

Research Article

Question Text Classification Method of Tourism Based on Deep Learning Model

Wanli Luo ¹ and Lei Zhang ²

¹College of Information and Engineering, Sichuan Tourism University, Chengdu, Sichuan 610000, China

²Personal Business Department, Sichuan Rural Credit, Chengdu, Sichuan 610000, China

Correspondence should be addressed to Wanli Luo; luowanli@sctu.edu.cn

Received 21 October 2021; Revised 12 November 2021; Accepted 14 December 2021; Published 5 January 2022

Academic Editor: Mohammad R Khosravi

Copyright © 2022 Wanli Luo and Lei Zhang. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The Internet of Things applications are diverse in nature, and a key aspect of it is multimedia sensors and devices. These IoT multimedia devices form the Internet of Multimedia Things (IoMT). Compared with the Internet of Things, it generates a large amount of text data with different characteristics and requirements. Aiming at the problems that machine learning and single structure deep learning model cannot effectively grasp the text emotional information in text processing, resulting in poor classification effect, this paper proposes a text classification method of tourism questions based on deep learning model. First, the corpus is trained with word2vec tool based on continuous word bag model to obtain the text word vector representation. Then, the attention mechanism is introduced into the long-short term network (LSTM), and the attention-based LSTM model is constructed for text feature extraction, which highlights the impact of different words in the input text on the text emotion category. Finally, the text features are input into the Softmax classifier to obtain the probability distribution of text categories, and the model is trained combined with the cross entropy loss function. The experimental results show that the average accuracy, recall, and F value are 0.943, 0.867, and 0.903, respectively, which has better classification effect than other methods.

1. Introduction

As one of the main ways of leisure and entertainment after the continuous improvement of China's social economy and people's material living standards, tourism has attracted more and more attention and favor, and "self-help tourism" has become the mainstream of tourism forms. In the process of self-help travel, problems such as route planning, catering, and accommodation and itinerary strategy are easy to occur. With the rapid development of the Internet, tourists mainly obtain tourism information through network query and Q & A. Access to information includes tourism information released by major tourism portals, tourism applications, and other media platforms. This kind of tourism information has the characteristics of popularization and generalization [1]. However, when obtaining tourism information, it is necessary to publish the questions and wait for the reply of other users, which has a delay. Moreover, the tourism Q & A community

usually classifies the questions according to the geographical location, which cannot fully cover all kinds of questions. In addition, the traditional tourism Q & A community generally uses manual annotation or machine learning model for problem classification, resulting in low classification efficiency and accuracy, unable to quickly and accurately locate the problem category of tourists, which affects the subsequent information retrieval [2–4]. Therefore, how to automatically classify all kinds of tourism questions quickly and efficiently has become an urgent problem to be solved.

The syntactic and semantic information of tourism question text mainly depends on the text composition and sequence order. On the one hand, the grammar of tourism questions consists of multiple question keywords and some network popular words. The words in the text sequence are modeled to form the low-level subspace structure information of the text sequence [5]. On the other hand, the semantic information and syntactic information of tourism

question text come from the text sequence itself. Compared with traditional machine learning technology, the existing deep learning technology can better capture the deep semantic information of the text and solve the error problem caused by manual design features, and the classification accuracy is higher [6]. However, most text classification methods are deep learning models based on a single structure or simply concatenate multiple models. When mining the deep features of text, a large amount of syntax and syntactic information will be lost and redundant information will be added [7]. Therefore, this paper proposes a text classification method of tourism questions using deep learning model. Its innovations are summarized as follows:

- (1) In order to overcome the problems of gradient explosion or disappearance problem of recurrent neural network (RNN), the proposed method uses the long-short term memory (LSTM) network to construct the text classification model and inputs the text word vector into the model to complete the feature extraction, so as to ensure the accuracy of tourism question text classification
- (2) In order to obtain the different influence weights of different words in emotion classification, a text emotion classification model based on LSTM with attention mechanism is proposed, which focuses on the emotional information of text data and further improves the expression ability of text features

The rest of this article is organized as follows. The second section introduces the relevant research progress in this field; the third section specifically introduces the proposed LSTM text classification model based on the attention mechanism; the fourth section compares with the current text classification model to realize the feasibility of the method proposed in this article and the optimality experiment simulation analysis; Section 5 is the conclusion of this paper.

2. Related Works

At present, there are many researches on the text classification methods of tourism questions at home and abroad. In addition to the early manual annotation methods, the traditional machine learning method is the main method of tourism question text classification in recent years. Early question text classification methods mainly used simple machine learning model to classify and recognize different types of question text. Ref. [8] proposed a text and document classification model of support vector machines (SVM). Different experimental results show that it has high classification accuracy on any kind of data set, but the classification efficiency needs to be improved. Ref. [9] proposed an active learning text question and answer classification method, which can potentially reduce the size of the training data set, but the prediction of model performance in active learning may be affected by statistical deviation, so there is still room to further improve the accuracy of text classification. Ref. [10] proposed a cost sensitive analysis

air valve, which is derived by differential evolution algorithm. The experimental results show that the algorithm has high classification accuracy, but the classification accuracy and classification efficiency of complex texts need to be improved.

Deep learning technology has developed rapidly in recent years and has been applied to question text classification tasks and achieved good results. Compared with traditional machine learning technology, it can capture the deep semantic information of text and solve the error problem caused by manual design features and has higher classification accuracy. Ref. [11] uses the classification algorithm of LSTM and convolutional neural network to improve the classification accuracy of problem data sets by changing the vector size and embedding type of combined architecture, but the data sets required for training are large and the training time is long. Ref. [12] evaluated the performance of shallow machine learning and deep learning in text classifiers and text classification embedded in small clinical data sets. Self-training and pretraining word embedding were used as input representation schemes to evaluate logistic regression and long-short term training methods. In the balanced data supported by pretraining embedding, the accuracy of deep learning method was better. Ref. [13] compares the text data classification algorithms of deep learning and traditional machine learning. The results show that the deep learning algorithm has better classification accuracy in some specific cases, but it needs more training data and training time to improve the accuracy. Ref. [14] proposed a unified learning framework of hierarchical cognitive structure learning model, which includes two submodules: attention ordered cyclic neural network and hierarchical two-way capsule. It has good text classification performance, but the simple series structure of the two models is difficult to mine deep-seated text features. Aiming at the shortcomings of the above methods, a text classification method of tourism questions based on deep learning model is proposed, and the attention mechanism is introduced into LSTM network to construct a high-performance text classification model.

3. Proposed Research Methods

3.1. Text Preprocessing. The text of tourism questions is different from the formal and standardized text published by traditional media. The text of tourism questions is usually very short and no more than 130 words at most, including punctuation, slang, abbreviations of specific terms, user nicknames, and other contents. These contents have brought great noise interference to the text emotion classification.

In order to remove unnecessary noise interference, the related technologies in natural language processing are used to preprocess the text. First, this paper uses Jieba word segmentation tool to segment each comment text; then, based on the stop words list provided by Baidu, the stop words are removed, and then the noise is removed. When removing noise, it mainly deals with slang, abbreviations of specific terms, user nicknames, punctuation, and other strings involved.

3.2. Attention-Based LSTM Text Classification Model. The research goal is to solve the problem of text classification, which is mainly divided into three parts: text data representation, text feature extraction, and text classifier. The structure of the attention-Based LSTM text classification model is shown in Figure 1, which is mainly composed of word vector representation part, feature extraction part, and classifier part.

Through the analysis of common text data representation technology, it is decided to use word embedding technology to complete the representation of text data. The word vector is obtained through the word embedding language model. In the feature extraction part, according to the characteristics of text classification corpus, this paper uses the attention-based LSTM model as the feature extraction model. The model uses the LSTM model as the coding model and adds the attention model mechanism to calculate the attention probability, i.e., influence weight, of the text sequence for the overall semantic information, and optimizes the feature vector. In the text classifier part, the logistic regression method is used as the classifier. The logistic regression classifier is simple and efficient and can be easily combined with the feature extraction model.

3.3. Text Word Vector Representation. The corpus is trained with Google's open source tool word2vec model to obtain the vector representation of text words. Word vector can capture the complex mapping from words in corpus to real dimensional vector space. Specify word vector space as φ , its size is $|\varphi| \times m$, each line in φ represents m dimensional word vector of a word, and $|\varphi|$ represents the number of words contained in word vector. A comment text T in the corpus can be expressed as the following sequence:

$$(c_1, c_2, \dots, c_n), \quad (1)$$

where n represents the number of words in text T ; c_i stands for the i ($1 \leq i \leq n$) word in T . If T is converted into a word vector matrix, first search the word vector corresponding to word c_i in φ . If it exists, select the corresponding word vector and represent it with C_i , otherwise, set the corresponding word vector $C_i = 0$. After finding the word vector corresponding to each word, stack each word vector to form a word vector characteristic matrix C , whose size is $n \times m$. Each line of C represents the word vector corresponding to a word in the corpus, which can be expressed as

$$(c_1, c_2, \dots, c_n) \Rightarrow (C_1, C_2, \dots, C_n)^T. \quad (2)$$

3.4. Feature Extraction. Text emotion classification is mainly based on the key emotional words expressing views, feelings, and attitudes in the text to judge the text emotion tendency, among which the words with strong emotional color play a key role in judging the text emotion tendency [15, 16]. In order to fully reflect the role of emotional keywords in the process of text emotion classification, this paper proposes a text emotion classification model based on LSTM with attention mechanism. The model adds attention mechanism on the LSTM based network, which mainly distributes the

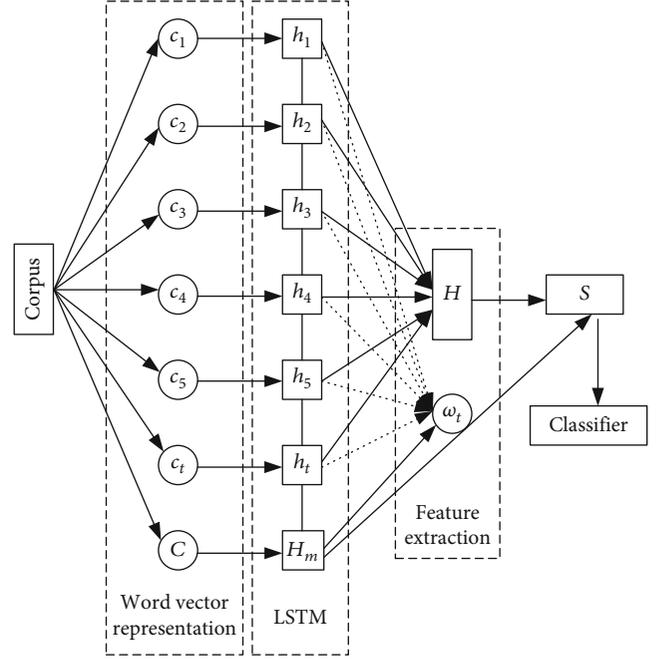


FIGURE 1: Structure of attention-based LSTM text classification model.

weight of emotional information of words and highlights the impact of different words in the input text on the emotional category of the text [17, 18].

3.4.1. LSTM Network Structure. In this paper, LSTM neural network structure is used as the core component of tourism question text emotion classification model. LSTM neural network structure not only has the advantages of traditional recurrent neural network (RNN), overcomes the problem of RNN gradient explosion or disappearance, but also can effectively process sequence data of arbitrary length and capture long-term dependence of data [19, 20]. The LSTM network structure is shown in Figure 2.

Taking the multifeature representation of words with emotion vector as the input of LSTM, the hidden layer state value corresponding to the input is obtained. The specific calculation of LSTM neural network memory cell is as follows:

$$\begin{aligned} f_t &= \delta(\omega_f[h_{t-1}, \mathbf{x}_{ct}] + \mathbf{b}_f) + \delta(\omega_f[h_{t-1}, \mathbf{x}_{qt}] + \mathbf{b}_f), \\ i_t &= \tanh(\omega_i[h_{t-1}, \mathbf{x}_{ct}] + \mathbf{b}_i) + \tanh(\omega_i[h_{t-1}, \mathbf{x}_{qt}] + \mathbf{b}_i), \\ C_t &= f_t \cdot C_{t-1} + i_t * \tilde{C}_t, \\ o_t &= \delta(\omega_o[h_{t-1}, \mathbf{x}_{ct}] + \mathbf{b}_o) + \delta(\omega_o[h_{t-1}, \mathbf{x}_{qt}] + \mathbf{b}_o), \\ h_t &= o_t * \tanh(C_t), \end{aligned} \quad (3)$$

where $\mathbf{x}_t = [\mathbf{x}_{ct}, \mathbf{x}_{qt}]$ represents the input at time t ; \mathbf{x}_{ct} and \mathbf{x}_{qt} represent the word meaning vector and emotion vector, respectively; h_t represents the output at time t ; i_t represents whether some information in the input door needs to be updated; f_t is the output matrix of the

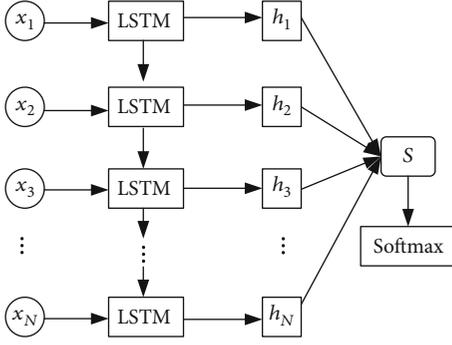


FIGURE 2: LSTM network structure.

forgetting gate; ω is the weight matrix; \mathbf{b} is the offset vector; δ is sigmoid nonlinear activation function.

3.4.2. Attention Mechanism. The emotional classification of text not only needs to consider the context relationship between words but also needs to consider which words are more prominent in the expression of text emotional classification. Words with greater emotional contribution should be given higher weight or attention [21, 22]. Aiming at the problem that the emotional features cannot be effectively highlighted in the process of text emotional classification, and the proposed method constructs an LSTM text classification model based on attention mechanism, which focuses on the emotional information of text data and further improves the expression ability of text features. In this model, the word emotion influence weight is determined based on the correlation between the output h_t of each hidden layer and the context vector \mathbf{s} . The calculation process of attention mechanism is shown in Figure 3.

The calculation of attention mechanism can be realized in two steps:

Step 1. Calculate the attention distribution on all input information, that is, take the context vector \mathbf{s} and the output h_t of the hidden layer as inputs, enter a single-layer perceptron, and obtain the implicit representation u_t of the result through calculation. The calculation formula is as follows:

$$u_t = \tanh(\alpha h_t + \beta \mathbf{s}), \quad (4)$$

where α and β are the weight matrix; h_t is the output of the hidden layer; \mathbf{s} is the query vector. Then, ϑ_t is obtained through softmax operation, which is calculated as follows:

$$\vartheta_t = \text{soft max}(u_t), \quad (5)$$

where the probability vector composed of ϑ_t is the emotional attention distribution of the word.

Step 2. Calculate the weighted sum of the input information according to the attention distribution ϑ_t , that is, the attention distribution ϑ_t represents the correlation between the time t information in the input information vector \mathbf{H} and the query \mathbf{s} when a query ϑ_t is given.

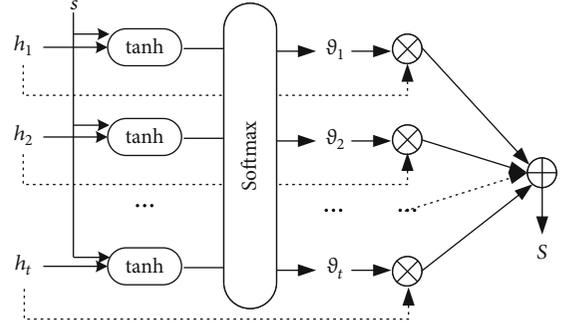


FIGURE 3: Calculation process of attention mechanism.

The input information is summarized by weighted summation to obtain the attention value. The specific calculation is as follows:

$$S = \sum_{t=1}^N \vartheta_t h_t. \quad (6)$$

3.5. Classifier. The text classification model based on attention-based LSTM uses softmax as the output layer for normalization calculation, and combined with the cross entropy loss function, the objective function is expressed as follows:

$$\text{Loss} = - \sum_{i=1}^K Y_i \log(y_i), \quad (7)$$

where K represents the number of texts in the corpus, \mathbf{Y}_i represents the real probability distribution vector of the current text category, \mathbf{y}_i represents the probability distribution vector of the current text predicted by the classification model, and the dimension of the vector is equal to the number of classification labels. By minimizing the objective function, the classification model can be obtained [23, 24].

The model based on attention mechanism generally includes two parts: one is the calculation process of attention probability distribution, and the other is the calculation process based on the final characteristics of attention distribution. In this model, the attention probability of the output data at time t to the final state is calculated as follows:

$$v_t = \frac{\exp(h'_t)}{\sum_{i=1}^N \exp(h'_i)}, h'_t = h_t^T \hat{\omega} F. \quad (8)$$

Softmax function is the calculation method of attention probability distribution, where N represents the number of input sequence elements; $\hat{\omega}$ is the weight matrix; F represents the sum of the final hidden layer state values in each independent direction in the LSTM; h_t represents the sum of the hidden layer state values at time t .

In the proposed model, the final feature F_{final} is obtained based on the attention distribution, and the calculation process is expressed as follows:

TABLE 1: Experimental environment.

Environmental parameters	Configuration
Operating system	Ubuntu 14.04.5
Development language	Python
Development framework	Tensorflow
Memory	256G
CPU	Intel(R) Xeon(R) CPU E5-2620
GPU	NVIDIA corporation GM200

TABLE 2: Statistical results of data sets.

Category	Training set	Test set	Validation set	Average text length
Place	1226	525	105	32
Time	1397	598	120	38
Entity	1329	569	114	63
Figures	1215	521	104	75
Description	1691	725	145	22
Character	143	61	12	2534

TABLE 3: LSTM parameter setting.

Parameter	Value	Parameter	Value
LSTM network layer	1 layer	numClasses	2
Batch_size	128	Dropout	0.7
LSTM_size	256	Loss function	Cross entropy
Learning rate	0.0001	Optimizer	RMSProp optimizer

$$F_{\text{final}} = \sum_{t=1}^N v_t h_t. \quad (9)$$

After obtaining the text feature vector F_{final} based on the attention mechanism, the probability distribution of the classification label is calculated through the Softmax function of the output layer. The calculation process is expressed as follows:

$$y = \text{soft max} \left(F'_{\text{final}} \right) = \frac{\exp \left(F'_{\text{final}(i)} \right)}{\sum_{j=1}^T \exp \left(F'_{\text{final}(j)} \right)}, \quad (10)$$

$$F'_{\text{final}} = \omega_o F_{\text{final}}$$

where D is the number of category labels; ω_o represents the weight matrix of the model output layer; $F'_{\text{final}(i)}$ represents the i component value in vector F'_{final} , and the vector length is equal to the number of classification labels. After the softmax function, the probability distribution y of text category based on the attention mechanism is obtained, and the cross entropy loss is calculated with the real category distribution Y , which is expressed as

$$E(Y, y) = -Y \log(y), \quad (11)$$

TABLE 4: Classification discrimination confusion matrix.

Real results	Prediction results	
	In category A	Not in category A
In category A	r	l
Not in category A	g	z

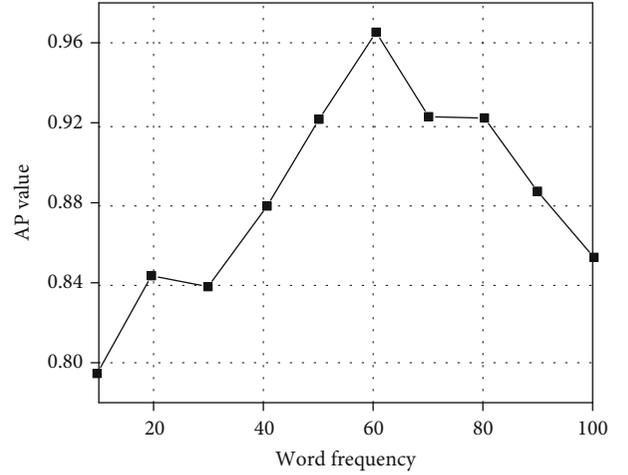


FIGURE 4: AP value change curve based on word frequency.

where Y represents the probability distribution of the real category; y represents the probability distribution of the category predicted by the model.

4. Experiment and Analysis

4.1. Experimental Setup

4.1.1. Hardware Environment. The experiment is implemented based on the deep learning framework TensorFlow, which is a deep learning framework based on graph calculation. It uses the data flow between nodes to transfer data and completes the calculation in the nodes. As an open source framework, Tensorflow integrates several models including convolutional neural network, RNN network, and LSTM model [25]. The emergence of Tensorflow framework makes the use of deep learning model simpler and convenient and reduces the difficulty of applying deep learning model [26]. The specific experimental environment is shown in Table 1.

4.1.2. Experimental Data Set. The tourism text data set is used as the experimental data set, which is a user-defined benchmark data set, mainly from tourism websites such as Ctrip, Tuniu, Ma honeycomb, and Tongcheng, including 6 categories of 10000 sample data such as tourism location, time, and people. Before the experiment, the data set needs to be preprocessed such as selection, cleaning, and stop words to reduce errors. In order to verify the effectiveness of the proposed method, 70% of the samples are randomly selected as the training set, the remaining 30% of the samples are used as the test set, and 20% of the samples in the training set are randomly divided as the cross-validation set (the experiment has conducted 6 cross-validation). The statistical results of the data set are shown in Table 2.

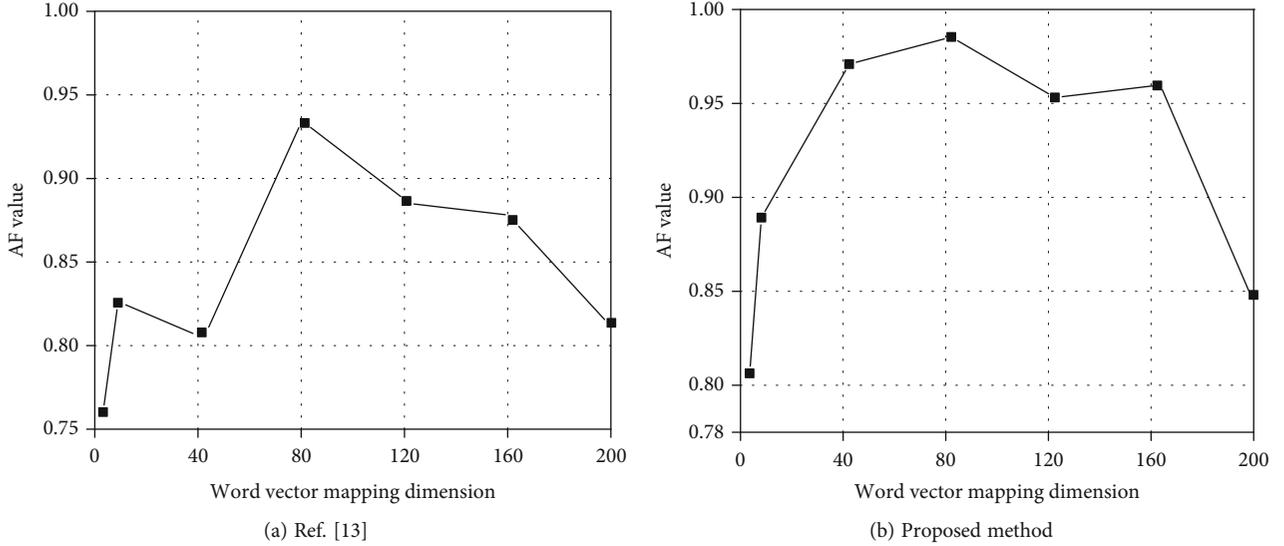


FIGURE 5: F value of text classification method for tourism questions with different word vector mapping dimensions.

4.1.3. LSTM Parameter Setting. Based on LSTM neural network structure, a text emotion classification model based on attention mechanism which can express word tag relationship is constructed. The LSTM neural network adopts one-layer network structure, the number of hidden nodes is 256, and the learning rate is 0.0001. The optimization algorithm adopts RMSPropOptimizer optimizer. The specific parameter settings are shown in Table 3.

4.2. Evaluating Indicator. The accuracy P , recall R , and F values are selected as the evaluation indexes, and the classification discrimination confusion matrix is shown in Table 4.

P represents the proportion of samples of real category among the samples predicted to be a category after emotion classification of the test set, that is

$$P = \frac{r}{r + g}. \quad (12)$$

R represents the proportion of a category predicted as a real category in all real categories in the test set, that is

$$R = \frac{r}{r + l}. \quad (13)$$

In order to comprehensively consider the accuracy P and recall R , the weighted harmonic average F of the two is used to measure the final classification effect, that is

$$F = \frac{2 \times P \times R}{P + R}. \quad (14)$$

The task of text emotion classification is oriented to multi classification. Therefore, after calculating the accuracy P and recall R corresponding to each category, the average accuracy (AP), average recall (AR), and average F (AF) corresponding to the three categories are used as the evaluation indexes to measure the performance of emotion classifier.

4.3. Model Training. When training the classification model, the number of iterations is set to 50, and the relationship between AP value and word frequency is shown in Figure 4.

As can be seen from Figure 4, as the word frequency increases, the AP value also increases gradually until the word frequency reaches the optimal value when the word frequency is 60. The AP value exceeds 0.96, and then the AP value decreases with the increase of word frequency. Therefore, when the word frequency is set to 60 in the training of deep learning model, its classification performance is the best.

4.4. Influence of Word Vector Dimension on Model Performance. The word vector mapping dimension plays an important role in the classification accuracy of the model. Therefore, we change the word vector mapping dimension to verify its impact on the text classification accuracy of tourism questions. At the same time, in order to demonstrate the classification accuracy of the proposed method, it is compared with Ref. [13]. The AF values of the two methods in tourism text dataset under different word vector embedding dimensions are shown in Figure 5.

As can be seen from Figure 5, with the increase of word vector mapping dimension, the AF value of the proposed tourism question text classification method first increases rapidly. When the word vector mapping dimension is greater than 80, the AF value stops increasing and begins to decrease, which indicates that too low word vector mapping dimension cannot better map the text to low-dimensional space, and high-dimensional embedding may lead to too sparse vector representation. Therefore, it cannot effectively improve the classification performance and will consume more training time. However, compared with Ref. [13], the AF value of the proposed method is higher as a whole, and when the word vector dimension is in the range of 40~140, the AF value fluctuates less. This is because it adopts LSTM network, which can better map text.

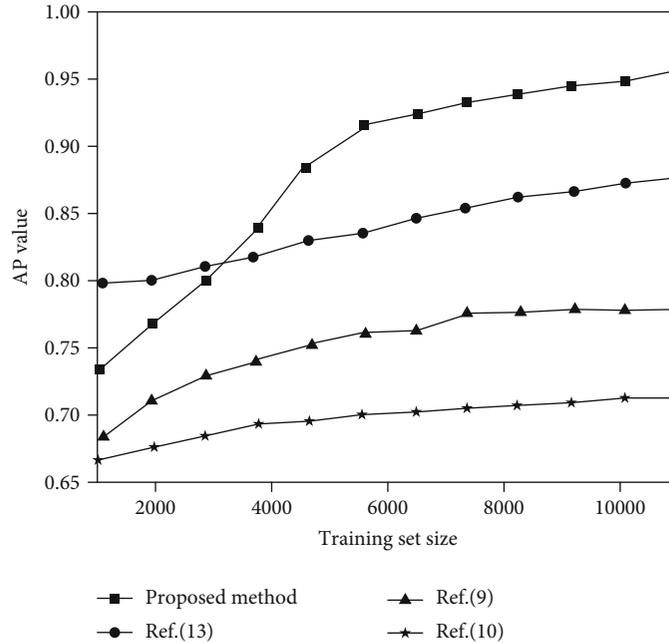


FIGURE 6: Comparison results of AP values of different methods.

TABLE 5: Comparison results of classification performance of different methods.

	Ref. [9]	Ref. [10]	Ref. [13]	Proposed method
AP	0.809	0.783	0.878	0.943
AR	0.752	0.664	0.839	0.867
AF	0.779	0.719	0.858	0.903

4.5. *Comparison of Classification Performance of Different Methods.* According to the size of the training set, the proposed method is compared with the methods in $T = \text{Ref. [9, 10] and [13]}$. The results are shown in Figure 6.

It can be found from Figure 6 that with the increase of training set, the AP value of various methods tends to be stable. Compared with other methods, the AP value of the proposed method is the highest and close to 0.952. The attention-based LSTM text classification model can effectively improve the classification effect of tourism question text and combined with the cross entropy loss function training model to further ensure the classification performance of the model. Ref. [13] uses a single deep learning model for text classification. Due to the lack of emotional consideration, the AP value is about 0.075 lower than the proposed method. Both Ref. [9] and Ref. [10] adopt traditional methods, so the overall text classification performance is poor when dealing with complex training sets, and the AP value is lower than 0.80.

In addition, the specific data of AP, AR, and AF obtained from the experiments by the four methods are shown in Table 5.

It can be seen from Table 5 that the overall classification performance of the proposed method is the best, and the values of AP, AR, and AF were 0.943, 0.867, and 0.903, respectively. The proposed method uses the LSTM network

to extract the depth feature vector, reduces the dimension of the output feature vector, introduces the attention mechanism to highlight the emotional role, significantly improves the classification performance of question text, and proves that it is robust to the emotional classification of tourism text. Ref. [9] uses the active learning model for text classification. The AF value of traditional machine learning is only 0.779, which is 0.124 lower than that of the proposed method. Ref. [10] uses differential evolution algorithm to realize text classification, but the algorithm does not fully analyze the characteristics of tourism text, and the algorithm performance is poor, so the AP value is only 0.783. Ref. [13] classifies text based on a single deep learning model. However, this method lacks emotional consideration, so its AF value is 0.858, and the whole performance needs to be further improved. It can be demonstrated that the proposed method has good text classification ability of tourism questions.

5. Conclusion

Under the background that self-help travel has become the mainstream form of tourism, tourists can obtain information through Q & A from the Internet platform, but there are problems of delay and inaccurate classification in self-help Q & A. Therefore, a text classification method of tourism questions based on deep learning model is proposed. The text word vector obtained based on the continuous word bag model is input into the attention-based LSTM model for feature extraction, and the probability distribution of text category is obtained by Softmax classifier. The proposed method is experimentally analyzed using the tourism text data set, and the results show that the LSTM model can effectively capture the relationship between word

vectors. When the word frequency is set to 60 and the word vector dimension is 80, the AP value of the model exceeds 0.96. The introduction of attention mechanism can better highlight the role of emotion and improve the accuracy of text classification of tourism questions. The AP, AR, and AF were 0.943, 0.867, and 0.903, respectively, which were better than other comparison methods. However, the proposed method uses Softmax function for task calculation. In the next research, some acceleration methods, such as hierarchical Softmax and negative sampling technology, can be considered to improve the overall performance of the classification model.

Data Availability

The data included in this paper are available without any restriction.

Conflicts of Interest

The authors declare that there is no conflict of interest regarding the publication of this paper.

Acknowledgments

We wish to express their appreciation to the reviewers for their helpful suggestions which greatly improved the presentation of this paper. This research was jointly funded by the Sichuan Science and Technology Program, Grant/Award numbers: 2019ZYZF0169; and the A Ba Achievements Transformation Program, Grant/Award numbers: 19CGZH0006 and 21CGZH0002.

References

- [1] C. D. Cottrill, "MaaS surveillance: privacy considerations in mobility as a service," *Transportation Research Part A: Policy and Practice*, vol. 131, no. 8, pp. 50–57, 2020.
- [2] W. Wang and A. Feng, "Self-information loss compensation learning for machine-generated text detection," *Mathematical Problems in Engineering*, vol. 2021, Article ID 6669468, 7 pages, 2021.
- [3] A. Mignan, "A preliminary text classification of the precursory accelerating seismicity corpus: inference on some theoretical trends in earthquake predictability research from 1988 to 2018," *Journal of Seismology*, vol. 23, no. 4, pp. 771–785, 2019.
- [4] F. Zhao, Y. Li, L. Bai, Z. Tian, and X. Wang, "Semi-supervised multi-granularity CNNs for text classification: an application in human-car interaction," *IEEE Access*, vol. 8, no. 99, pp. 68000–68012, 2020.
- [5] S. G. Burdisso, M. Errecalde, and M. Montes-y-Gómez, "A text classification framework for simple and effective early depression detection over social media streams," *Expert Systems with Applications*, vol. 133, no. 11, pp. 182–197, 2019.
- [6] M. Sokolowska, M. Mazurek, M. Majer, and M. Podpora, "Classification of user attitudes in twitter -beginners guide to selected machine learning libraries," *IFAC-PapersOnLine*, vol. 52, no. 27, pp. 394–399, 2019.
- [7] M. M. Mironczuk, J. Protasiewicz, and W. Pedrycz, "Empirical evaluation of feature projection algorithms for multi-view text classification," *Expert Systems with Applications*, vol. 130, no. 4, pp. 97–112, 2019.
- [8] X. Luo, "Efficient english text classification using selected machine learning techniques," *Alexandria Engineering Journal*, vol. 60, no. 3, pp. 3401–3409, 2021.
- [9] A. Varghese, T. Hong, C. Hunter, G. Agyeman-Badu, and M. Cawley, "Active learning in automated text classification: a case study exploring bias in predicted model performance metrics," *The Environmentalist*, vol. 39, no. 3, pp. 269–280, 2019.
- [10] C. Padurariu M. E. Breaban et al., "Dealing with data imbalance in text classification," *Procedia Computer Science*, vol. 159, pp. 736–745, 2019.
- [11] S. Yilmaz and S. Toklu, "A deep learning analysis on question classification task using Word2vec representations," *Neural Computing and Applications*, vol. 32, no. 7, pp. 2909–2928, 2020.
- [12] M. Oleynik, A. Kugic, Z. Kasáč, and M. Kreuzthaler, "Evaluating shallow and deep learning strategies for the 2018 n2c2 shared task on clinical text classification," *Journal of the American Medical Informatics Association*, vol. 26, no. 11, pp. 1247–1254, 2019.
- [13] A. Varghese, G. Agyeman-Badu, and M. Cawley, "Deep learning in automated text classification: a case study using toxicological abstracts," *Environment Systems and Decisions*, vol. 40, no. 4, pp. 465–479, 2020.
- [14] B. Wang, X. Hu, P. Li, and P. S. Yu, "Cognitive structure learning model for hierarchical multi-label text classification," *Knowledge-Based Systems*, vol. 218, no. 3, pp. 106876–106887, 2021.
- [15] D. Petschke and T. Staab, "A supervised machine learning approach using naive Gaussian Bayes classification for shape-sensitive detector pulse discrimination in positron annihilation lifetime spectroscopy (PALS)," *Section A, Accelerators, Spectrometers, Detectors and Associated Equipment*, vol. 947, no. 12, pp. 162742–162742.9, 2019.
- [16] B. Zhong, X. Xing, P. Love et al., "Convolutional neural network: deep learning-based classification of building quality problems," *Advanced Engineering Informatics*, vol. 40, no. 7, pp. 46–57, 2019.
- [17] Z. Chen and J. Ren, "Multi-label text classification with latent word-wise label information," *Applied Intelligence*, vol. 51, no. 2, pp. 966–979, 2021.
- [18] Y. Zhu, W. Zheng, and H. Tang, "Interactive dual attention network for text sentiment classification," *Computational Intelligence and Neuroscience*, vol. 2020, Article ID 8858717, 11 pages, 2020.
- [19] K. Purwandari, J. W. C. Sigalingging, T. W. Cenggoro, and B. Pardamean, "Multi-class weather forecasting from twitter using machine learning approaches," *Procedia Computer Science*, vol. 179, no. 4, pp. 47–54, 2021.
- [20] A. Mohasseb, M. Bader-El-Den, and M. Cocea, "A customised grammar framework for query classification," *Expert Systems with Applications*, vol. 135, no. 11, pp. 164–180, 2019.
- [21] B. Stasak, J. Epps, and R. Goecke, "Automatic depression classification based on affective read sentences: opportunities for text-dependent analysis," *Speech Communication*, vol. 115, no. 6, pp. 1–14, 2019.
- [22] B. André Sumithra et al., "Text classification to inform suicide risk assessment in electronic health records," *Studies in Health Technology and Informatics*, vol. 264, no. 3, pp. 40–44, 2019.

- [23] T. Henry, D. Banks, D. Owens-Oas, and C. Chai, "Modeling community structure and topics in dynamic text networks," *Journal of Classification*, vol. 36, no. 2, pp. 322–349, 2019.
- [24] Y. Baghdadi, A. Bourrée, A. Robert et al., "Automatic classification of free-text medical causes from death certificates for reactive mortality surveillance in France," *International Journal of Medical Informatics*, vol. 131, no. 11, article 103915, 2019.
- [25] M. Hashemi, "Web page classification: a survey of perspectives, gaps, and future directions," *Multimedia Tools and Applications*, vol. 79, no. 17-18, pp. 11921–11945, 2020.
- [26] F. Beretta, Á. L. Rodrigues, R. Peroni, and J. F. C. L. Costa, "Using UAV for automatic lithological classification of open pit mining front," *REM-International Engineering Journal*, vol. 72, no. 1, Supplement 1, pp. 17–23, 2019.