# A CNN-Based Framework for Predicting Public Emotion and Multi-Level Behaviors Based on Network Public Opinion

Hangfeng Lin[1]* and Naiqing Bu[2]

[1] School of Political Science and Public Administration, East China University of Political Science and Law, Shanghai, China,
[2] School of Sociology, Sanya University, Sanya, China

Analysis of network public opinion can help to effectively predict the public emotion and the multi-level government behaviors. Due to the massive and multidimensional characteristics of network public opinion data, the in-depth value mining of public opinion is one of the research bottlenecks. Based on Term Frequency-Inverse Document Frequency (TF-IDF) and deep learning technologies, this paper proposes an advanced TF-IDF mechanism, namely TF-IDF-COR, to extract text feature representations of public opinions and develops a CNN-based prediction model to predict the tendency of publics' emotion and mental health. The proposed method can accurately judge the emotional tendency of network users. The main contribution of this paper is as follows: (1) based on the advantages of TF-IDF mechanism, we propose a TF-IDF-COR mechanism, which integrates the correlation coefficient of word embeddings to TF-IDF. (2) To make the extracted feature semantic information more comprehensive, CNN and TF-IDF-COR are combined to form an effective COR-CNN model for emotion and mental health prediction. Finally, experiments on Sina-Weibo and Twitter opinion data sets show that the improved TF-IDF-COR and the COR-CNN model have better classification performance than traditional classification models. In the experiment, we compare the proposed COR-CNN with support vector machine, k-nearest neighbors, and convolutional neural network in terms of accuracy and F1 score. Experiment results show that COR-CNN performs much better than the three baseline models.

Keywords: network public opinion analysis, mental health of netizens, emotional tendency prediction, convolutional neural network, TF-IDF, multi-level government behaviors

## INTRODUCTION

In the new media era, publics tend to express their views on social events through the Internet. Yong people are extremely active in the network. Their enthusiasm and desire to participate in social life and the expression of subject consciousness are becoming stronger and stronger. They are used to overseeing the "voice box" of the network and have the appeal of expressing their own opinions and interests on social, cultural, economic, and other issues. They take the new media as a way and means to pursue their own value and express their individual consciousness. It is necessary to do intelligent analysis of the content, characteristics, and problems of young people's online public opinion in the new media era. We can improve the collection channels of young people's online public opinion information to achieve this.

Due to the massive and multidimensional characteristics of network public opinion data, the in-depth value mining of public opinion data has always been one of the bottlenecks research areas. In recent years, the rise and practicability of artificial intelligence technology has provided a new means and path for us to realize the automation, intelligence, and accuracy of network public opinion analysis. Therefore, some researchers have also made useful explorations. For example, some researchers used wavelet analysis to decompose the development process of public opinion (Gao et al., 2021), and then applied artificial neural network to model and predict the trend of public opinion (Punjabi et al., 2019). In addition, neural network simulation is adopted to simulate the development process of public opinion (Jelodar et al., 2020). Gray prediction and pattern recognition are employed to predict the trend of public opinion (Eagly et al., 2020).

The analysis and prediction of public opinion data is the most critical step in the network public opinion analysis technology. By designing appropriate algorithms to analyze the public opinion data, we can explore the hot topics, evaluate their communication impact and public opinion level, and adopt reasonable means to guide and control public opinion. In terms of public opinion analysis, the commonly used technical means include Bayesian classifier, support vector machine (SVM), random depth, and neural networks (De Gelder, 2010). Public opinion prediction is an important step after data analysis. This process is mainly to provide important reference for public opinion monitoring and public opinion early warning and help formulate relevant response measures for government departments and enterprises at all levels. At present, the research on network public opinion prediction is relatively extensive, and there are many methods to use (Eagly et al., 2020). Public opinion information mostly comes from short text comment information. Its text is separated from the written language, the structure becomes more concise and lack of standardization, which often makes it difficult to extract text features. Traditional emotion analysis methods often rely on emotion dictionary and feature extraction. With the continuous updating and iteration of Internet culture and data volume, many manuals updating of emotion dictionary is required, otherwise semantic features will be lost, and classification will be inaccurate.

Based on the analysis of TF-IDF (Term Frequency-Inverse Document Frequency) and deep learning technologies, this paper proposes a TF-IDF-COR (Term Frequency-Inverse Document Frequency-Correlation) mechanism to extract text representations of public opinions, and a Convolutional neural network (CNN)-based prediction model (Pavlova, 2017) to predict the risk of publics' mental health. The proposed method can accurately judge the emotional tendency of network users. The main work of this paper are as follows: (1) When there are TF (word frequency) and IDF (inverse document frequency), we multiply the two words to get the TF-IDF value of a word. The larger the TF-IDF of a word in the article, the higher the importance of the word in the article. Therefore, by calculating the TF-IDF of each word in the article, the top words are the keywords of the article. Based on the advantages of TF-IDF mechanism, we propose a TF-IDF-COR mechanism,

which integrates the coefficients of word embeddings to TF-IDF. (2) CNN can better extract the local features of the text. To make the extracted feature semantic information more comprehensive, the two are combined to form a COR-CNN model. The COR-CNN model with the optimal parameters is obtained by comparing several groups of model parameters, which improves the classification performance compared with the traditional CNN models. Finally, experiments on Sina-Weibo and Twitter data set (Rodríguez et al., 2020) show that the improved TF-IDF-COR and the COR-CNN model has better classification performance than traditional classification models. In the experiment, we compare the proposed COR-CNN with SVM, KNN (K-Nearest Neighbor) and CNN models in terms of accuracy, recall and F1 score (Kobylińska and Kusev, 2019). Experiment results show that COR-CNN performs much better than the three baseline models. The main contribution of this work is as follows:

(1) We propose a new feature extraction and representation mechanism, called TF-IDF-COR. It integrates the coefficients of word embeddings to TF-IDF to find top words of the keywords of an article more efficiently.
(2) We propose a CNN-based model, called COR-CNN to extract feature semantic information by considering the coefficients of word embeddings.
(3) We conduct comprehensive experiments based on Sina-Weibo and Twitter data set for evaluate the proposed TF-IDF-COR and COR-CNN. They are compared with three classical machine learning algorithms. The experiment proves that the proposed methods are effective and efficient.

The structure of the remaining paper is: Section Related work introduces the related work about the machine learning algorithms for risk analysis of network public opinions and the machine learning algorithms for mental health evaluation and prediction. Section Text Representation Learning and Risk Prediction of Network Public Opinions introduces the proposed CNN-based framework for emotion and multi-level behavior prediction based on the network public opinions. Section Experiment Setting and Result Analysis shows the experiment design and result comparison with baseline models. Section Conclusion presents the paper conclusion and the future work.

## RELATED WORK

Lorenz-Spreen et al. (2020) proposed two classes of behavioral interventions-nudging and boosting-that enlist these cues to redesign online public opinion to promote deliberate cognition and autonomous choice. Pickett (2019) reviewed evidence for the effects of public opinion on court decision-making, capital punishment policy and use, correctional expenditures, and incarceration rates. D'Andrea et al. (2019) automatically inferred trends in the public opinion regarding the stance toward the vaccination topic. It enables the detection of significant opinion shifts. These opinion shifts can be possibly explained with the occurrence of specific social context-related events. This study (Han et al., 2020) uses random forest algorithm explored public

opinion in the early stages of COVID-19 in China by analyzing Sina-Weibo texts in terms of space, time, and content. In controlling the crisis, accurate response countermeasures should be formulated following public help demands. The study (Jia and Chen, 2020) applies emotional analysis is to the evolution analysis of network public opinion, and the change of network public opinion characteristics with time series is obtained, which can get the change of emotional characteristics of public opinion participants with time series.

The research work (Srividya et al., 2018) proposes to apply various machine learning algorithms such as support vector machines, decision trees, naïve bayes classifier, K-nearest neighbor classifier and logistic regression to identify state of mental health in a target group. Shatte et al. (2019) point out there is significant room for the application of machine learning to other areas of psychology and mental health. The challenges of using machine learning techniques are discussed. The opportunities to improve and advance the machine learning algorithms for text analysis is also analyzed. The work (Wilkinson et al., 2017) is a study of postnatal depression in women using a decision tree model, postnatal depression screening and treatment is a cost-effective intervention that should be considered as part of routine postnatal care. The work (Alonso et al., 2018) applies data mining techniques to mental health disorders such as dementia, schizophrenia, and depression. Such techniques can be of great help in clinical decision-making, diagnostic prediction and improving the quality of life of patients. This paper (Garcia-Ceja et al., 2018) surveys recent research works in mental health monitoring systems (MHMS) using sensor data and machine learning.

# TEXT REPRESENTATION LEARNING AND RISK PREDICTION OF NETWORK PUBLIC OPINIONS

In this section we propose a CNN-based framework, called COR-CNN for network public opinions classification. At first, we need to extract text summaries from public opinions and news. Then we use the external corpus to train the word vector model, convert the text representations extracted from the opinions into a vectorized representation based on the word vector, and then concatenate the vectorized representation of the sentence into a vectorized representation of the entire text. Finally, the convolutional neural network is trained, and the trained network model is used to predict the emotion of netizens.

## Learning Text Representations in Public Opinions

Compared with the data in the form of online comments, news texts are generally longer. And the difference between the lengths of the texts is also large, so the texts must be processed before use. Sentences can more accurately grasp the semantics of text than words, so we solve this problem by extracting text summaries of texts. For text classification, the traditional TF-IDF algorithm will lose some key classification criteria. Therefore, we will design an enhanced TF-IDF algorithm to solve this problem. Then, the

top K sentences with the highest scores are selected as the text summary of the text, which reduces the data dimension and eliminates noises.

## TF-IDF-COR Algorithm

TF-IDF considers the text set as a whole, and its IDF part does not consider the inter-class distribution information of feature items. If the entry $t_a$ has a high frequency of occurrence in a certain category $c_a$, which leads to the occurrence of the entry $t_a$ in more texts. Although the entry $t_a$ appears less in other categories, the weight calculated according to the IDF algorithm will be too small. The entry $t_a$ will be mistaken for an entry with poor ability to distinguish between categories. Obviously, this is not in line with the actual situation. Words that only appear frequently in a certain category or categories are the most one or several categories of iconic words have high information value for text classification, so they should be given high weights. Correspondingly, if the entry $t_a$ appears only in a small amount of text. The frequency of occurrence in each class is relatively uniform. This kind of unimportant word will be given very high weight by the IDF algorithm, so it is not in line with the actual situation. When the distribution of a word under different topic news is quite different, it means that the word has a strong representativeness for a certain category or categories of news and can be used as a key basis for classification. The above-mentioned defects are because the traditional TF-IDF algorithm does not consider the distribution information of words between classes, which obviously loses part of the accuracy when performing text classification. Therefore, we measure the inter-class distribution information of words by adding an inter-class correlation coefficient (Wang et al., 2012).

## Learn Text Representations

For the task of text representation extraction, the goal is to extract the most important set of sentences in the text. First, the weights of the terms are calculated by an improved TF-IDF algorithm. Suppose D is a text collection. For any word $w_a$ in the text $D_b$, the word frequency is expressed by Equation 1.

$$tf_{ab} = \frac{n_{ab}}{\sum_k n_{kb}} \quad (1)$$

where $n_{ab}$ represents the number of occurrences of the word $w_a$ in the text $D_b$, and m represents the length of the dictionary.

The inverse document frequency is expressed by Equation 2.

$$idf_a = \log \frac{|D|}{1 + |\{D_b : t_a \in D_b\}|} \quad (2)$$

The inter-class dispersion of words is expressed by Equation 3.

$$ICD_a = \frac{\sqrt{\frac{1}{c-1} \sum_{r=1}^{c} \left( tf_{ak}(w_{ar}) - \overline{tf_{ak}(w_{ar})} \right)^2}}{tf_{ar}(w_{ar})} \quad (3)$$

In the formula: $\overline{tf_{ar}} = \frac{1}{c} \sum_{a=1}^{n} tf_{ar}(w_{ar})$, where $tf_{ar}$ represents the number of occurrences of the word $w_a$ in the L class, and c represents the number of text classes.

The product of Equations 1–3 is the TF-IDF-COR value of the corresponding term.

After obtaining the weights of all terms, the importance of the sentence is represented by the accumulation of the weights of the terms in the sentence. In addition, considering the different lengths of sentences, the results need to be processed accordingly to prevent the selected results from being biased toward long sentences. For a given sentence T, it can be represented by the terms contained in T: $T = (o_1, o_2 \cdots o_n)$, the importance of T is defined in Equation 4.

$$I(T) = \frac{\sum_{a=1}^{m} tf_{ab}^* idf_a^* ICD_a}{\log(|T|+1)} \tag{4}$$

where $tf_{ab}$ represents the terms of occurrences of the word $w_a$ in the sentence $T_b$, and $ICD_a$ represents the inter-class dispersion of words represents the length of the dictionary (see Equation 3).

Finally, to obtain text representations, all sentences of the text are ranked. The top K sentences for scoring are selected as the representation of the current text.

## Determine the Embedding Distribution of Words in Public Opinions

To describe the semantic embedding distribution of words in each paper, an abstract distributed convolutional layer is designed. The text representation I(T) of the words in the sub-word is used as input. A cross entropy function is used to generate the embedding distributed representation of each word. Formally, using semantic information provided by $I_a$, the expression order of abstract semantic field distribution of sub-words, $K = \{k_1, k_2, \ldots, k_N\}$, can be calculated by Equations 5–7.

$$g_a = W_r I_a + b_r \tag{5}$$

$$r_a = softmax(g_a) \tag{6}$$

$$r_{a,b} = \frac{e^{m_{ab}}}{\sum_{k=1}^{d_2} e^{m_{i,k}}} \tag{7}$$

where, $W_r$ and $b_r$ are parameters, $W_r \in \mathbb{R}^{d_3 \times d_4}, m_i, b_r \in \mathbb{R}^{d_3}$.

There is naturally a gap inconsistency between sematic words and non-semantic words. Moreover, the greater the inconsistency between the semantic words and the semantic area of the text, the greater the possibility that the sematic words will be determined. Using the total distributed distance between semantic words and non-semantic words in the sub-word, the semantic area of the sub-word is not consistent. That is, the minimum area word embedding distribution distance between each semantic word and non-semantic words. The semantic area inconsistency of sematic embeddings is relatively high, and the goal is to maximize the semantic area inconsistency of sub-word. Choosing the minimum area distance can ensure that the local inconsistency between each pair of semantic words and non-semantic words is modeled correctly. The semantic field of the sentence T semantic area inconsistency SAI(T) is calculated by Equation 8.

$$SAI(T) = \sum_{j=1}^{J} \min_{i \geq u \geq U} \{Dista(r_j^{(m)}, r_u^{(u)})\} \tag{8}$$

where, $r_k^{(met)}$, $r_u^{(un)}$ represent the semantic area of the embedding word $t_i$ and the non-metaphorical word $t_u$ in the sentence, respectively. $Dis(r_j^{(m)}, r_u^{(u)})$ expresses distance function of $r_j^{(m)}$ and $r_u^{(u)}$. More than one distance measures exist in real life.

Distance 1: Kullback-Leibler dispersion distance is calculated by Equation 9.

$$D(r_{ki}, r_{uj}) = \sum_{m=1}^{d_2} r_{k_{im}}^* \log(\frac{r_{i_m}}{r_{j_m}}) \tag{9}$$

Distance 2: Euclidean distance is calculated by Equation 10.

$$D(r_{ki}, r_{uj}) = \sqrt{\sum_{m=1}^{d_2} (r_{ki_m} - r_{j_m})^2} \tag{10}$$

Distance 3: Korrigierter Kosinusabstand is calculated by Equation 11.

$$D(r_i, r_j) = -\frac{\sum_{m=1}^{d_2} r_{i_m} r_{j_m}}{\sqrt{\sum_{m=1}^{d_2} (r_{i_m})^2} \sqrt{\sum_{m=1}^{d_2} (r_{j_m})^2}} \tag{11}$$

Distance 4: Gaussian distance is calculated by Equation 12.

$$D(r_{ki}, r_{uj}) = \exp(-\frac{\sqrt{\sum_{m=1}^{d_2} (r_{ki_m} - r_{j_m})^2}}{2\sigma^2}) \tag{12}$$

## Tensorization Words

Tensorization words refers to the transformation of words in a language into digital representations that are easy for computers to process. The NNLM and the Log-Linear model are both outstanding representatives of using neural networks to obtain word vectors. The well-known Word2vec model is borrowed from these two and is a more concise and efficient word vector model. Word2vec technology is a key breakthrough in the application of deep learning technology in autologous language processing.

### Word2vec and Skip Gram

Word2vec is a group of related models used to generate word vectors. These models are shallow and double-layer neural networks used for training to reconstruct linguistic word texts. The network is represented by words. The input words in adjacent positions need to be guessed. Under the assumption of word bag model in word2vec, the order of words is not important. After training, word2vec model can be used to map each word to a vector, which can be used to represent the relationship between words. This vector is the hidden layer of neural network.

Skip gram is a simple but very practical model. In natural language processing, the selection of corpus is a very important problem. At first, the corpus must be sufficient. On the one hand, the number of words in the dictionary should be large enough. On the other hand, it should include as many sentences reflecting the relationship between words as possible. For

example, only when the sentence pattern of "fish swimming in the water" is as many as possible in the corpus, can the model learn the semantic and grammatical relationship in the sentence. This is the same as human learning natural language. If people repeat more times, they will imitate. Second, the corpus must be accurate. In other words, the selected corpus can correctly reflect the semantic and grammatical relationship of the language, which seems not difficult to achieve. For example, in Chinese, the corpus of people's daily is more accurate. However, more often, it is not the selection of corpus that raises concerns about accuracy, but the method of processing. In the n-ary model, due to the limitation of window size, the relationship between words beyond the scope of the window and the current word cannot be correctly reflected in the model. If the window size is simply expanded, it will increase the complexity of training. The skip gram model solves these problems.

After extracting the text representations from the text, the sentences are tensorized based on the word embedding vector model. The word2vec and corpus tool are used to train the model. For a sentence $T = (t_1, t_2 \cdots t_n)$, for any term $M_i = (m_1, m_2 \cdots m_k)$, where k represents the dimension of the word. The tensorized representation of statement T is $T = (t_1 \oplus t_2 \oplus \cdots \oplus t_n)$, where $\oplus$ is the concatenation operator. Therefore, sentence T is tensorized to a series of ordered vectors. Similarly, for each text $P = (T_1 \oplus T_2 \oplus \cdots \oplus T_l)$, where L represents the order of importance of sentence T. That is to say, $T$ is the most important sentence of the text P. After converting the text into vector form, the vectorized data can be used to train the neural network.

### A CNN Model for Word Semantic Analysis in Public Opinions

Convolutional neural networks are composed of convolutional layers, pooling layers, and fully connected layers. The convolutional layer extracts the features of the data through convolutional computation. The pooling layer selects the optimal features from the features provided by the convolutional layer, and then outputs them to the fully connected layer for processing. The input of the input layer is a matrix representing the text, k represents the dimension of the word vector, and m represents the number of word vectors contained in each data. The convolution layer involves the convolution kernel $k \in R^{lj}$, l represents the size of the convolution window, and j is the convolution dimension, which is equal to the dimension of the word vector. In general, $T_{i:i+h}$ represents the word $T_i, T_{i+1} \cdots T_{i+h}$. The expression to generate a text feature is $F = f(t \bullet T_{i:i+h+a})$, where a is the bias and f is a non-linear function. Applying this convolution kernel to $(T_{1:h}, T_{2:h+1} \cdots T_{N-h+1:N})$ generates a feature map $p = (p_1, p_2 \cdots p_{N-h+1})$. The pooling layer uses the maximum pooling method to sample the feature map, and only retains the most important features of each feature group: $p_{max} = p_i$. The multiple feature vectors output by the pooling layer are spliced and input to the input of the fully connected layer. The loss function of this COR-CNN model is defined in

**TABLE 1 |** Performance comparison of different classification models.

| Classification model | Precision | Recall | F1 |
| --- | --- | --- | --- |
| KNN | 0.832 | 0.798 | 0.827 |
| SVM | 0.897 | 0.875 | 0.882 |
| NB | 0.8333 | 0.7838 | 0.8052 |
| LSTM | 0.9221 | 0.8986 | 0.9147 |
| TF-IDF-COR | 0.957 | 0.949 | 0.953 |

**TABLE 2 |** Influence of different representation learning mechanism on risk classification results.

| | Precision | | Recall | | F1 | |
| --- | --- | --- | --- | --- | --- | --- |
| | TF-IDF | TF-IDF-COR | TF-IDF | TF-IDF-COR | TF-IDF | TF-IDF-COR |
| Sports | 0.841 | 0.903 | 0.884 | 0.947 | 0.857 | 0.928 |
| Health | 0.910 | 0.956 | 0.924 | 0.987 | 0.911 | 0.978 |
| Travel | 0.876 | 0.939 | 0.839 | 0.926 | 0.845 | 0.935 |
| Military | 0.912 | 0.978 | 0.933 | 0.945 | 0.925 | 0.959 |
| Education | 0.903 | 0.929 | 0.892 | 0.917 | 0.884 | 0.902 |
| Economy | 0.942 | 0.983 | 0.926 | 0.976 | 0.933 | 0.980 |

Equation 13.

$$Loss = \frac{1}{L} \sum_{k=1}^{L} \sum_{i=1}^{N_k} \gamma_{y_i} (1 - \hat{y}_i)^\delta \log(\hat{y}_i) - \mu H(T_L) \quad (13)$$

where,

$$H(T) = \begin{cases} SAI(S_l) & S_l \text{ is semantic} \\ -SAI(S_l) & S_l \text{ is Non semantic} \end{cases} \quad (14)$$

## EXPERIMENT SETTING AND RESULT ANALYSIS

### Experiment Settings

The data used in this experiment comes from the Twitter public opinion corpus (Pak and Paroubek, 2010) and Sina public opinion corpus (Xue et al., 2014). Some opinion data are selected, covering six categories of education, economy, health, military, tourism, and sports. For data preprocessing, a word segmentation tool is used to segment the dataset and remove stop words. After that, the text representations are extracted by using the TF-IDF-COR algorithm. If the text and sentences are too long or too short, they are truncated to a fixed length or filled with blank data. The word vector adopts the skip-gram model in word2vec and is trained using the Wikipedia Chinese and English corpus.

Experiments are designed to evaluate the performance of the sentence-level news classification scheme based on the convolutional neural network model. To evaluate the effect of the classification scheme, the accuracy, recall, and F1 score are used as matrices to measure the performance of the CNN model. To reflect the superiority of the model and verify the model, a group of experiments using the SVM model for text classification was
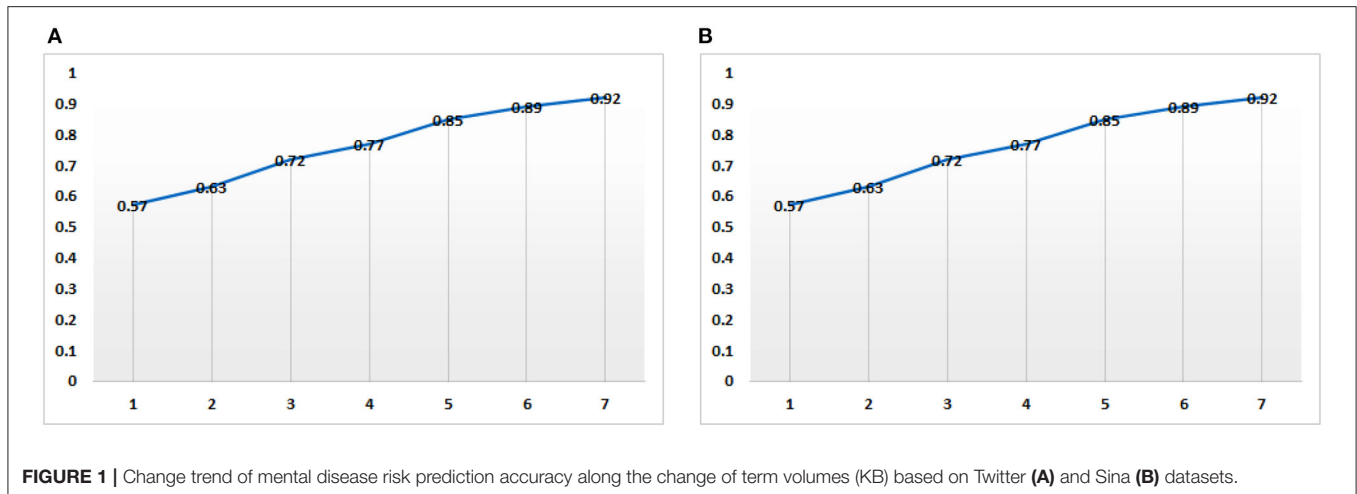
**FIGURE 1 |** Change trend of mental disease risk prediction accuracy along the change of term volumes (KB) based on Twitter **(A)** and Sina **(B)** datasets.

set as a baseline group. In addition, to illustrate the advantages of using the TF-IDF-COR algorithm in the information extraction stage of this mechanism, a set of comparative experiments using the TF-IDF-COR are set up on the same data set.

## Experiment Result

**Table 1** is the comparison results of the proposed mechanism with two traditional classification models KNN, SVM, Naïve Bayes (NB) and Long Short-tem Memory network (LSTM) in terms of precision (Precision), recall (Recall), and F1 score.

As we can see from **Table 1**, the opinion classification algorithm based on the proposed TF-IDF-COR and CNN model is more accurate and more stable than traditional KNN, SVM, NB, and LSTM. On the one hand, the convolutional neural network model can extract richer classification features by increasing the convolution kernel. On the other hand, it can also extract higher-level classification features by increasing the number of convolutional layers. To put it simply, convolutional neural networks can extract richer features horizontally and more levels of features vertically, which is incomparable to traditional machine learning models.

**Table 2** shows the comparison results of the precision rate (Precision), recall rate (Recall), and F1 score of the CNN model by using the traditional TF-IDF and the TF-IDF-COR, respectively, when learning text representations.

From **Table 2**, we can see that compared with the traditional TF-IDF algorithm, the average accuracy of the TF-IDF-COR is higher. The accuracy of the categories is more balanced. The higher accuracy in the original mechanism is only slightly reduced in the new mechanism, and the lower accuracy in the original mechanism is greatly improved in the new mechanism. Therefore, by improving the TF-IDF algorithm to capture different distributions of word embeddings between classes, the performance of CNN in text classification can be further improved.

As can be seen from **Figure 1**, with the change of term volumes (from 1 to 6 k terms), the risk prediction accuracy on both Twitter and Sina datasets are gradually increasing. With the accumulation of the public opinion data, the effect of predicting

the public mental diseases based on the public's opinions will become better and better. On the other hand, it is also reflected that with the increase of the public's time online, the potential occurrence of public mental diseases in the future will become more and more clear.

## CONCLUSION

This paper proposed an advanced TF-IDF-COR mechanism to extract text representations of public opinions, and a CNN-based prediction model to predict the risk of publics' mental health. The proposed method can accurately judge the emotional tendency of network users and government behaviors. The proposed TF-IDF-COR mechanism integrates the correlation coefficients of word embeddings to TF-IDF. CNN and TF-IDF-COR are combined to form a COR-CNN model. Finally, experiments on Sina-Weibo and Twitter data sets prove that the improved TF-IDF-COR and the COR-CNN model have better classification performance than traditional classification models. In the experiment, we compare the proposed COR-CNN with SVM, KNN and CNN models in terms of accuracy and F1 score. Experiment results show that COR-CNN performs much better than the three baseline models. As we used a CNN model to extract features, which is time consuming. The time and space efficiency of the proposed methods should be future improved. In the future, we will take more lightweight models like ResNet to improve the CNN-based feature extraction model to improve its time and space efficiency.

## DATA AVAILABILITY STATEMENT

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

## AUTHOR CONTRIBUTIONS

HL designed the framework, collected data, and did the experiment. NB proof read the paper. Both authors contributed to the article and approved the submitted version.

# REFERENCES

Alonso, S. G., de la Torre-Díez, I., Hamrioui, S., López-Coronado, M., Barreno, D. C., Nozaleda, L. M., et al. (2018). Data mining algorithms and techniques in mental health: a systematic review. *J. Med. Syst.* 42, 1–15. doi: 10.1007/s10916-018-1018-2

D'Andrea, E., Ducange, P., Bechini, A., Renda, A., and Marcelloni, F. (2019). Monitoring the public opinion about the vaccination topic from tweets analysis. *Exp. Syst. Applic.* 116, 209–226. doi: 10.1016/j.eswa.2018.09.009

De Gelder, B. (2010). The grand challenge for frontiers in emotion science. *Front. Psychol.* 1, 187. doi: 10.3389/fpsyg.2010.00187

Eagly, A. H., Nater, C., Miller, D. I., Kaufmann, M., and Sczesny, S. (2020). Gender stereotypes have changed: a cross-temporal meta-analysis of US public opinion polls from 1946 to 2018. *Am. Psychol.* 75, 301. doi: 10.1037/amp0000494

Gao, R., Du, L., Duru, O., and Yuen, K. F. (2021). Time series forecasting based on echo state network and empirical wavelet transformation. *Appl. Soft Comput.* 102, 107111. doi: 10.1016/j.asoc.2021.107111

Garcia-Ceja, E., Riegler, M., Nordgreen, T., Jakobsen, P., Oedegaard, K. J., and Tørresen, J. (2018). Mental health monitoring with multimodal sensing and machine learning: a survey. *Pervas. Mobile Comput.* 51, 1–26. doi: 10.1016/j.pmcj.2018.09.003

Han, X., Wang, J., Zhang, M., and Wang, X. (2020). Using social media to mine and analyze public opinion related to COVID-19 in China. *Int. J. Environ. Res. Public Health* 17, 2788. doi: 10.3390/ijerph17082788

Jelodar, H., Wang, Y., Orji, R., and Huang, S. (2020). Deep sentiment classification and topic discovery on novel coronavirus or COVID-19 online discussions: NLP using LSTM recurrent neural network approach. *IEEE J. Biomed. Health Inform.* 24, 2733–2742. doi: 10.1109/JBHI.2020.3001216

Jia, F., and Chen, C. C. (2020). Emotional characteristics and time series analysis of Internet public opinion participants based on emotional feature words. *Int. J. Adv. Robot. Syst.* 17, 1729881420904213. doi: 10.1177/1729881420904213

Kobylińska, D., and Kusev, P. (2019). Flexible emotion regulation: how situational demands and individual differences influence the effectiveness of regulatory strategies. *Front. Psychol.* 10, 72. doi: 10.3389/fpsyg.2019.00072

Lorenz-Spreen, P., Lewandowsky, S., Sunstein, C. R., and Hertwig, R. (2020). How behavioural sciences can promote truth, autonomy and democratic discourse online. *Nat. Hum. Behav.* 4, 1102–1109. doi: 10.1038/s41562-020-0889-7

Pak, A., and Paroubek, P. (2010). "Twitter as a corpus for sentiment analysis and opinion mining," in *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)* (Valletta).

Pavlova, M. A. (2017). Emotion science in the twenty-first century. Time, sex, and behavior in emotion science: over and above. *Front. Psychol.* 8, 1211. doi: 10.3389/fpsyg.2017.01211

Pickett, J. T. (2019). Public opinion and criminal justice policy: theory and research. *Ann. Rev. Criminol.* 2, 405–428. doi: 10.1146/annurev-criminol-011518-024826

Punjabi, V. D., More, S., Patil, D., Bafna, V., Shah, S., and Bachhav, H. (2019). A survey on trend analysis on Twitter for predicting public opinion on ongoing events. *Int. J. Comput. Applic.* 180, 13–17. doi: 10.5120/ijca2018916596

Rodríguez, C. P., Carballido, B. V., Redondo-Sama, G., Guo, M., Ramis, M., and Flecha, R. (2020). False news around COVID-19 circulated less on sina weibo than on twitter. how to overcome false information? *Int. Multidiscipl. J. Soc. Sci.* 9, 107–128. doi: 10.17583/rimcis.2020.5386

Shatte, A. B. R., Hutchinson, D. M., and Teague, S. J. (2019). Machine learning in mental health: a scoping review of methods and applications. *Psychol. Med.* 49, 1426–1448. doi: 10.1017/S0033291719000151

Srividya, M., Mohanavalli, S., and Bhalaji, N. (2018). Behavioral modeling for mental health using machine learning algorithms. *J. Med. Syst.* 42, 1–12. doi: 10.1007/s10916-018-0934-5

Wang, Z., Li, T., Xiong, N., and Pan, Y. (2012). A novel dynamic network data replication scheme based on historical access record and proactive deletion. *J. Supercomput.* 62, 227–250. doi: 10.1007/s11227-011-0708-z

Wilkinson, A., Anderson, S., and Wheeler, S. B. (2017). Screening for and treating postpartum depression and psychosis: a cost-effectiveness analysis. *Matern. Child Health J.* 21, 903–914. doi: 10.1007/s10995-016-2192-9

Xue, B., Fu, C., and Shaobin, Z. (2014). "A study on sentiment computing and classification of sina weibo with word2vec," in *2014 IEEE International Congress on Big Data* (Anchorage, AK: IEEE), 358–363. doi: 10.1109/BigData.Congress.2014.59