

Two-Level LSTM for Sentiment Analysis With Lexicon Embedding and Polar Flipping

Ou Wu^{ib}, Tao Yang^{ib}, Mengyang Li, and Ming Li

Abstract—Sentiment analysis is a key component in various text mining applications. Numerous sentiment classification techniques, including conventional and deep-learning-based methods, have been proposed in the literature. In most existing methods, a high-quality training set is assumed to be given. Nevertheless, constructing a high-quality training set that consists of highly accurate labels is challenging in real applications. This difficulty stems from the fact that text samples usually contain complex sentiment representations, and their annotation is subjective. We address this challenge in this study by leveraging a new labeling strategy and utilizing a two-level long short-term memory network to construct a sentiment classifier. Lexical cues are useful for sentiment analysis, and they have been utilized in conventional studies. For example, polar and negation words play important roles in sentiment analysis. A new encoding strategy, that is, ρ -hot encoding, is proposed to alleviate the drawbacks of one-hot encoding and, thus, effectively incorporate useful lexical cues. Moreover, the sentimental polarity of a word may change in different sentences due to label noise or context. A flipping model is proposed to model the polar flipping of words in a sentence. We compile three Chinese datasets on the basis of our label strategy and proposed methodology. Experiments demonstrate that the proposed method outperforms state-of-the-art algorithms on both benchmark English data and our compiled Chinese data.

Index Terms—Flipping, lexicon embedding, long short-term memory (LSTM), sentiment analysis, text classification.

I. INTRODUCTION

TEXT is important in many artificial intelligence applications. Among various text mining techniques, sentiment analysis is a key component in applications, such as public opinion monitoring and comparative analysis. Sentiment analysis can be divided into three problems according to input texts, namely, sentence, paragraph, and document levels. This study focuses on sentence and paragraph levels.

Manuscript received October 1, 2019; accepted July 25, 2020. Date of publication September 23, 2020; date of current version May 19, 2022. This work was supported in part by Tianjin NSF under Grant 19JCZDJC41231300; in part by the AI Key Project of Tianjin under Grant 19ZXZNGX0050; in part by the Frontier Science and Technology Innovation Project under Grant 2019QY2404; in part by NSFC under Grant 61673377; and in part by the Zhejiang Lab Fund under Grant 2019KB0AB03. This article was recommended by Associate Editor S. Ozawa. (Corresponding author: Ou Wu.)

Ou Wu is with the Center for Applied Mathematics, Tianjin University, Tianjin 300072, China (e-mail: wuou@tju.edu.cn).

Tao Yang, Mengyang Li, and Ming Li are with the School of Computer Science, Civil Aviation University of China, Tianjin 300300, China (e-mail: yangtao087@gmail.com; limengyang99@gmail.com; liming2960@gmail.com).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TCYB.2020.3017378>.

Digital Object Identifier 10.1109/TCYB.2020.3017378

Text sentiment analysis is usually considered a text classification problem. Almost all existing text classification techniques are applied to text sentiment analysis [1]. Typical techniques include bag-of-words (BOW)-based [2], topic model-based [3], deep learning-based [4], and lexicon-based (or rule-based) methods [5].

Although many achievements have been made and sentiment analysis has been successfully used in various commercial applications, its accuracy can be further improved. The construction of a high-accuracy sentiment classification model usually entails the challenging compilation of training sets with numerous samples and sufficiently accurate labels. The reason behind this difficulty is two-fold. First, the sentiment is somewhat subjective, and a sample may receive different labels from different users. Second, some texts contain complex sentiment representations, and a single label is difficult to provide. We conduct a statistical analysis of public Chinese sentiment text sets in GitHub. The results show that the average label error is larger than 10%. This error value reflects the degree of difficulty of sentiment labeling.

Negation and interrogative sentences are difficult to classify when deep-learning-based methods are applied. Although lexicon-based methods can deal with particular types of negation sentences, their generalization capability is poor.

We address the above issues with a new methodology. First, we introduce a two-stage labeling strategy for sentiment texts. In the first stage, annotators are invited to label a large number of short texts with relatively pure sentiment orientations. Each sample is labeled by only one annotator. In the second stage, a relatively small number of text samples with mixed sentiment orientations are annotated, and each sample is labeled by multiple annotators. Second, we propose a two-level long short-term memory (LSTM) [6] network to achieve two-level feature representation and classify the sentiment orientations of a text sample to utilize two labeled datasets. Third, in the proposed two-level LSTM network, lexicon embedding is leveraged to incorporate linguistic features used in lexicon-based methods. Finally, the labels in a word-polar dictionary usually contain noise and the polarity of a word can also change in different contexts. A flipping model is proposed to model the sentiment polarity flipping of a word in a sentence.

Three Chinese sentiment datasets are compiled to investigate the performance of the proposed methodology. The experimental results demonstrate the effectiveness of the proposed methods. Our work is new in the following aspects.

- 1) A highly effective labeling strategy is adopted. Labeling a high-quality training set is difficult in sentiment

analysis. In our labeling strategy, samples are divided into ones with relatively pure sentiment orientations and ones with relatively complex sentiment orientations. This procedure is easily performed in practice.

- 2) A two-level LSTM network is proposed. Our labeling procedure yields two training sets with different sentiment levels; therefore, we propose a two-level LSTM network that can effectively utilize the two datasets.
- 3) Lexicon embeddings are introduced based on a new encoding strategy. To incorporate useful cues that are usually used in lexicon-based methods, an effective encoding strategy, namely, ρ -hot encoding, is proposed in this work to address the limitations of the classical one-hot encoding.
- 4) A flipping model is proposed to model polar flipping of words. Polar words are particularly important in sentiment analysis. However, the polar labels in a polar-word dictionary are not definitely right because the dictionary is very likely to contain label noise and the polarity of a word changes according to its context. For example, the polarity of “heavy” is positive in “the fish is heavy,” whereas it is negative in “the cell phone is heavy.” To this end, a flipping model¹ is established to describe the flipping of the polarity of words in texts.

The remainder of this article is organized as follows. Section II briefly reviews related work. Section III describes our methodology. Section IV reports the experimental results and Section V concludes the study.

II. RELATED WORK

A. Text Sentiment Analysis

Sentiment analysis aims to predict the sentiment polarity of an input text sample. Sentiment polarity can be divided into negative, neutral, and positive in many applications.

The existing sentiment classification methods can be roughly divided into two categories, namely, lexicon-based and machine-learning-based methods [9]. Lexicon-based methods [10] construct polar and negation word dictionaries. A set of rules for polar and negation words is compiled to judge the sentiment orientation of a text document. This method cannot effectively predict implicit orientations. Machine-learning-based methods [11], [12] utilize a standard binary or multiclassification approach. Different feature extraction algorithms, including BOW [13] and part of speech (POS) [11], are used. Word embedding and deep neural networks have recently been applied to sentiment analysis, and promising results have been obtained [14], [15].

B. Lexicon-Based Sentiment Classification

Lexicon-based methods are actually implemented in an unsupervised manner. They infer the sentiment categories of input texts on the basis of polar and negation words. The primary advantage of these methods is that they do not require labeled training data. The key to lexicon-based methods is

¹In theory, our model should be useful in any arbitrary method that leverages the polar labels of words as supervised information, such as [7] and [8].

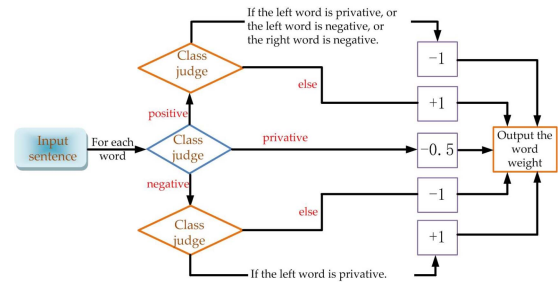


Fig. 1. Lexicon-based approach for sentiment classification.

the lexical resource construction, which maps words into a category (positive, negative, neutral, or negation). Senti-WordNet [16] is a lexical resource for English text sentiment classification. For Chinese texts, Senti-HowNet is usually used.

Fig. 1 characterizes a typical lexicon-based sentiment classification approach. The approach iteratively checks each word in an input sentence from left to right. The weight score of each word is calculated according to the procedure shown in Fig. 1. The final sentiment score is the average score of the words with weight scores. The scores of positive, neutral, and negative sentiments are denoted as “+1,” “0,” and “−1,” respectively. According to the lexicon-based algorithm shown in Fig. 1, the sentiment score of “it is not bad” is 0.25, and the sentiment score of “it is good” is 1. However, the score of “it is not so bad” is −0.75, and this score is definitely wrong. Therefore, machine-learning (including feature learning) methodologies have become the mainstream in sentiment analysis.

C. Deep-Learning-Based Sentiment Classification

Deep learning (including word embedding [17]) has been applied to almost all text-related applications, such as translation [18], quality assurance [19], recommendation [20], and categorization [21]. Popular deep neural networks are divided into convolutional neural networks (CNNs) [22] and recurrent neural networks (RNNs) [23], [24]. Both are utilized in sentiment classification [25], [26]. Kim [4] investigated the use of CNN in sentence sentiment classification and achieved promising results. LSTM [27], a classical type of RNN, is the most popular network used for sentiment classification. A binary-directional LSTM [28] with an attention mechanism is demonstrated to be effective in sentiment analysis.

Deep-learning-based methods rarely utilize the useful resources adopted in lexicon-based methods. Qian *et al.* [8] incorporated lexicon-based cues into the training of an LSTM-based model. Their proposed method relies on a new loss function that considers the relationships between polar or certain types of words (e.g., negation) and those words next to them in input texts. Our study also combines lexical cues into LSTM. Nevertheless, unlike Qiao *et al.*'s study that implicitly used lexical cues, the present work explicitly uses lexical cues in the LSTM network. Shin *et al.* [29] combined the lexicon embeddings of polar words with word embeddings for sentiment classification. The difference between our approach and the method proposed by Shin *et al.* is discussed in Section III-C5.

Numerous studies on aspect-level sentiment analysis exist [30]–[32]. This problem is different from the sentiment classification investigated in this study.

III. METHODOLOGY

This section first introduces our two-stage labeling procedure. A two-level LSTM is then proposed. Lexicon embedding is finally leveraged to incorporate lexical cues.

A. Two-Stage Labeling

As stated earlier, the sentiment is subjective, and texts usually contain mixed sentiment orientations. Therefore, texts' sentiment orientations are difficult to label. In our study, three sentiment labels, namely, positive, neutral, and negative, are used. The following sentences are taken as examples.

S1: The service is poor. The taste is good, but the rest is not so bad.

S2: The quality of the phone is good, but the appearance is just so-so.

In user annotation, the labels of these two sentences depend on users. If a user is concerned about service, then the label of *S1* may be “negative.” In contrast, for another user who does not care about service, the label may be “positive.” Similarly, a user may label *S2* as positive if he cares about quality. Another user may label it as negative if the conjunction “but” attracts the user’s attention more. Another user may label it as “neutral” if they are concerned about quality and appearance.

The underlying reason is that sentiment is more subjective than semantics. In related research on subjective categorization, such as visual aesthetics, each sample is usually repeatedly annotated by multiple annotators, and the average label is taken as the final label of the sample. This labeling strategy can also be applied to text sentiment annotation. However, we argue that this strategy is unsuitable for a (relatively) large number of samples. The reason lies in the following two aspects.

- 1) Multiple annotators for a large number of datasets require a large budget.
- 2) In our practice, annotators claim that their judgment criteria on sentiment become fused on texts with mixed sentiment orientations (e.g., *S1* and *S2*) over time during labeling, and they become bored accordingly.

A two-stage labeling strategy is adopted in this study. In the first stage, each sentence/paragraph is divided into several clauses according to punctuation. The sentiment of each partitioned clause is relatively easy to annotate; therefore, each clause is labeled by only one user. In the second stage, a relatively small-sized sentence/paragraph set is labeled, and each sentence is labeled by all involved annotators. We still take the two sentences, *S1* and *S2*, as examples. *S1* and *S2* are split into clauses, as shown below.

- 1) *S1*:
 - a) *S1.1*: The service is poor.
 - b) *S1.2*: The taste is good.
 - c) *S1.3*: But the rest is not so bad.
- 2) *S2*:
 - a) *S2.1*: The quality of the phone is good.

- b) *S2.2*: But the appearance is just so-so.

Each of the above clauses is labeled by only one annotator. The annotation in the first stage is easy to perform; thus, the number of clauses can be larger than the number of sentences used in the second labeling stage.

B. Two-Level LSTM

Given two training datasets (denoted by *T1* and *T2*), a new learning model should be utilized. LSTM² is a widely used deep neural network in deep-learning-based text classification.

LSTM is a typical RNN model for short-term memory, which can last for a long period of time. An LSTM is applicable to classify, process, and predict time-series information with given time lags of unknown size. A common LSTM block is composed of a cell, an input gate, an output gate, and a forget gate.

When LSTM is used to classify an input sentence, the hidden vectors (h_t) of each input vector are summed to form a dense vector that can be considered the feature representation of the input sentence, that is

$$v_t = \sum_t h_t. \quad (1)$$

In many applications, a bidirectional LSTM (Bi-LSTM) structure is usually used. In Bi-LSTM, forward and backward information are considered for information at time t ; hence, the context is modeled. Bi-LSTM is thus significantly reasonable for text processing tasks. In our two-level LSTM, Bi-LSTM is used at each level.

The output hidden state at time t of a Bi-LSTM block can be described as follows:

$$\begin{aligned} \vec{h}_t &= \vec{o}_t \otimes \tanh(\vec{c}_t) \\ \overleftarrow{h}_t &= \overleftarrow{o}_t \otimes \tanh(\overleftarrow{c}_t) \\ h_t &= [\vec{h}_t, \overleftarrow{h}_t] \end{aligned} \quad (2)$$

where \vec{h}_t , \vec{o}_t , and \vec{c}_t are the corresponding vectors at time t in the forward LSTM block and \overleftarrow{h}_t , \overleftarrow{o}_t , and \overleftarrow{c}_t are the corresponding vectors at time t in the backward LSTM block. $H = \{h_1, \dots, h_T\}$. When attention is used, the dense feature vector γ of an input sentence is calculated as follows:

$$\begin{aligned} \alpha &= \text{softmax}(\mathbf{w}^T H) \\ \gamma &= H\alpha^T \end{aligned} \quad (3)$$

where α is the vector that consists of attention weights.

Our proposed network consists of two levels of the LSTM network. In the first level, a Bi-LSTM is used and learned on the basis of the first training set *T1*. This level is a conventional sentiment classification process. The input of this level is a clause, and the input x_t is the embedding of the basic unit of the input texts.³ The network is shown in Fig. 2(a).

In the second level, a Bi-LSTM is also used and learned on the basis of the second training set *T2*. The input of this

²CNN is another widely used text classification model. Our idea can also be applied to CNN.

³In English texts, the basic unit is usually a word; in Chinese texts, the basic unit is a Chinese word or character.

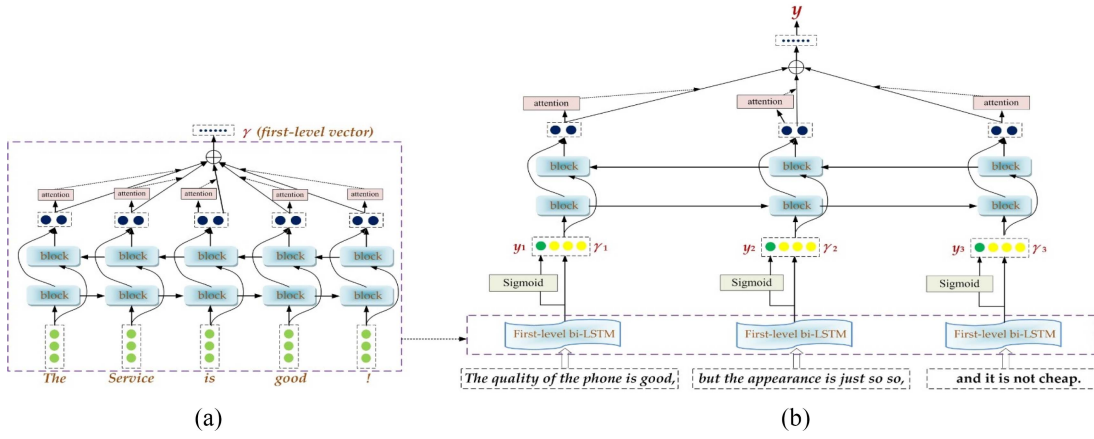


Fig. 2. Proposed two-level LSTM network (this network does **not** including the lexicon embedding and the polar flipping model which are introduced in Figs. 5 and 7). (a) First-level bi-LSTM. (b) Whole two-level structure.

level is a sentence or a paragraph. The input x_t consists of two parts.⁴ The first part is the feature vector of the t -th clause. The feature vector is generated by the first-level network. In other words, the dense feature shown in Fig. 2(a) (γ) is used. The second part is the sentiment score (not predicted label) output by the first-level network. The sentence S1 (The service is poor. The taste is good, but the rest is not so bad.) used in Section III-A is taken as an illustrative example. S1 contains three clauses. Let σ be the sigmoid function. The input vector of S1 can be represented by

$$S1 : \{\eta_1, \eta_2, \eta_3\}$$

where

$$\begin{aligned} \eta_i &= \{y_i^{(1)}, \gamma_i^{(1)}\} \\ y_i^{(1)} &= \sigma(W\gamma_i^{(1)} + b) \\ i &= 1, 2, 3 \end{aligned} \quad (4)$$

where $y_i^{(1)}$ is the output score of the i th clause by the first-level LSTM and $\gamma_i^{(1)}$ is the feature representation of the i th clause by the first LSTM. The network of the entire two-level network is shown in Fig. 2(b). The loss function of the entire network is defined as follows:

$$l = \sum_n \left[\text{loss}(y_n, y'_n) + \frac{\lambda}{n_l} \sum_i \text{loss}(y_{ni}, y'_{ni}) \right] \quad (5)$$

where y_n and y'_n are the true and predicted labels of the n th sample, respectively, y_{ni} and y'_{ni} are the true and predicted label of the i th clause of the n th sample, respectively; λ is the parameter, and n_l is the number of clauses in the i th sample.

C. Lexical Embedding

Lexicon embedding aims to integrate a wide range of lexical cues into the two-level LSTM network. Based on our empirical analysis and previous studies, key lexical words, POS, and conjunction cues are considered. The proposed lexicon embedding is based on ρ -hot encoding. Therefore, ρ -hot encoding is first described.

⁴The third part is lexicon embedding which will be introduced in Section III-C4.

1) ρ -Hot Encoding: For categorical data, one-hot encoding is the most widely used encoding strategy when different categories are independent.⁵ For example, if one-hot encoding is used to represent three categories, namely, positive, neutral, and negative, the encoding vectors for the three categories are $[1, 0, 0]^T$, $[0, 1, 0]^T$, and $[0, 0, 1]^T$, respectively.

In this work, many lexical cues are categorical data, and different categories are independent. These lexical cues can directly be represented by one-hot encoding. The encoded vectors for lexical cues are then concatenated with other vectors, such as character/word embedding. Based on our empirical evaluations, a more effective encoding is proposed based on the conventional one-hot encoding. The encoding strategy is defined as follows:

$$\rho\text{-hot encoding} : v_{\rho,n} = \rho \cdot v_1 \otimes \mathbf{1}_{n \times 1} \quad (6)$$

where $v_{\rho,n}$ is the ρ -hot encoded vector, ρ is the proportion parameter, v_1 is the one-hot encoded vector, $\mathbf{1}_{n \times 1}$ is an n -dimensional vector, and \otimes is the tensor product. If both ρ and n are equal to 1, then ρ -hot encoding is reduced to one-hot encoding. The parameter n is applied to increasing the length of the final encoded vector.

When ρ equals 0.6 and n equals 4. For a 3-D one-hot encoded vector v_1 , we have

$$\begin{aligned} \text{if } v_1 &= [1, 0, 0]^T \\ v_{\rho} &= 0.6 \cdot v_1 \otimes \mathbf{1}_{4 \times 1} \\ &= [0.6, 0.6, 0.6, 0.6, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]^T \\ \text{if } v_1 &= [0, 1, 0]^T \\ v_{\rho} &= 0.6 \cdot v_1 \otimes \mathbf{1}_{4 \times 1} \\ &= [0, 0, 0, 0, 0.6, 0.6, 0.6, 0.6, 0, 0, 0, 0, 0, 0, 0]^T \\ \text{if } v_1 &= [0, 0, 1]^T \\ v_{\rho} &= 0.6 \cdot v_1 \otimes \mathbf{1}_{4 \times 1} \\ &= [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0.6, 0.6, 0.6, 0.6, 0]^T. \end{aligned} \quad (7)$$

⁵When different categories are correlated, sophisticated encoding strategies can be utilized. For example, one-hot is a traditional encoding for words. Many word-embedding methods are proposed with consideration of the relation among words.

The numerical examples in (7) indicate that the obtained ρ -hot encoded vector (v_ρ) can have varied values for nonzero elements and longer lengths compared with the corresponding one-hot encoded vector (v_1). Similar yet simple tricks for the increasing of the vector length have been successfully used in previous studies [33], [34].

So far, there is no accurate theoretical explanation for our proposed ρ -hot encoding. Nevertheless, one-hot encoding presents two main limitations when the encoded vector is concatenated with other vectors.

- 1) The value difference between the elements of one-hot encoded vectors and those of other encoded vectors (e.g., word embedding vectors) may be large. Fig. 3 shows the histogram of the values of the elements of the word embedding vectors. The magnitudes of most elements are smaller than 1.
- 2) The lengths of one-hot encoded vectors are usually shorter than those of other encoded vectors. Consequently, the proportion of the one-hot encoded part is small in the concatenated vectors.

The above two limitations may affect the final sentiment analysis performance, whereas our proposed ρ -hot encoding alleviated these two limitations.

2) *Embedding for Key Lexical Words*: Most lexicon-based sentiment methods rely on four types of words, namely, positive, negative, neutral, and negation. These words are useful cues for predicting the sentiment labels of input texts. The incorporation of these words should also be useful. A previous study has shown that a typical document comprises approximately 8% of such words [35]. Sentiments expressed in a conditional sentence can be difficult to determine due to the semantic condition. The sentiment polarities of interrogative sentences are also difficult to classify according to our empirical study. Because the automatically judging whether a sentence is conditional or interrogative is challenging, we directly consider suppositive and interrogative words.

Five types of words, namely, positive (Pos), negative (Neg), negation (Nt), suppositive (Sup), and interrogative (Int), are represented by the proposed encoding method. The remaining words, which do not belong to any of the above five types, are called “others (Oth)” instead of “neutral” because some words, such as “the,” are unrelated to “sentiment.” The value of n in (6) is set as 10. The encoded vectors are as follows:

$$\begin{aligned} \text{Pos} &: [\rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}]^T \\ \text{Neg} &: [\mathbf{0}_{1 \times 10}, \rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}]^T \\ \text{Nt} &: [\mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}]^T \\ \text{Sup} &: [\mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}]^T \\ \text{Int} &: [\mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \rho_{1 \times 10}, \mathbf{0}_{1 \times 10}]^T \\ \text{Oth} &: [\mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \rho_{1 \times 10}]^T. \end{aligned}$$

In the proposed ρ -hot embedding, the parameter ρ can be learned during training.

Certain types (e.g., positive, negative, and negation) of words should play more important roles than other words do in texts; therefore, their embeddings are also used in the attention layer. A new LSTM based on our lexicon embedding is

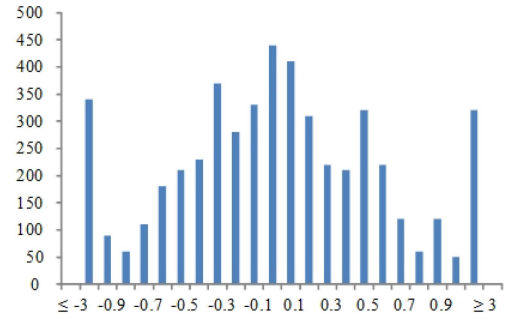


Fig. 3. Histogram of the values in word embedding vectors. Most values are smaller than 1.

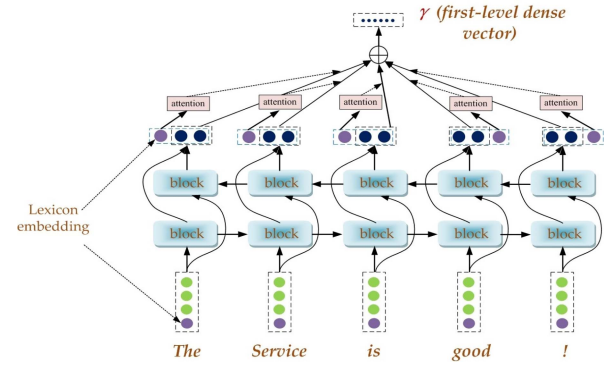


Fig. 4. First-level LSTM with lexicon embedding in both the input and attention layers.

proposed, as shown in Fig. 4. The attention layer and final dense vector of the network in Fig. 2(a) are calculated as follows:

$$\begin{aligned} \alpha_t &= \text{softmax} \left(W_\alpha \begin{bmatrix} h_t \\ l_t \end{bmatrix} + b_i \right) \\ \gamma &= \sum_t \alpha_t \begin{bmatrix} h_t \\ l_t \end{bmatrix} \end{aligned} \quad (8)$$

where α_t is the attention weight for the t -th input, l_t is the lexicon embedding for key lexical words for the t -th input, and γ is the final dense vector. Equation (7) is used in the first-level LSTM.

3) *Embedding for POS*: POS is usually used as a key cue in sentiment analysis [36]. Intuitively, adjective and adverb words are very likely to be more important than some other types of words, such as the adjective and article. To this end, we use additional lexicon embedding based on POS information. This new lexicon embedding is also applied to the attention layer. The motivation lies in that certain types of POS should play important roles in sentiment.

The proposed ρ -hot embedding is still applied to POS types in this study. According to our initial case studies, eight POS types are considered. They are noun, adjective, verb, pronoun, adverb, preposition, accessory, and others. The eight POS types are represented by the proposed ρ -hot encoding. We let n in (6) be 10. The first three POS types are as follows:

$$\begin{aligned} \text{Noun} &: [\rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \dots, \mathbf{0}_{1 \times 10}]^T \\ \text{Adj} &: [\mathbf{0}_{1 \times 10}, \rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \dots, \mathbf{0}_{1 \times 10}]^T \end{aligned}$$

$$\text{Verb} : [\mathbf{0}_{1 \times 10}, \mathbf{0}_{1 \times 10}, \rho_{1 \times 10}, \mathbf{0}_{1 \times 10}, \dots, \mathbf{0}_{1 \times 10}]^T.$$

When POS embedding is used, the attention layer and final outputs of the network in (3) become

$$\begin{aligned} \alpha_t &= \text{softmax} \left(W_\alpha \begin{bmatrix} h_t \\ l_t \\ \eta_t \end{bmatrix} + b_i \right) \\ \gamma &= \sum_t \alpha_t \begin{bmatrix} h_t \\ l_t \end{bmatrix} \end{aligned} \quad (9)$$

where η_t is the lexicon embedding for POS of the t -th input.

4) *Embedding for Conjunction*: Conjunction words play important roles in sentiment analysis [37]. For example, conjunctions, such as “but” and “moreover,” usually indicate the focus of texts and attract readers attention. These words may be useful cues for attention inference in the second level of our network. Therefore, conjunctions are considered in the input of the second-level LSTM.

Once a set of conjunction words is compiled, ρ -hot embedding is used. In our experiments, the number of conjunction words (including thesaurus) is 13 for Chinese texts and 26 for English texts. Therefore, the parameter n for Chinese conjunction words in (6) is set as 1.

When conjunction embedding is used for the second-level layer, the attention layer and final outputs of the network in Fig. 2(b) are calculated as follows:

$$\begin{aligned} \beta_t &= \text{softmax} \left(W_\beta \begin{bmatrix} y_t^{(1)} \\ h_t^{(2)} \\ \omega_t^s \\ \omega_t^e \end{bmatrix} + b'_i \right) \\ \gamma^{(2)} &= \sum_t \beta_t \begin{bmatrix} h_t^{(2)} \\ y_t^{(1)} \end{bmatrix} \end{aligned} \quad (10)$$

where β_t is the attention weight for the t -th input clause, $h_t^{(2)}$ is the hidden vector of the second-level LSTM, ω_t^s and ω_t^e are the conjunction embeddings for the first and last words in the t -th input clause, respectively, and $\gamma^{(2)}$ is the final dense vector used for the final classification.

The two-level network with lexicon embedding is shown in Fig. 5. The lexicon embedding is used in both input and attention layers.

5) *Differences Between Our and Existing Lexicon Embedding*: Shin *et al.* [29] also embedded lexical information into sentiment analysis. Three major differences exist between our method and the method proposed by Shin *et al.* [29].

- 1) The lexicon embedding proposed by Shin *et al.* uses one-hot encoding, whereas the proposed method uses a new encoding strategy that can be considered a soft one-hot encoding.
- 2) The lexicon embedding proposed by Shin *et al.* extends the lengths of raw encoded vectors. However, the extension aims to keep the lengths of lexical and word embeddings equal. Their extension method also only relies on zero padding and is thus different from the proposed method.

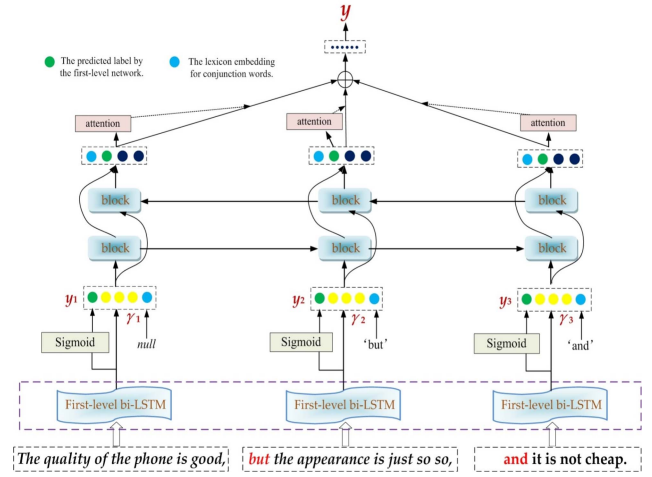


Fig. 5. Whole two-level LSTM network with lexicon embedding in both the input and attention layers.

- 3) Only sentimental words are considered in the lexicon embedding proposed by Shin *et al.* On the contrary, sentimental words, POS, and conjunctions are considered in our work.

D. Polar Flipping Model

The polarity of words is crucial in sentiment analysis. Most rule-based methods mainly rely on polar and negation words. The word polarity is also usually applied as (partial) supervised information in learning-based methods. Qian *et al.* [8] leveraged the polar labels to regularize the hidden state of each word in training texts. They utilized the KL divergence between the polar label and the predicted label of a word as the regularization term during training. Although word-level polar information is proven to be useful in sentiment analysis, the application of word-level polar information suffers from the following defects.

- 1) Because the number of words usually exceeds 10 000, the labeling of all words' polarity is not an easy task. Labeling error is thus unavoidable.
- 2) The polarities of some words highly depend on their contexts. Therefore, the polarity of a word may vary in different texts.

In this work, the word-level polar labels are considered noisy labels. In noisy-label learning [38], a flipping model is usually used to model the relationships between (possibly) noisy and true labels. Inspired by noisy-label learning, a polar flipping model is proposed. Let P_1 , P_2 , and P_3 represent the three polar labels “positive,” “negative,” and “neutral,” respectively. Without loss of generalization, the label of a given word is assumed to be P_1 . Let x be the representation of the given word and C be the representation of the context information. Then, the polar flipping model is described as follows:

$$P = \sigma_1(x, C)P_1 + [1 - \sigma_1(x, C)]\{\sigma_2(x, C)P_2 + [1 - \sigma_2(x, C)]P_3\} \quad (11)$$

where both σ_1 and σ_2 are sigmoid functions. Equation (10) describes the flipping of the polarity of a word given its

context. If $\sigma_1(x, C) == 1$, then $P = P_1$; else if $\sigma_2(x, C) == 1$, then $P = P_2$; else $P = P_3$.

In this work, the hidden state vector h of the word in LSTM is used to represent (x, C) in (10). Therefore, (10) becomes

$$P = \sigma_1(h)P_1 + [1 - \sigma_1(h)]\{\sigma_2(h)P_2 + [1 - \sigma_2(h)]P_3\}. \quad (12)$$

In noisy-label learning, the noisy rate is usually assumed to be small. In our model, to control the flipping rate, we take the flipping loss as an additional regularization term. Consequently, the loss function defined in (5) becomes

$$l_{\text{total}} = \sum_n \left[\text{loss}_1(y_n, y'_n) + \frac{\lambda_1}{n_I} \sum_i \text{loss}_2(y_{ni}, y'_{ni}) + \frac{\lambda_2}{n_{IK}} \sum_k \text{loss}_3(P_{nik}, P'_{nik}) \right] \quad (13)$$

where P_{nik} and P'_{nik} are the true and the flipped polar labels of the k th word in the i th clause in the n th sample, λ_1 and λ_2 are balance parameters, and n_I is the number of words in the i th clause of the n th sample. The functions loss_1 and loss_2 in (12) are the cross-entropy loss, and the function loss_3 is calculated as follows:

$$\text{loss}_3(P_{nik}, P'_{nik}) = 1 - \sigma_1(h_{nik}) \quad (14)$$

where h_{nik} is the hidden vector of the k th word in the i th clause in the n th sample. According to the loss function defined in (12), our method can be viewed as being added supervised information in the middle layers (both word and clause levels) of the entire network as shown in Fig. 6. The first-level LSTM in the proposed two-level network described in Fig. 2(b) is shown in Fig. 7.

Based on Fig. 7, there are three main differences between our's and existing hierarchal networks.

- 1) Our network depends on a relatively new labeling strategy, namely, two-stage labeling.
- 2) Each level of our network is associated with supervised information, whereas only the top level is with supervised information in almost all existing hierarchal networks.
- 3) Each level of our network is associated with a loss term in our network, whereas there is usually only one loss in existing hierarchal networks.
- 4) Plentiful lexical cues are embedded into our network, whereas there is usually only one type of lexical cues considered in existing networks.

E. Learning Details

The algorithmic steps of the entire learning procedure for the proposed ρ -hot lexicon embedding-based two-level LSTM (called ρ TI-LSTM) are shown in Algorithm 1. In Algorithm 1, T_1 refers to the training data that consist of clauses and the labels obtained in the first-stage labeling procedure. T_2 refers to the training data that consist of sentences and the labels obtained in the second-stage labeling procedure.

The proposed two-level LSTM can be applied to texts with arbitrary languages. Word information is required in lexical

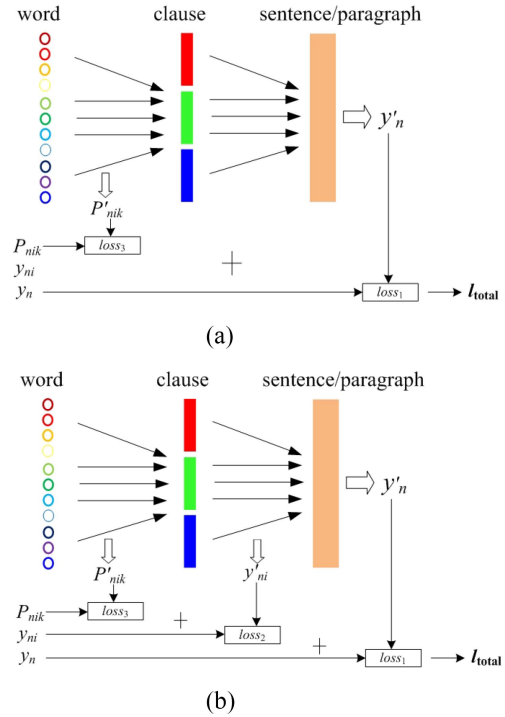


Fig. 6. Illustrative network structure with supervised information in middle layers. (a) Loss in existing methods. (b) Loss in our method.

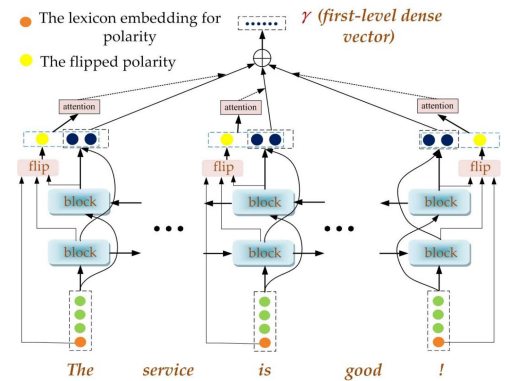


Fig. 7. First-level network with the flipping model. (To simplify the illustration, only polar information is shown in the lexicon embedding.)

construction regardless of whether a character or word embedding is used. The reason is that the three types of lexicon embeddings are performed at the word level. Therefore, when character embedding is used, the lexicon embedding of each character is the lexicon embedding of the word containing it.

IV. EXPERIMENTS

This section shows the evaluation of the proposed methodology in terms of the two-level LSTM network and each part of the lexicon embedding. Both English benchmark text corpora and Chinese text corpora are used.

A. Results on Chinese Corpora

In the experiments, three competing algorithms, namely, BOW, CNN, and (conventional) LSTM, are used. The lexicon

Algorithm 1 ρ T1-LSTM

Input: Training sets T1 and T2; dictionary of key lexical words; POS for each word; dictionary of conjunction words; character/word embeddings for each character/word; parameters λ_1 and λ_2 .

Output: A trained two-level LSTM for sentiment classification.

Steps:

1. Construct the ρ -hot-based embedding vector for each word (including punctuation) in the clauses in T1. The embeddings include the character/word and lexicon embeddings of each character/word;
2. Construct the embedding vector for each conjunction word in each clauses as the input for the second-level LSTM;
3. Train the two-level LSTM on the basis of the input embedding vectors and labels of polar words, the T1 text clauses, and the T2 text samples; the loss function is defined in (12).

embedding-based method proposed by Shin *et al.* [29] discussed in Section III-C5 is also compared.

For BOW, term frequency-inverse document frequency is utilized to construct features. Ridge regression [39] is used as a classifier. For CNN, a three-channel CNN is used. For LSTM, one-layer and two-layer Bi-LSTM with attention are adopted, and the results of the network with superior performance are presented. CNN and LSTM are performed on TensorFlow, and default parameter settings are followed.

The key parameters are searched as follows. The embedding dimensions of characters and words are searched in [100, 150, 200, 250, 300]. The parameter n in ρ -hot encoding is searched in [1, 3, ..., 15]. The parameters λ_1 and λ_2 are searched in [0.001, 0.01, 0.1, 1, 10]. Baidu Chinese word segmentation API is used.

1) *Experimental Data and Labeling:* We compile three Chinese text corpora from online data for three domains, namely, “hotel,” “mobile phone (mobile),” and “travel.” All texts are about user reviews. Each text sample collected is first partitioned into clauses according to Chinese tokens.⁶ Three clause sets are subsequently obtained from the three text corpora.

The labels “+1,” “0.5,” and “0” correspond to the three sentiment classes positive, neutral, and negative, respectively. The text data are labeled according to our two-stage labeling strategy.

- 1) In the first stage, only one user is invited to label each clause sample as the sentiment orientations for clauses (or subsentences) are easy to label.
- 2) In the second stage, five users⁷ are invited to label each text sample in the three raw datasets. The average score of the five users in each sample is calculated. Samples with average scores located in [0.6, 1] are labeled as

⁶Token-based token is inaccurate for English text partition. Nevertheless, the segment results for Chinese texts are acceptable. A more reasonable way will be investigated in our future work.

⁷Five graduate students, including three males and two females, were invited to label the data.

TABLE I
DETAILS OF THE THREE DATA CORPORA. EACH CORPUS CONSISTS OF RAW SAMPLES (SENTENCES OR PARAGRAPHS) AND PARTITIONED CLAUSES (SUBSENTENCES)

Data corpus		raw	clauses
Travel	Pos.	1567	3490
	Neu.	576	5168
	Neg.	1957	2633
	Total	4100	11291
Hotel	Pos.	1586	3987
	Neu.	401	2123
	Neg.	1838	5154
	Total	3825	11264
Mobile	Pos.	1400	2788
	Neu.	589	2375
	Neg.	1494	2955
	Total	3483	8118

TABLE II
NUMBERS OF FIVE TYPES OF KEY LEXICAL WORDS

Data corpus	Travel	Hotel	Mobile
Positive	366	254	358
Negative	327	194	382
Negation	61	61	61
Interrogative	48	48	48
Suppositive	18	18	18

positive. Samples with average scores located in [0, 0.4] are labeled as negative. Others are labeled as neutral.

The details of the labeling results are shown in Table I.

All the training and test data and the labels are available online.⁸

In our experiments, the five types of key lexical words introduced in Section III-C2 are manually constructed. The details of the five types of words are listed in Table II.⁹ The conjunction words are also manually constructed. The number of conjunction words used in the experiments is 169.

In each experimental run, the training set is compiled on the basis of the training data listed in Table I. The compiling rule is specified before each experimental run. The test data are fixed to facilitate experimental duplication and comparison by other researchers. The data with the fixed split are available at our Github page.

2) *Results of Existing Methods:* In this section, each of the three raw datasets (associated with their labels) shown in Table I is used. The clause data are not used. In other words, the training data used in this section are the same as those used in previous studies. For each data corpus, 1000 raw data samples are used as the test data, and the rest are used as the training data. The involved algorithms are detailed as follows.

- 1) *CNN-C* denotes the CNN with (Chinese) character embedding.
- 2) *CNN-W* denotes the CNN with (Chinese) word embedding.

⁸<https://github.com/Tju-AI/two-stage-labeling-for-the-sentiment-orientation/tree/master/data>

⁹The five types of key lexical words are also available and introduced in our Github project page.

TABLE III
CLASSIFICATION ACCURACIES OF EXISTING
ALGORITHMS ON RAW SAMPLES

Data corpus	Travel	Hotel	Mobile
CNN-C	0.723	0.698	0.727
CNN-W	0.731	0.729	0.748
CNN-Lex-C	0.744	0.734	0.731
CNN-Lex-W	0.758	0.764	0.755
Bi-LSTM-C	0.754	0.753	0.805
Bi-LSTM-W	0.746	0.785	0.809
Lex-rule	0.556	0.539	0.684
BOW+LR	0.713	0.678	0.702
BOW+SVM	0.746	0.735	0.718
BOW+NB	0.698	0.712	0.687
BOW+RF	0.733	0.756	0.743

- 3) *CNN-Lex-C* denotes the algorithm that also integrates polar words in CNN, which is proposed by Shin *et al.* [29]. The (Chinese) character embedding is used.
- 4) *CNN-Lex-W* denotes the algorithm that also integrates polar words in CNN, which is proposed by Shin *et al.* [29]. The (Chinese) word embedding is used.
- 5) *Bi-LSTM-C* denotes the BI-LSTM with (Chinese) character embedding.
- 6) *Bi-LSTM-W* denotes the Bi-LSTM with (Chinese) word embedding.
- 7) *Lex-rule* denotes the rule-based approach shows in Fig. 1. This approach is unsupervised.
- 8) *BOW* denotes the conventional algorithm that is based on BOWs features. Four shallow classifiers are used, namely, logistic regression (LR), support vector machine (SVM), naive Bayes (NB), and random forest (RF). For SVM, the parameters C and g are searched via five-fold cross-validation from $\{0.1, 1, 5, 10, 100\}$ and $\{0.01, 0.1, 1, 5, 10\}$, respectively. For RF, the number of trees is searched via five-fold cross-validation from $\{10, 20, 50, 100, 200, 500\}$.

The accuracies of the above algorithms are listed in Table III. Overall, Bi-LSTM significantly outperforms CNN and BOW ($p < 0.01$) based on the t -test. This conclusion is in accordance with the conclusion that RNN performs efficiently against CNN in a broad range of natural language processing (NLP) tasks on the basis of extensive comparative studies [40]. In addition, CNN-Lex outperforms CNN under both character and word embeddings ($p < 0.05$), which suggests that lexicon cues are useful in sentiment analysis. Lex-rule achieves the lowest accuracies on all the three datasets. Considering that the performances of (traditional) CNN, Lex-rule, and BOW are relatively poor, they are not applied in the remaining parts.

3) *Results of Two-Level LSTM Without Lexicon Embedding:* In this experimental comparison, the proposed two-level LSTM is evaluated, whereas lexicon embedding is not used in the entire network. The primary goal is to test whether the introduced two-stage labeling and the two-level network structure are useful for sentiment analysis.

The raw and clause data listed in Table I are used to perform the two-level LSTM. TI-LSTM denotes the two-level LSTM.

TABLE IV
CLASSIFICATION ACCURACIES OF COMPETING ALGORITHMS

Data corpus	Travel	Hotel	Mobile
TI-LSTM-C(R+C)	0.801	0.813	0.820
TI-LSTM-W(R+C)	0.770	0.772	0.820
CNN-Lex-C (R+C)	0.755	0.770	0.749
CNN-Lex-W (R+C)	0.773	0.776	0.786
Bi-LSTM-C (R+C)	0.781	0.784	0.817
Bi-LSTM-W (R+C)	0.762	0.789	0.813
Baseline (best in Table 3)	0.758	0.785	0.809

TABLE V
CLASSIFICATION ACCURACIES OF TWO-LEVEL LSTM
WITH LEXICON EMBEDDING

Data corpus	Travel	Hotel	Mobile
ρ TI-LSTM-C without flipping	0.816	0.837	0.826
ρ TI-LSTM-W without flipping	0.800	0.810	0.841
ρ TI-LSTM-C	0.817	0.837	0.837
ρ TI-LSTM-W	0.804	0.821	0.847

“R+C” refers to the mixed data of raw and clause data. The test data are still the 1000 samples used in Section IV-A2 for each corpus. Table IV shows the classification accuracies. To ensure that the results differ from those in Table III, we explicitly add R+C after each algorithm in Table IV. In the last line of Table IV, the base results for each corpus in Table III are also listed.

On all the three data corpora, the proposed two-level network (without lexicon embedding) with character embedding, TI-LSTM-C, does not significantly outperform other involved competing methods by conducting t -test. However, it achieves the highest accuracies on all three data corpora. The results in Table IV indicate that the performances of TI-LSTM on the mixed training and test data (R+C) are better than those of Bi-LSTM. This comparison indicates that the proposed two-level LSTM is useful.

In addition, for the involved algorithms, most accuracies achieved on R+C are higher than the best results only achieved on “R” listed in Table III. This comparison suggests that the introduced two-stage labeling is useful.

4) *Results of the Entire Two-Level LSTM:* In this experimental run, both lexicon embedding and the flipping model are used in the proposed two-level LSTM or ρ TI-LSTM. Table V shows the results. In order to assess the usefulness of the flipping model, the results of ρ TI-LSTM without considering flipping are also present.

The performances of TI-LSTM with both lexicon embedding and the flipping model (i.e., ρ TI-LSTM) are consistently better than those of TI-LSTM without lexicon embedding (i.e., TI-LSTM) listed in Table IV and ρ TI-LSTM without flipping ($p < 0.05$). ρ TI-LSTM with flipping also slightly outperforms ρ TI-LSTM without flipping indicating that the flipping modular is useful in the network. The improved accuracies of ρ TI-LSTM over TI-LSTM on the three data corpora are listed in Table VI.

TABLE VI
ACCURACY IMPROVEMENT OF TWO-LEVEL LSTM WHEN LEXICON EMBEDDING WAS USED OVER THOSE OF TWO-LEVEL LSTM WITHOUT LEXICON EMBEDDING

Travel	Hotel	Mobile
+1.6%	+2.4%	+1.7%
+3.4%	+3.9%	+2.7%

TABLE VII
ACCURACIES OF THE COMPETING METHODS ON ENGLISH CORPORA IN ρ -HOT ENCODING

Method	MR	SST
CNN	0.815	0.469
Bi-LSTM	0.793	0.465
Tree-LSTM	0.807	0.481
NCSL	0.829	0.471
LR-Bi-LSTM	0.821	0.486
ρ TI-LSTM	0.851	0.498

B. Results on English Corpora

Two benchmark datasets are used for evaluating the proposed models: 1) movie review (MR) [41] and 2) Stanford sentiment treebank (SST) [42]. The former consists of 10 662 sentences with binary classes (positive and negative); while the latter consists of 11 885 sentences with five classes {very negative, negative, neutral, positive, very positive}.

The involved competing algorithms are CNN, LSTM/Bi-LSTM, Tree-LSTM [43], NCSL [7], and LR-Bi-LSTM [8]. The last two methods are two state-of-the-art methods that also utilize the word-level polar information. The results of these competing algorithms are obtained directly from the results presented in Qian *et al.*'s work [8]. The polar and negation words are compiled by following the method used in Qian *et al.*'s work [8]. The number of conjunction words in the dictionary is 207 which are attributed to 13 classes of words.

The proposed method achieves the highest accuracies among all the competing methods, including state-of-the-art NCSL and LR-Bi-LSTM, shown in Table VII. The increased accuracies on both sets are smaller than those on datasets reported in Section IV-A. The main reason lies in that the average numbers of clauses on both MR and SST are less than 1.5.

C. Discussion

The experimental evaluation discussed in Section IV-B verifies the effectiveness of the proposed method, ρ TI-LSTM, on fixed training/testing data split. In this section, we conducted more experiments via 10-cross-validation on each data corpus to further evaluate the proposed method. Moreover, unlike the conventional RNN, ρ TI-LSTM contains lexicon embedding that consists of new technique and components, including ρ -hot encoding, embedding for polar words, embedding for POS, and embedding for conjunctions. Therefore, this section evaluates the performances of the involved technique and embeddings separately.

TABLE VIII
CLASSIFICATION ACCURACIES OF COMPETING METHODS VIA TEN-FOLD CROSS-VALIDATION

Data corpus	Travel	Hotel	Mobile
CNN-Lex-C (R+C)	0.751	0.755	0.751
CNN-Lex-W (R+C)	0.766	0.769	0.763
Bi-LSTM-C (R+C)	0.768	0.765	0.801
Bi-LSTM-W (R+C)	0.751	0.778	0.803
ρ TI-LSTM-C	0.812	0.829	0.831
ρ TI-LSTM-W	0.801	0.815	0.841

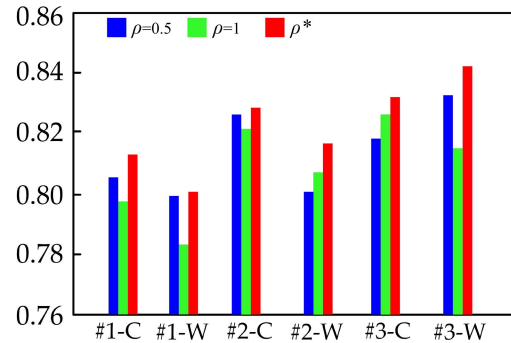


Fig. 8. Classification accuracies under different ρ values. #1-C and #1-W represent ρ TI-LSTM-C and ρ TI-LSTM-W on the first (travel) corpus, respectively; #2-C and #2-W represent ρ TI-LSTM-C and ρ TI-LSTM-W on the second (hotel) corpus, respectively; and #3-C and #3-W represent ρ TI-LSTM-C and ρ TI-LSTM-W on the third (hotel) corpus, respectively. ρ^* is the searched optimal value.

1) *Ten-Fold Cross-Validation Results on the Chinese Corpora:* The ten-fold cross-validation results for the main competing methods in Section IV-A2 are listed in Table VIII. Based on the *t*-test, ρ TI-LSTM significantly outperforms the involved competing methods ($p < 0.01$) on all the three datasets. Overall, the performances of the methods with character embedding are comparable to those of the methods with word embedding.

2) *Effect of Different Parameters on ρ -Hot Encoding:* Our ρ -hot encoding differs from one-hot encoding in two aspects. The first aspect is that the nonzero values in one-hot encoding are only equal to 1, whereas the nonzero values in ρ -hot encoding are ρ . The second aspect is that only one element in one-hot encoding is nonzero, whereas n elements in ρ -hot encoding are nonzero.

In this experiment, we test whether ρ -hot encoding is useful in two experimental runs. In the first run, the value of ρ is manually set to 0.5 and 1 in the experimental run without optimization. The parameter n in (6) is set as 15. The classification accuracies vary according to different ρ values on all the three data corpora. When ρ equals 1, the accuracies are the lowest in most cases shown in Fig. 8.

The results shown in Fig. 8 indicate that the value of ρ does affect the performance of the entire network. Consequently, the classical one-hot encoding, which fixes the value of nonzero elements as 1, is ineffective. In our experiments, the learned value of ρ is approximate 0.4.

In the second run, the performances under different n (i.e., 1, 5, 10, and 15) are tested. Table IX shows the comparison results. The value of n does affect the performance of

TABLE IX
ACCURACIES OF ρ TL-LSTM WITH DIFFERENT n
VALUES IN ρ -HOT ENCODING

	n	Travel	Hotel	Mobile
ρ Tl-LSTM-C	1	0.790	0.814	0.809
	5	0.803	0.821	0.815
	10	0.798	0.815	0.820
	searched	0.812	0.829	0.831
ρ Tl-LSTM-W	1	0.796	0.798	0.820
	5	0.802	0.809	0.829
	10	0.800	0.801	0.820
	searched	0.801	0.815	0.841

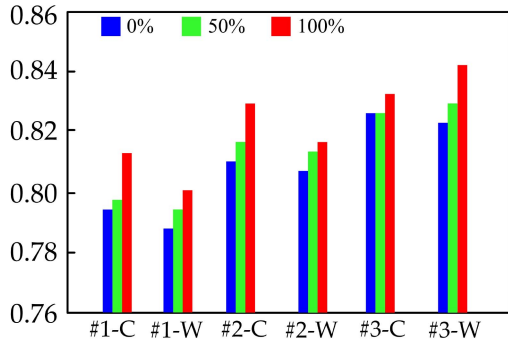


Fig. 9. Classification accuracies under different proportions of polar words. #1-C and #1-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the first (travel) corpus, respectively; #2-C and #2-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the second (hotel) corpus, respectively; and #3-C and #3-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the third (hotel) corpus, respectively.

the entire network, thereby indicating that the introduction of the n -duplicated strategy in encoding is effective. In the experiments, when n is increasing, the accuracies first increase and then decrease. The main reason may lie in the fact that when n becomes large, the proportion of lexicon embedding becomes large accordingly. An overlength input feature vector may incur “curse of dimensionality” and thus weaken the performance of the proposed two-level network.

3) *Effect of Polar Words*: In this experimental run, we test whether the labeled polar (negative and positive) words do affect the performance of the entire method when they are used in lexicon embedding. To this end, we order the polar words according to their frequencies in the training data. Top 0%, 50%, and 100% polar words are used. The corresponding classification accuracies are depicted in Fig. 9.

In most cases, the accuracies are the lowest when no polar words are used in the lexicon embedding. When all polar words are used, the proposed network achieves the highest accuracies.

In the experiment, only one user is invited to manually compile the dictionary for a data corpus. One and a half hour is needed for each data corpus. In our viewpoint, it is worth manually compiling the polar words for sentiment analysis by considering the performance improvement and time consumption.

4) *Effect of POS Cues*: In this experimental run, we test whether POS cues do play positive roles in the entire model.

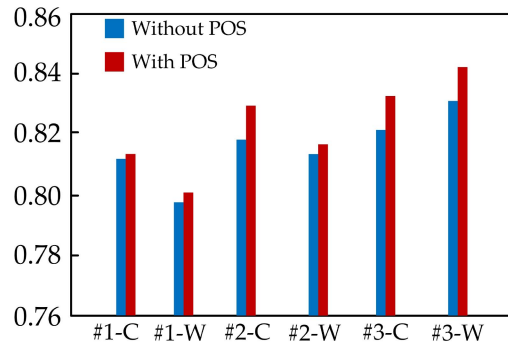


Fig. 10. Classification accuracies with and without POS in lexicon embedding. #1-C and #1-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the first (travel) corpus, respectively; #2-C and #2-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the second (hotel) corpus, respectively; and #3-C and #3-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the third (hotel) corpus, respectively.

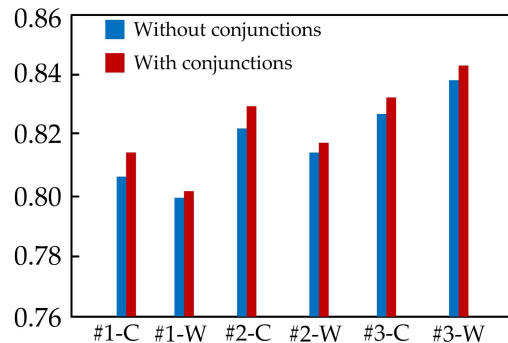


Fig. 11. Classification accuracies with and without conjunction in lexicon embedding. #1-C and #1-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the first (travel) corpus, respectively; #2-C and #2-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the second (hotel) corpus, respectively; and #3-C and #3-W represent ρ Tl-LSTM-C and ρ Tl-LSTM-W on the third (hotel) corpus, respectively.

To this end, we remove POS in the lexicon embedding of the proposed method. The results are shown in Fig. 10.

In most cases, the accuracies with POS embedding are greater than those without POS embedding, thereby indicating that the application of POS to lexicon embedding is useful.

5) *Effect of Conjunction Cues*: In this experimental run, we test whether conjunction cues do play positive roles in the entire model. To this end, the lexicon embedding for conjunction words is also removed from the proposed method. The results are shown in Fig. 11.

The algorithm with conjunction embedding outperforms that without conjunction embedding consistently, thereby indicating that the application of conjunction to lexicon embedding is useful.

6) *Effect of Different Training Tricks*: In this experimental run, we test the main competing methods under different training tricks, including dropout, batch normalization (BN), and pooling. Table X presents the accuracies of ρ Tl-LSTM based on 10-cross-validation. These two training tricks do not improve the performances. Table XI presents the accuracy comparison between average pooling and max pooling when CNN-lex is used. The average pooling significantly outperforms max pooling ($p < 0.01$).

TABLE X
CLASSIFICATION ACCURACIES WITH DIFFERENT TRAINING TRICKS

Data corpus	Travel	Hotel	Mobile
ρ Tl-LSTM-C+BN	0.800	0.817	0.830
ρ Tl-LSTM-W+BN	0.791	0.804	0.825
ρ Tl-LSTM-C+dropout	0.796	0.820	0.819
ρ Tl-LSTM-W+dropout	0.795	0.801	0.832
ρ Tl-LSTM-C+BN+dropout	0.793	0.808	0.823
ρ Tl-LSTM-W+BN+dropout	0.794	0.809	0.832

TABLE XI
CLASSIFICATION ACCURACIES WITH TWO POOLING STRATEGIES

Data corpus	Travel	Hotel	Mobile
CNN-Lex-C (R+C) with average	0.768	0.767	0.767
CNN-Lex-C (R+C) with max	0.751	0.755	0.751
CNN-Lex-W (R+C) with average	0.779	0.773	0.778
CNN-Lex-W (R+C) with max	0.766	0.769	0.763

7) *Summarization of the Experimental Comparisons:* According to the above experimental comparisons, our proposed methodology has the following advantages.

- 1) The two-stage labeling strategy can provide more supervised information, which is useful for model training.
- 2) The proposed ρ -hot encoding is more flexible than the one-hot encoding and is thus more useful for cues embedding.
- 3) The embedding of more lexical cues integrates more useful information related to sentiment orientations. The flipping module is also helpful in performance improvement.

V. CONCLUSION

High-quality labels are crucial for learning systems. Nevertheless, texts with mixed sentiments are difficult for humans to label in text sentiment classification. In this study, a new labeling strategy was introduced to partition texts into those with pure and mixed sentiment orientations. These two categories of texts were labeled using different processes. A two-level network was accordingly proposed to utilize the two labeled data in our two-stage labeling strategy. Lexical cues (e.g., polar words, POS, and conjunction words) are particularly useful in sentiment analysis. These lexical cues were used in our two-level network, and a new encoding strategy, that is, ρ -hot encoding, was introduced. ρ -hot encoding was motivated by one-hot encoding. However, the former alleviates the drawbacks of the latter. Due to labeling noise or context, the polarity of a word varied in different texts. A flipping model was proposed to model the polarity flipping process. Three Chinese sentiment text data corpora were compiled to verify the effectiveness of the proposed methodology. Our proposed method achieved the highest accuracies on these three data corpora. On English data corpora, the proposed method outperformed state-of-the-art algorithms.

The proposed two-level network and lexicon embedding can also be applied to other types of deep neural networks. In our future work, we will extend our main idea into several networks and text mining applications.

ACKNOWLEDGMENT

The experimental datasets and partial codes are publicly available at Github: <https://github.com/Tju-AI/two-stage-labeling-for-the-sentiment-orientations>.

REFERENCES

- [1] B. Liu, *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. Cambridge, U.K.: Cambridge Univ. Press, 2015.
- [2] A. Pak and P. Paroubek, "Twitter as a corpus for sentiment analysis and opinion mining," in *Proc. Int. Conf. Lang. Resources Eval. (LREC)*, 2010, pp. 17–23.
- [3] F. Li, S. Wang, S. Liu, and M. Zhang, "SUIT: A supervised user-item based topic model for sentiment analysis," in *Proc. AAAI Conf. (AAAI)*, 2014, pp. 1636–1642.
- [4] Y. Kim, "Convolutional neural networks for sentence classification," in *Proc. Int. Conf. Empirical Methods Nat. Lang. (EMNLP)*, 2014, pp. 1746–1751.
- [5] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-based methods for sentiment analysis," *Comput. Linguist.*, vol. 37, no. 2, pp. 267–307, 2011.
- [6] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [7] Z. Teng, D.-T. Vo, and Y. Zhang, "Context-sensitive lexicon features for neural sentiment analysis," in *Proc. EMNLP*, 2016, pp. 1629–1638.
- [8] Q. Qian, M. Huang, J. Lei, and X. Zhu, "Linguistically regularized LSTMs for sentiment classification," in *Proc. Annu. Meeting Assoc. Comput. Linguist. (ACL)*, 2017, pp. 1679–1689.
- [9] W. Zhao *et al.*, "Weakly-supervised deep embedding for product review sentiment analysis," *IEEE Trans. Knowl. Data Eng.*, vol. 30, no. 1, pp. 185–197, Sep. 2018.
- [10] M. Hu and B. Liu, "Mining and summarizing customer reviews," in *Proc. ACM SIGKDD Int. Conf. Knowl. Disc. Data Min. (KDD)*, 2004, pp. 168–177.
- [11] T. Mullen and N. Collier, "Sentiment analysis using support vector machines with diverse information sources," in *Proc. Int. Conf. Empirical Methods Nat. Lang. (EMNLP)*, 2004, pp. 412–418.
- [12] S. Liu, X. Cheng, F. Li, and F. Li, "TASC: Topic-adaptive sentiment classification on dynamic tweets," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 6, pp. 1696–1709, Dec. 2015.
- [13] G. Paltoglou and M. Thelwall, "A study of information retrieval weighting schemes for sentiment analysis," in *Proc. Annu. Meeting Assoc. Comput. Linguist. (ACL)*, 2010, pp. 1386–1395.
- [14] X. Glorot, A. Bordes, and Y. Bengio, "Domain adaptation for large-scale sentiment classification: A deep learning approach," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2011, pp. 513–520.
- [15] D. Tang, B. Qin, X. Feng, and T. Liu, "Effective LSTMs for target dependent sentiment classification," in *Proc. Int. Conf. Comput. Linguist. (COLING)*, 2016, pp. 3298–3307.
- [16] A. Esuli, S. Baccianella, and F. Sebastiani, "SENTI-WordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining," in *Proc. Int. Conf. Lang. Resources Eval. (LREC)*, 2010, pp. 2200–2204.
- [17] S. Lai, K. Liu, L. Xu, and J. Zhao, "How to generate a good word embedding," *IEEE Intell. Syst.*, vol. 31, no. 6, pp. 5–14, Jul. 2015.
- [18] K. Cho *et al.*, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," in *Proc. Int. Conf. Empirical Methods Nat. Lang. (EMNLP)*, 2014, pp. 1724–1734.
- [19] M.-C. Yang, N. Duan, M. Zhou, and H.-C. Rim, "Joint relational embeddings for knowledge-based question answering," in *Proc. Int. Conf. Empirical Methods Nat. Lang. (EMNLP)*, 2014, pp. 645–650.
- [20] A. M. Dai, C. Olah, and Q. V. Le, "Document embedding with paragraph vectors," 2015. [Online]. Available: arxiv.org/abs/1507.07998.
- [21] D. Tang, B. Qin, and T. Liu, "Document modeling with gated recurrent neural network for sentiment classification," in *Proc. Int. Conf. Empirical Methods Nat. Lang. (EMNLP)*, 2015, pp. 1422–1432.
- [22] N. Kalchbrenner, E. Grefenstette, and P. Blunsom, "A convolutional neural network for modelling sentences," in *Proc. Annu. Meeting Assoc. Comput. Linguist. (ACL)*, 2014, pp. 655–665.
- [23] R. Socher, C. C.-Y. Lin, A. Y. Ng, and C. D. Manning, "Parsing natural scenes and natural language with recursive neural networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2011, pp. 129–136.
- [24] M. Zhang, Y. Zhang, and D. Vo, "Gated neural networks for targeted sentiment analysis," in *Proc. AAAI Conf.*, 2016, pp. 3087–3093.

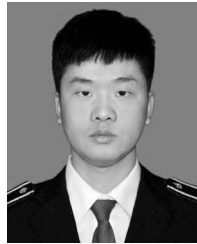
- [25] C. Santos, N. Dos, and M. Gattit, "Deep convolutional neural networks for sentiment analysis of short texts," in *Proc. Int. Conf. Comput. Linguist. (COLING)*, 2014, pp. 69–78.
- [26] K. Schouten, O. van der Weijde, F. Frasincar, and R. Dekker, "Supervised and unsupervised aspect category detection for sentiment analysis with co-occurrence data," *IEEE Trans. Cybern.*, vol. 48, no. 4, pp. 1263–1275, Apr. 2018.
- [27] F. A. Gers, J. Schmidhuber, and F. Cummins, "Learning to forget: Continual prediction with LSTM," *Neural Comput.*, vol. 12, no. 10, pp. 2451–2471, 2000.
- [28] X. Zhou, X. Wan, and J. Xiao, "Attention-based LSTM network for cross-lingual sentiment classification," in *Proc. Int. Conf. Empirical Meth. Nat. Lang. (EMNLP)*, 2016, pp. 247–256.
- [29] B. Shin, T. Lee, and J. D. Choi, "Lexicon integrated CNN models with attention for sentiment analysis," in *Proc. WAS-SA@EMNLP*, 2017, pp. 149–158.
- [30] K. Schouten and F. Frasincar, "Survey on aspect-level sentiment analysis," *IEEE Trans. Knowl. Data Eng.*, vol. 28, no. 3, pp. 813–830, Oct. 2016.
- [31] K. Schouten, O. van der Weijde, F. Frasincar, and R. Dekker, "Supervised and unsupervised aspect category detection for sentiment analysis with co-occurrence data," *IEEE Trans. Cybern.*, vol. 48, no. 4, pp. 1263–1275, Aug. 2018.
- [32] Y. Wang, M. Huang, X. Zhu, and L. Zhao, "Attention-based LSTM for aspect-level sentiment classification," in *Proc. EMNLP*, 2016, pp. 606–615.
- [33] L. He, K. Lee, M. Lewis, and L. Zettlemoyer, "Deep semantic role labeling: What works and what's next," in *Proc. Annu. Meeting Assoc. Comput. Linguist. (ACL)*, 2017, pp. 473–483.
- [34] B. Wang and W. Lu, "Learning latent opinions for aspect-level sentiment classification," in *Proc. AAAI Conf.*, 2018, pp. 5537–5544.
- [35] R. Narayanan, B. Liu, and A. Choudhary, "Sentiment analysis of conditional sentences," in *Proc. Int. Conf. Empirical Methods Nat. Lang. (EMNLP)*, 2009, pp. 180–189.
- [36] E. Cambria, S. Poria, A. F. Gelbukh, and M. Thelwall, "Sentiment analysis is a big suitcase," *IEEE Intell. Syst.*, vol. 32, no. 6, pp. 74–80, Nov./Dec. 2017.
- [37] X. Ding and B. Liu, "The utility of linguistic rules in opinion mining," in *Proc. 30th Annu. Int. ACM SIGIR Conf. (SIGIR)*, 2007, pp. 811–812.
- [38] Y. Duan and O. Wu, "Learning with auxiliary less-noisy labels," *IEEE Trans. Neural Netw. Learning Syst.*, vol. 28, no. 7, pp. 1716–1721, Apr. 2017.
- [39] E. Dobriban and S. Wager, "High-dimensional asymptotics of prediction: Ridge regression and classification," *Ann. Stat.*, vol. 46, no. 1, pp. 247–279, 2018.
- [40] W. Yin, K. Kann, M. Yu, and H. Schütze, "Comparative study of CNN and RNN for natural language processing," 2017. [Online]. Available: arxiv.org/abs/1702.01923.
- [41] B. Pang and L. Lee, "Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales," in *Proc. Annu. Meeting Assoc. Comput. Linguist. (ACL)*, 2005, pp. 115–124.
- [42] R. Socher *et al.*, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proc. EMNLP*, 2013, pp. 1631–1642.
- [43] K. S. Tai, R. Socher, and C. D. Manning, "Improved semantic representations from tree-structured long short-term memory networks," in *Proc. Annu. Meeting Assoc. Comput. Linguist. (ACL)*, 2015, pp. 1556–1566.



mining and machine learning.

Ou Wu received the B.Sc. degree in electrical engineering from Xi'an Jiaotong University, Xi'an, China, in 2003, and the M.Sc. and Ph.D. degrees in computer science from the National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences, Beijing, China, in 2006 and 2012, respectively.

In 2007, he joined NLPR as an Assistant Professor. In 2017, he joined the Center for Applied Mathematics, Tianjin University, Tianjin, China, as a Full Professor. His research interests include data



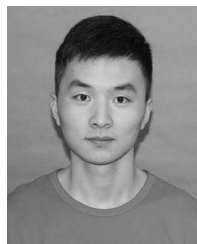
Tao Yang received the B.Eng. degree in electronic information engineering from the Civil Aviation University of China, Tianjin, China, in 2016, where he is currently pursuing the M.Eng. degree in control engineering.

He is currently a Research Intern with the Center for Applied Mathematics, Tianjin University, Tianjin. His research interests include text mining and deep learning.



Mengyang Li received the B.Eng. degree in automation from the Zhengzhou University of Aeronautics, Zhengzhou, China, in 2015. He is currently pursuing the M.Eng. degree in pattern recognition and intelligent control with the Electronic Information and Automation College, Civil Aviation University of China, Tianjin, China.

He is currently a Research Intern with the Center for Applied Mathematics, Tianjin University, Tianjin. His research interests include text mining and deep learning.



Ming Li received the B.Eng. degree in automation from the Zhengzhou University of Aeronautics, Zhengzhou, China, in 2015. He is currently pursuing the M.Eng. degree in control science and engineering with the Civil Aviation University of China, Tianjin, China.

He is currently a Research Intern with the Center for Applied Mathematics, Tianjin University, Tianjin. His research interests include text mining and deep learning.