



Extraction Method of Emotional Elements of Online Learning Text Information Based on Natural Language Processing Technology

Haolin Song¹(✉), Dawei Song¹, and Yankun Zhen²

¹ School of Computer Science and Technology, Beijing Institute of Technology, Beijing 100081, China

songhaolin_123@126.com

² College of Science, Xi'an Shiyou University, Xi'an 710065, China

Abstract. The current methods of extracting emotional elements of text information generally adopt the principle of template matching, the algorithm is complex. Due to the limitations of the selected template, the network learning text information emotion elements cannot be comprehensively extracted, so the extraction accuracy and efficiency are low. In order to solve the above problems, this paper studies the emotional element extraction method of online learning text information based on natural language processing technology. Preprocess the online learning text information and find new words; Split the preprocessed text into sentences to generate transaction items; Frequent noun items are mined by association rules, irrelevant nouns are filtered by filtering algorithm, and the emotional elements of text information are extracted; Using the credibility analysis algorithm to judge the emotional polarity of text, and using the RNN neural network algorithm in natural language processing technology, the emotional elements of online learning text information are extracted. The test data show that the extraction time of the proposed feature extraction method is reduced by at least 35%, and the extraction accuracy of the method is improved to 80%, and the extraction result is more reliable.

Keywords: Natural language processing · Online learning · Text information · Emotional elements · Emotion extraction · Neural network

1 Introduction

Online learning is an inevitable product of the information age. It breaks the time and space constraints of traditional forms of education. It is a kind of education anytime and anywhere. The education industry has a considerable demand for products in the field of artificial intelligence. Especially with the development of online education industry and people's recognition of relevant concepts, it can be predicted that in the near future, a large part of education work will be completed by computers, such as automatic marking, intelligent error correction, human like teaching, etc.

The accelerated landing of artificial intelligence industry in the education industry can not only reduce the workload of teachers, but also generate targeted teaching tasks for different types of students, which is of great significance to the improvement of personalized education system. Natural language processing is the intersection of computer science and Linguistics and an important branch in the field of artificial intelligence.

In recent years, with the improvement of computer computing ability, natural language processing technology has played a more and more important role to help solve various problems encountered in people's actual production and life. Deliver high-quality teaching resources through the Internet, realize students' personalized learning and promote students' all-round and free development. At present, for online learning, in order to meet the needs of students for timely and continuous feedback, teachers need to invest a lot of energy and uninterrupted online guidance.

However, the working hours of teachers' feedback are limited and their role is limited. Therefore, it is necessary to promote learning with the effective integration of technical support and feedback methods. It is urgent to explore the scientific feedback strategy under technical support [1]. The text of students' online learning contains a lot of information. Learning these information can analyze learners' evaluation of online learning and other effect feedback. Emotional elements are generally composed of three parts: emotional source, receptor and emotional tendency. The emotional elements in the information play an important role in judging whether the feedback is positive or not. Therefore, extracting emotional elements from text information is very necessary to improve the quality of online teaching.

Literature [2] used the domain dictionary to analyze the constituent elements of scenic spots, which provided a more accurate method for garden management optimization. Literature [3] uses Gru neural network to fuse multiple features to realize text emotion classification, which has better generalization ability than traditional machine learning. Literature [4] designed an emotion analysis model integrating attention mechanism and bigru. Experiments on public data sets show that the model can improve the performance of emotion analysis.

When extracting emotional elements from online learning text information, the above emotional element analysis method needs a priori template matching. The processing process is complex and cumbersome, which reduces the extraction efficiency, and the accuracy of the extraction method is difficult to meet the actual needs.

Therefore, this paper studies the emotional element extraction method of online learning text information based on natural language processing technology. Preprocess online learning text information and find new words; Split the preprocessed text into sentences to generate transaction items; Mining frequent noun items by association rule algorithm; The innovation of this paper is to use the filtering algorithm to filter irrelevant nouns and extract the emotional elements of text information; The credibility analysis method is used to judge the emotional polarity of text, and the RNN neural network in natural language processing technology is used to extract the emotional elements of online learning text information. The experimental results show that the proposed method shortens the extraction time, improves the extraction accuracy, and the actual extraction results are more reliable.

2 Research on Emotional Element Extraction Method of Online Learning Text Information Based on Natural Language Processing Technology

2.1 Online Learning Text Information Preprocessing and New Word Discovery

When dealing with Chinese corpus, especially internet corpus, there are too many colloquial expressions. In addition, most of the existing word segmentation tools are trained based on the classic people's daily corpus, so it is difficult to identify new words in internet corpus. Therefore, by strengthening the performance of word segmentation, we can pave the way for improving the next complex natural language processing tasks. In this paper, a method based on generalized suffix tree is used to extract high-frequency strings [5].

In the field of natural language processing, text features include the grammar of the language itself, the semantics combined with the context, the position in the text, and the meaning of the overall expression.

TF-IDF is a method based on word frequency statistics, which is used to evaluate the influence of a word on the text [6].

The TF-IDF values are divided into two parts: TF values and IDF values, defined as the following:

$$\begin{cases} TF = \frac{time_{L_i}}{N} \\ IDF = \frac{N_w}{W_{L_i} + 1} \end{cases} \quad (1)$$

In formula (1), $time_{L_i}$ is the number of times the word L_i appears in the online learning text information; N is the total number of words in the online learning text information; N_w is the total number of text information in the split corpus; W_{L_i} is the total number of text containing the word L_i .

TF value reflects the importance of a word, IDF value is used to filter some stop words with high probability, and the final TF-IDF value is obtained by multiplying the two.

By calculating the likelihood ratio between the words of each constituent word, determine whether the assumption of independence between these words is tenable, so as to exclude the string formed by some useless high-frequency words [7]:

$$\lg L(\xi) = \max[\lg P(D_2) - \lg P(D_1)] \quad (2)$$

In formula (2), ξ is a string, it is a word string composed of multiple single words in $\xi_1 \xi_2$; While D_1 is a formalization of the independence hypothesis, representing the meaning of $p(D_2|D_1) \neq p(D_1|D_2)$, which represents the situation where D_2 and D_1 are independent. Similarly, the case of D_2 represents a formalization of the non-independence assumption.

Use the formula to express it as follows:

$$p(D_1|D_2) \neq p(D_2|D_1) \quad (3)$$

In information theory, mutual information has two characteristics: nonnegativity and symmetry. It is a measure of common information between two random variables. When two variables are independent, their mutual information is 0. When there is a dependency between them, their mutual information is related to their dependency and the moisture content of the variable itself.

In the task of new word discovery, this measure can effectively measure the correlation program between characters in the string. If the correlation between words is high, it may be a new word, otherwise, it is less likely. Define mutual information as the following form:

$$MH = \log_2 \frac{p(x)}{\prod_{i=1}^n p(l_i)} \quad (4)$$

In formula (4), $p(x)$ is the probability that the string $x = l_1, l_2, \dots, l_n$ appears in the corpus, $\prod_{i=1}^n p(l_i)$ represents the probability product of characters within a string.

Glove is selected as the benchmark of text representation. Glove combines the advantages of global statistical information of matrix decomposition and local context window, makes frequency statistics of all words as global a priori statistical information, combines it with the context information obtained by sliding window, and allocates relative weight to train word vector.

The string information is obtained through preprocessing, the string information is constructed into a generalized suffix tree, and the information is merged to obtain the sub strings that appear repeatedly in the corpus. In these sub strings, the information of each sub string is found as the standard to judge whether it is a new word, so as to obtain a new word. After the online learning text information is processed according to the above process, the features in the information are extracted.

2.2 Text Information Feature Extraction

The text after lexical analysis in the above preprocessing is divided into sentences and stored in branches to generate transaction items as the input of the association rule mining system.

This paper uses the method of association rules to mine frequent noun items from the online learning text information database, that is, those nouns and noun patterns that appear more often. These nouns and their patterns may appear as noun phrases or product features in the comment text. The specific process is as follows: firstly, the text in the comment database is divided into sentences of different lengths by punctuation for part of speech tagging; the reserved part of speech labels the Class Nouns in the sentence, generates several transaction items, and forms a transaction data set. Then, the Apriori algorithm is used to find the frequent noun item set in the transaction data set. Feature filtering algorithm is used to eliminate inaccurate single nouns and noun patterns that can not form noun phrases.

The basic model of the association rules can be described as follows:

Let $X = \{x_1, x_2, \dots, x_i\}$ be a collection of all projects. S is the transaction database. Each transaction has a unique transaction identity TID, and the transaction W is a subset of items.

Let A be a collection of items, called an item set. The transaction W contains a set of items A . If the item set A contains m items, it is called the m item set. The number of occurrences of itemset A in transaction database S as a percentage of the total transactions in S is called itemset support. If the support of an itemset exceeds the minimum support threshold given by the user, the itemset is called a frequent itemset.

Apriori algorithm divides the process of discovering association rules into two steps. In the first step, all frequent itemsets in the transaction database are retrieved iteratively, that is, itemsets with support not lower than the threshold set by the user; The second step uses frequent itemsets to construct rules that meet the minimum trust of users. In the task of feature extraction, only the first step is needed, and the found frequent itemsets are used as feature set candidates.

There are redundant and inaccurate noun items in the feature candidate set obtained by association rule mining. This paper uses the filtering algorithm to filter the irrelevant nouns by calculating the domain correlation of nouns, so as to obtain the corresponding final feature set.

The probability $P_f = -2 \ln(L)$ of being calculated for each candidate feature. If $y_2 < y_1$, order $-2 \ln(L) = -2 \times Ly$, otherwise 0. The calculation formula of Ly , y_1 , y_2 , y is shown below.

$$\left\{ \begin{array}{l} Ly = (X_{11} + X_{21}) \lg y + (X_{12} + X_{22}) \lg(1 - y) \\ \quad - X_{11} \lg y_1 - X_{21} \lg y_2 - X_{12} \lg(1 - y_1) - X_{22} \lg(1 - y_2) \\ y_1 = \frac{X_{11}}{X_{11} + X_{12}} \\ y_2 = \frac{X_{21}}{X_{21} + X_{22}} \\ y = \frac{X_{11} + X_{21}}{X_{11} + X_{12} + X_{21} + X_{22}} \end{array} \right. \quad (5)$$

In the above formula, X_{11} and X_{12} are the number of documents containing feature words in W_{F1} and W_{F2} , respectively; W_{F1} is a collection of documents that contain a certain topic; W_{F2} is collection of documents without a particular topic. X_{21} and X_{22} represent the number of documents without feature words in W_{F1} and W_{F2} . As shown in Table 1 below:

Table 1. Type markers for learning text information online

	Theme	Non-theme
Characteristic	X_{11}	X_{12}
Non-characteristic	X_{21}	X_{22}

Tight filtering is to eliminate noun patterns with low tightness, because such patterns are difficult to form noun phrases in the text, and are likely not a feature. Therefore, if the interval between two nouns is not more than two words, we think that these two words may form noun phrases and can be used as candidates for characteristics.

After filtering the independent features in the feature set by the information entropy, the emotional polarity classification of the text information was performed.

2.3 Classification of Emotional Polarity of Text Information

Emotional information classification is an important basic research task of text emotion analysis, mainly including text subjective and objective classification and emotional polarity classification (also known as praise and derogatory classification). Subjective and objective classification of text is the premise of the study of the emotion analysis of text, through subjective and objective classification can identify the text containing tendentious views, which are known as subjective text. The classification of emotional polarity can be divided into word level, sentence level and chapter level. This paper uses credibility analysis technology to classify of text information.

The main idea of the classifier fusion strategy based on credibility analysis is to take the specified single classifier with a high speed as the main classifier, and set a certain number of auxiliary classifiers with high accuracy to vote on the category within a certain range.

Assuming that a credibility function $c(x)$ is determined according to the classification principle of the main classifier and the discrimination threshold γ ($\gamma > 0$) is set, when the main classifier classifies the samples, if $c(x) > \gamma$, it is considered that the discrimination result of the main classifier has high credibility, and the result can be used as the final classification result of the samples. The discrimination procedure is shown in Fig. 1.

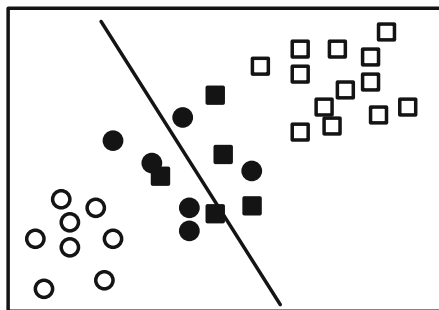


Fig. 1. Schematic representation of the binary emotional polarity discrimination

When hollow samples are easily distinguishable samples, the classification results of the main classifier can be taken directly taken as the output result of the combined classifier, and solid samples are not easily distinguishable samples, requiring the main and auxiliary classifier vote, and then the vote as the output result of the combined classifier.

The discriminative process for the combined classifier is shown in Fig. 2.

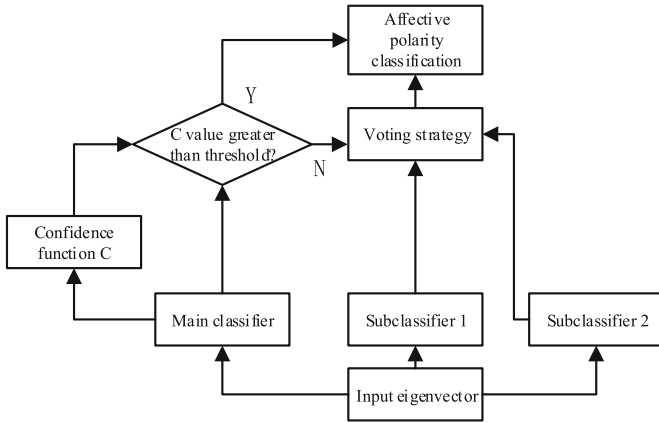


Fig. 2. Discriminative process of text emotional polarity

The confidence function was defined based on the distance of the test sample and the class center in the main classifier:

$$R(x) = \frac{d_1}{d_2} \tag{6}$$

In the formula (6), d_1 is the distance between the clustering center and the text vector for the main emotional polarity feature; d_2 is cluster the distance between the center and the text vector for secondary affective polarity features. When the cosine similarity is used to calculate the distance between the sample and the class center, if $d_1 > d_2$, indicates that the test sample and the nearest class center are closer and far from the subnearest class center, which is considered to be easily distinguishable, the hollow sample in Fig. 1; Assuming that the emotional polarity feature τ_i appears the number of word frequency in the emotional polarity category CA_j is F_{ij} and the text frequency is f_{ij} , the intra-class distribution is:

$$S(\tau_{ij}) = \frac{N \tau_{ij} \log(f \tau_{ij} + 1)}{\sqrt{\sum_{i=1}^n [N \tau_{ij} \log(f \tau_{ij} + 1)]^2}} \tag{7}$$

In the formula (7), $N \tau_{ij}$ is the ratio of the number of words to the number of times the feature appears in all emotion categories.

When the difference between the components of the distribution vector between classes is large, the feature has a strong ability to distinguish the categories; When the components of inter class distribution are similar or the same, the ability of features to distinguish categories will be very low. Set the resolution threshold for emotion polarity classification.

Combined with the problem of emotional polarity classification, in addition to the feature selection based on Category attribute analysis method, we should also consider the emotional features related to polarity classification. Affective word feature is an important feature of polarity classification problem, so we add affective word feature to the feature set, and comprehensively consider the feature selection process of affective polarity classification based on Category attribute analysis method, which is described as follows:

Input: Training set T

Output: The feature set, the set F

1. Obtain the evaluation words as candidate features according to the evaluation word dictionary, and input T;
2. Segment the documents in the training set T, and then count the word frequency information;
3. The feature selection method based on Category attribute analysis is adopted to sort the words in the order of feature weight from large to small;
4. The evaluation words with large weight are selected from the candidate features of polar words and added to the feature set F;
5. The non evaluation words with large weight are selected from the sorted words and added to the feature set to form the output of mixed feature vector space.

From the above feature selection process, we can see that the feature selection part is carried out in the evaluation word set and non evaluation word set respectively. By increasing the number of evaluation words in the feature set, we can effectively strengthen the influence of evaluation words in emotional polarity classification. After the classification of text emotional polarity, the emotional tendency of text information is judged.

2.4 Realize the Emotional Element Extraction of Text Information

According to the above processing process, the neural network algorithm in natural language processing technology is used to extract emotional elements. A four layer RNN neural network is established to extract emotional elements. The number of output layer nodes of neural network is consistent with the number of online learning text information vectors of input network.

The emotional element features of the text information extracted above are used as the element screening feature template of the hidden layer 1, and the emotional polarity features are used as the element screening feature template of the hidden layer 2. After the two hidden layers are processed, the output layer outputs the emotional elements of the online learning text information. The neural network fitness functions used in this paper are as follows:

$$fitness = \frac{1}{\sum_{i=1}^N \sum_{j=1}^K (Y_j(i) - \bar{Y}_j(i))^2} \quad (8)$$

In the formula (8), $Y_j(i)$ is the actual output of the data i at the output node j is trained during training for a hybrid neural network; $\bar{Y}_j(i)$ is the expected output of training data i at output node j during training; K is the number of output nodes for the neural network; N is number of samples.

The neural network was trained using the emotion element training sample set, and it was calculated through multiple iterations, with the minimum output error as the neural network training stop criterion.

So far, the emotional elements in the text information are obtained according to the above process, and the extraction and processing of the emotional elements of the online learning text information by means of the natural language processing technology is completed.

3 Experimental Study

This section will conduct an experimental study on the practical application of the emotional element extraction method of online learning text information based on natural language processing technology.

3.1 Experimental Contents

The experiment is carried out on the data set of known emotional elements. The emotional element extraction method of online learning text information based on natural language processing technology proposed in this paper is compared with the emotional element extraction method based on domain dictionary and the emotional element extraction method based on attention mechanism. The advantages and disadvantages of the three methods are evaluated by comparing the extraction time, the accuracy of the extracted elements and the recall rate of the three methods. The experimental environment of the control group and the experimental group was ensured consistent.

3.2 Experimental Results

Figure 3 shows the comparison results between the accuracy and recall curve of emotional element extraction by three methods on the same data set.

According to Fig. 3, the accuracy-recall curve of this method is located above the other two method curves. The recall rate and accuracy rate of this method are above 0.9, while the accuracy rates of the other two methods are lower than 0.9. According to the meaning of accurate recall curve, the result of this method is more reliable. According to the meaning of the precision-recall curve, the results extracted in this method are more reliable.

Figure 4 shows the time-consuming comparison of the method for element extraction of different datasets.

It can be seen from the information in Fig. 4 that by increasing the number of groups from the experimental data to 80 groups, the time of extracting emotional elements in this method is less than 150 ms, which is shorter than that in the other two methods.

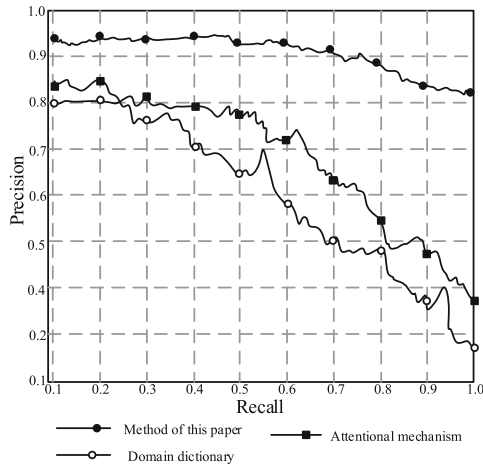


Fig. 3. Method accuracy – recall curve comparison

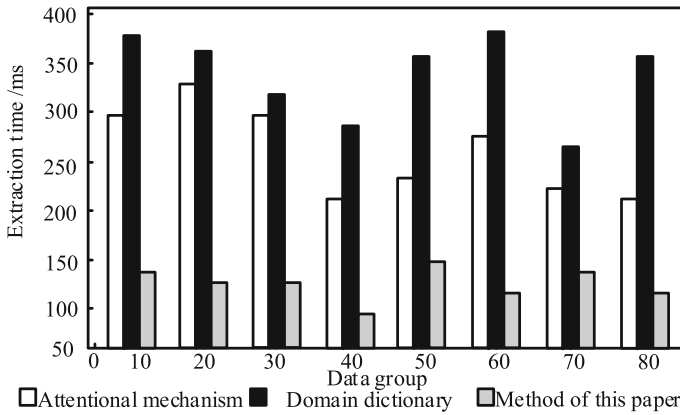


Fig. 4. Time-consuming comparison of emotional element extraction

Numerical results show that compared with the other two methods, the extraction time of this method is reduced by at least 35% and the efficiency is significantly improved.

According to the summary of the above experimental analysis, the emotional element extraction method of online learning text information based on natural language processing technology proposed in this paper has the advantages of accuracy, high efficiency and strong reliability.

4 Conclusion

In today’s society, because modern information technology and the Internet are developing rapidly, and people pay more and more attention to education, the teaching method of online learning has become a topic of concern, and it is also a new way for people to try to

learn knowledge. A large amount of text information is generated in the process of online learning, which contains a large number of emotional elements, which is very important to improve the teaching effect. This paper studies the emotional element extraction method of online learning text information based on natural language processing technology. Preprocess the text information of online learning; Split the preprocessed text into sentences, find new words and generate transaction items; Frequent noun items are extracted according to the theory of association rules; The filtering algorithm is used to filter irrelevant nouns and mine the emotional elements of text information; According to the credibility analysis theory, the text emotional polarity is judged, and the RNN neural network algorithm in natural language processing technology is innovatively used to extract the emotional elements of online learning text information. The test data show that when using this method to extract emotional elements from text information of online learning, the efficiency and accuracy of the method are significantly improved, which is helpful to optimize the effect of online learning. However, due to the limited conditions, this paper has some shortcomings, such as the scale of extracting targeted online learning text information. Future research could improve the scale of online learning text information.

References

1. Jiang, T., Li, Q.-X., Li, Q.-M.: Fine grained text sentiment classification method based on improved capsule network. *Comput. Simul.* **38**(10), 466–470 (2021)
2. Liu, W.-L., Huang, W.: Sentiment clustering of landscape constituents of Liuyuan garden based on domain dictionary. *Sci. Technol. Eng.* **21**(08), 3174–3179 (2021)
3. Wang, G.-S., Huang, X.-J., Min, L.: GRU neural network text emotion classification model based on multi-feature fusion. *J. Chin. Comput. Syst.* **40**(10), 2130–2138 (2019)
4. Yang, Q., Zhang, Y.-W., Zhu, L., et al.: Text sentiment analysis based on fusion of attention mechanism and BiGRU. *Comput. Sci.* **48**(11), 307–311 (2021)
5. Xu, J., Yan, C.: Slime mold foraging inspired feature selection algorithm and its application in sentiment recognition. *J. Nanjing Univ. Sci. Technol.* **45**(05), 596–605 (2021)
6. Han, P., Liu, S., Jia, Y., et al.: Sentiment analysis of semi-supervised Weibo text based on variational self-encoding. *Comput. Appl. Softw.* **38**(12), 280–285 (2021)
7. Yan, L., Zhu, X., Chen, X.: Emotional classification algorithm of comment text based on two-channel fusion and BiLSTM-attention. *J. Univ. Shanghai Sci. Technol.* **43**(6), 597–605 (2021)