2019). Geotourism, geoconservation, and geodiversity along the belt and road: A case study of Dunhuang UNESCO global Geopark in China. *Proceedings of the Geologists' Association*, 130, 232–241.
- Wei, R. (2006). Analysis of the research subject of information science based on the keyword. *Information Science*, 24, 1400–1404.
- Wong, C. U. I., & Qi, S. (2017). Tracking the evolution of a destination's image by text-mining online reviews - The case of Macau. *Tourism Management Perspectives*, 23, 19–29.
- Xiang, Z., Du, Q., Ma, Y., & Fan, W. (2017). A comparative analysis of major online review platforms: Implications for social media analytics in hospitality and tourism. *Tourism Management*, 58, 51–65.
- Xinhua News. (2013). Three famous scenic spots in China are investigated by UNESCO. Retrieved from http://www.gov.cn/jrzq/2013-01/12/content_2310677.htm.
- Xinhua News. (2019). Country's geopark construction has become an important starting point for poverty alleviation. Retrieved from http://www.gov.cn/xinwen/2019-10/18/content_5442054.htm.

- Xinhua News. (2020). Continue to exceed the limit Huangshan Scenic Area stops tourists from entering the park on the same day. Retrieved from <http://js.people.com.cn/GB/n2/2020/0406/c359574-33929021.html>.
- Ye, Q., Law, R., & Gu, B. (2008). The impact of online user reviews on hotel room sales. *International Journal of Hospitality Management*, 28(1), 180–182.
- Yuan, H., Lau, R. Y. K., & Xu, W. (2016). The determinants of crowdfunding success: A semantic text analytics approach. *Decision Support Systems*, 91, 67–76.
- Zhang, D., Xu, H., Su, Z., & Xu, Y. (2015). Chinese comments sentiment classification based on word2vec and SVMperf. *Expert Systems with Applications*, 42, 1857–1863.
- Zhao, Y., Xu, X., & Wang, M. (2019). Predicting overall customer satisfaction: Big data evidence from hotel online textual reviews. *International Journal of Hospitality Management*, 76, 111–121.
- Zhou, C., Lan, W., Qiang, Z., & Wei, X. (2013). Face recognition based on PCA image reconstruction and LDA. *Optik - International Journal for Light and Electron Optics*, 124, 5599–5603.
- Zou, K. H., O'Malley, A. J., & Mauri, L. (2007). Receiver-operating characteristic analysis for evaluating diagnostic tests and predictive models. *Circulation*, 115, 654–657.
- Zouros, N. C. (2010). Lesvos Petrified Forest geopark, Greece: geoconservation, Geotourism and Local development. *The George Wright Forum*, 27, 19–28.



Yuyan Luo is an associate professor at College of Management Science, Chengdu University of Technology, Chengdu, China. She received the Ph.D. in management science and engineering from Business School of Sichuan University, Chengdu, China, in 2013. Her current research interests focus on tourism management, system simulation and data mining. She has published more than 20 papers in journals including *Knowledge-Based System*, *Journal of Cleaner Production*, *Discrete & Continuous Dynamical Systems - S*, *Chaos Solitons & Fractals*, *Journal of Intelligent & Fuzzy Systems*, *Symmetry*, etc



Jinjie He is a postgraduate student at College of Management Science, Chengdu University of Technology, Chengdu, China. Her research interests are natural language processing, machine learning and sentiment analysis.



Yu Mou is a postgraduate student at College of Management Science, Chengdu University of Technology, Chengdu, China. His research interests include system dynamics and behavior research on personalized recommendation.



Jun Wang* is an associate professor at Business School, Sichuan Normal University, Chengdu, China. He received the Ph.D. in management science and engineering from Business School of Sichuan University, Chengdu, China, in 2014. His current research areas include applied predictive modeling, complex scenic system modeling analysis, and uncertain multi-attribute decision theory and methods. He has published more than 15 papers in journals including *Knowledge-Based System*, *European Journal of Operational Research*, *Chaos Solitons & Fractals*, etc



Tao Liu is an undergraduate student at College of Management Science, Chengdu University of Technology, Chengdu, China. His research interests include data visualization and system dynamics.