

Convolutional Gated Recurrent Units for Medical Relation Classification

Bin He^a, Yi Guan^{a,*}, Rui Dai^b

^aResearch Center of Language Technology, Harbin Institute of Technology, Harbin, China

^bDepartment of Mathematics, Harbin Institute of Technology, Harbin, China

Abstract

Convolutional neural network (CNN) and recurrent neural network (RNN) models have become the mainstream methods for relation classification. We propose a unified architecture, which exploits the advantages of CNN and RNN simultaneously, to identify medical relations in clinical records, with only word embedding features. Our model learns phrase-level features through a CNN layer, and these feature representations are directly fed into a bidirectional gated recurrent unit (GRU) layer to capture long-term feature dependencies. We evaluate our model on two clinical datasets, and experiments demonstrate that our model performs significantly better than previous single-model methods on both datasets.

Keywords: Relation classification; Clinical record; Convolutional neural network; Gated recurrent unit.

1. Introduction

Relation classification, a natural language processing (NLP) task which identifies the relation between two entities in a sentence, is an important technique in many subsequent NLP applications, such as question answering and knowledge base completion. In the clinical domain, Informatics for Integrating Biology and the Bedside (i2b2) released an annotated relation dataset on clinical records and attracted considerable attention [1]. Identifying relations in clinical records is a challenging task because one sentence from clinical records may contain more than two medical concepts and a concept may contain several words. Figure 1 illustrates relation samples in this task.

Due to the powerful feature learning ability, convolutional neural network (CNN) and recurrent neural network (RNN) are the mainstream architectures in the relation classification task [2–9]. In order to utilize the advantages of these two neural networks simultaneously, combinations of CNN and RNN turn into a research trend. The most direct way is to use the voting

Sentence
<i>Pain control</i> was initiated with <i>morphine</i> but was then changed to <i>demerol</i> , which gave the patient better relief of <i>his epigastric pain</i> .
Relations
(<i>pain control</i> , <i>his epigastric pain</i> , type=TrIP)
(<i>morphine</i> , <i>his epigastric pain</i> , type=TrAP)
(<i>demerol</i> , <i>his epigastric pain</i> , type=TrIP)

Fig. 1. An example of medical relations in a sample sentence. TrIP, treatment improves medical problem; TrAP, treatment is administered for medical problem.

scheme [10]. The second combination way is to feed features extracted by a RNN architecture into CNN [11, 12], which can be seen as generating new input representations by RNN. The third way is to stack RNN on CNN. Even though this architecture has not been applied to identify medical relations from clinical text, its variants have achieved remarkable results in many other classification tasks [13–16].

Deep learning methods have presented satisfactory results [2–7, 17, 18] and make the models less dependent on manual feature engineering. Moreover, some researchers proposed models only with word representations as input features [9, 19],

*Corresponding author

Email addresses: hebin_hit@hotmail.com (Bin He),

guanyi@hit.edu.cn (Yi Guan), 13B912003@hit.edu.cn (Rui Dai)

which achieved outstanding results. Similarly, our goal is to propose a model for relation classification on clinical records, without using any external feature set. In this work, we follow the third combination way and design a two-layer architecture: input representations (word-level) are fed into a CNN layer to learn n-gram features (phrase-level), and these feature representations are directly used as the input of a bidirectional gated recurrent unit (GRU) [20] layer to achieve the final sample representation (sentence-level). Our main contributions are as follows: (1) we propose a unified architecture to identify medical relations in clinical records, which has the ability to capture both local features (extracted by a CNN layer) and sequential correlations among these features (extracted by a bidirectional GRU layer); (2) we also explore training our model with attention mechanism (C-BGRU-Att) and compare the performance with the model using the conventional max-pooling operation (C-BGRU-Max); (3) experiments show our model achieves better performance than previous single-model methods, with only word embedding features.

2. Methodology

Figure 2 describes the architecture of our model for medical relation classification on clinical records. This model learns a distributed representation for each relation sample, and calculates final scores with relation type representations. More details will be discussed in the following sections.

2.1. Word representation layer

With reference to a previous study on relation classification [21], word position features capture information of the relative position between words and target concepts. Therefore, an word embedding matrix $W^w \in \mathbb{R}^{d^w \times |V^w|}$ and an word position embedding matrix $W^{wp} \in \mathbb{R}^{d^p \times |V^p|}$ are given in this work, where V^w is the vocabulary, V^p is the word position set, and d^w and d^p are pre-set embedding sizes. Every word in the relation sample is mapped to a column vector \mathbf{x}_i^w to represent the word feature. In addition, relative distances between the

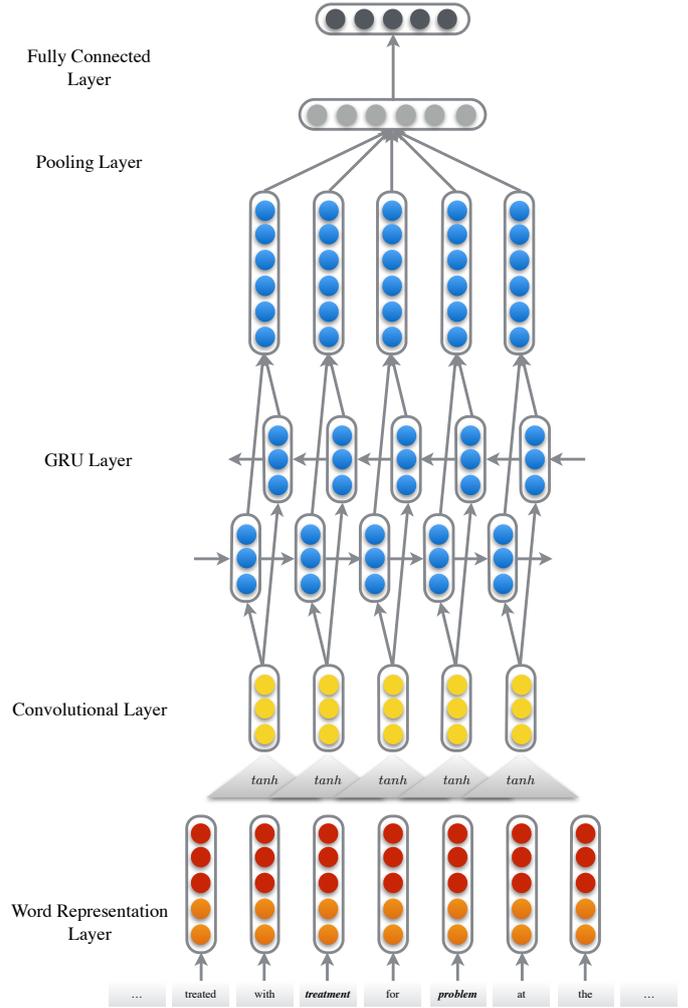


Fig. 2. Architecture of our model for medical relation classification. In the input of this architecture, concept contents in the relation sample “she was treated with [steroids]_{treatment} for [this swelling]_{problem} at the outside hospital , and these were continued .” are replaced by their concept types.

current word and the target concepts are mapped to word position vectors \mathbf{x}_i^{p1} and \mathbf{x}_i^{p2} . Based on the above features, each word can be represented by $\mathbf{x}_i' = [(\mathbf{x}_i^w)^T, (\mathbf{x}_i^{p1})^T, (\mathbf{x}_i^{p2})^T]^T$, and $\mathbf{x}_i' \in \mathbb{R}^{d^x}$, where $d^x = d^w + 2d^p$.

2.2. Convolutional layer

The semantic representations of n-grams are valuable features to the relation classification task, and convolution operation can capture this information by combining word embedding features in a fixed window. Given the input representation $\mathbf{x}' = (\mathbf{x}'_1, \mathbf{x}'_2, \dots, \mathbf{x}'_n)$ and a context window size k , concatenation of successive words in this window size can be defined

as $X_j = [\mathbf{x}'_j{}^T, \dots, \mathbf{x}'_{j+k-1}{}^T]^T$, and the representation of this relation sample can be reformatted as $X = (X_1, \dots, X_{n-k+1})$. Given a weight matrix of the convolutional filters W^{conv} and a linear bias \mathbf{b} , the local feature representations are computed:

$$C_j = \tanh(W^{conv} \cdot X_j + \mathbf{b}), \quad (1)$$

where $W^{conv} \in \mathbb{R}^{d^c \times d^x k}$, $\mathbf{b} \in \mathbb{R}^{d^c}$, and \tanh denotes the hyperbolic tangent function.

Generally, this convolutional result will be fed into a max-pooling operation to extract the most significant features. However, these extracted features are independent, and the correlation information among the local features are not captured. GRU has the ability to make up for this deficiency by using a gating mechanism to capture short-term and long-term dependencies. Therefore, in this study, a GRU layer is stacked on top of the convolutional layer to continue the feature extraction work.

2.3. GRU layer

Similar to the long short-term memory (LSTM) unit with a memory cell and three gating units [22, 23], GRU is much simpler to compute because only two gating units are used to adaptively capture dependencies over different time scales: one is the reset gate \mathbf{r}_j , which controls how much information from the previous hidden state is kept in the candidate hidden state; another is the update gate \mathbf{z}_j , which decides how much previous information contributes and how much information from the candidate hidden state is added. The computational process are demonstrated by the following equations:

$$\mathbf{r}_j = \sigma(W_r \cdot C_j + U_r \cdot \mathbf{h}_{j-1} + \mathbf{b}_r), \quad (2)$$

$$\mathbf{z}_j = \sigma(W_z \cdot C_j + U_z \cdot \mathbf{h}_{j-1} + \mathbf{b}_z), \quad (3)$$

$$\tilde{\mathbf{h}}_j = \tanh(W_h \cdot C_j + \mathbf{r}_j \odot (U_h \cdot \mathbf{h}_{j-1}) + \mathbf{b}_h), \quad (4)$$

$$\mathbf{h}_j = (1 - \mathbf{z}_j) \odot \mathbf{h}_{j-1} + \mathbf{z}_j \odot \tilde{\mathbf{h}}_j, \quad (5)$$

where σ is the logistic sigmoid function, \odot stands for the element-wise multiplication, C_j is the current n-gram feature representation (mentioned in Section 2.2), \mathbf{h}_{j-1} and $\tilde{\mathbf{h}}_j$ are the previous and the candidate hidden state, respectively, and

$\mathbf{h}_j \in \mathbb{R}^{d^h}$ is the current hidden state. $W_r, U_r, \mathbf{b}_r, W_z, U_z, \mathbf{b}_z, W_h, U_h$ and \mathbf{b}_h are weight matrices to be learned.

We use a bidirectional GRU [20] to encode the n-gram feature representations, which contains a forward GRU and a backward GRU. A sequence of forward hidden states $(\vec{\mathbf{h}}_1, \dots, \vec{\mathbf{h}}_{n-k+1})$ and a sequence of backward hidden states $(\overleftarrow{\mathbf{h}}_1, \dots, \overleftarrow{\mathbf{h}}_{n-k+1})$ are obtained. The final j -th hidden state can be achieved by concatenating the j -th forward and backward hidden state: $\mathbf{h}_j = [\vec{\mathbf{h}}_j{}^T, \overleftarrow{\mathbf{h}}_j{}^T]^T$, which contains the dependencies of the preceding and the following n-gram features.

2.4. Pooling layer

Two different kinds of pooling schemes are adopted to generate the semantic representation of the relation sample \mathbf{rs} .

Max pooling can be seen as a down-sampling operation that aims to extract the most significant features. After using this operation in our network, the i -th feature value \mathbf{rs}_i is calculated by

$$\mathbf{rs}_i = \max([\mathbf{h}_1]_i, \dots, [\mathbf{h}_{n-k+1}]_i), \quad (6)$$

where $[\mathbf{h}_j]_i$ denotes the i -th element in vector \mathbf{h}_j . And all these features constitute the semantic representation of the relation sample $\mathbf{rs} = (\mathbf{rs}_1, \dots, \mathbf{rs}_{d^h})^T$.

Attentive pooling Given the output of the GRU layer $H = [\mathbf{h}_1, \dots, \mathbf{h}_{n-k+1}]$, we follow the attention mechanism used in [9], and the representation \mathbf{rs} is formed:

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{v}^T \cdot \tanh(H)), \quad (7)$$

$$\mathbf{rs} = \tanh(H \cdot \boldsymbol{\alpha}^T), \quad (8)$$

where \mathbf{v} is a model parameter vector and $\boldsymbol{\alpha}$ is a weight vector to measure which parts of the GRU output are relatively significant for the relation classification.

2.5. Fully connected layer

We apply a softmax classifier to achieve the confidence scores with a class embedding matrix W^{cs} :

$$\mathbf{s}_\theta = \text{softmax}(W^{cs} \cdot \mathbf{rs}), \quad (9)$$

where θ is the model parameter set. s_{θ}^y is the confidence score of the true relation type y , and the loss function can be defined as

$$\mathcal{L} = -\frac{1}{m} \sum_{i=1}^m \log s_{\theta}^y + \beta \|\theta\|^2, \quad (10)$$

where m is the sample size and β is the l_2 regularization parameter.

3. Experiments

3.1. Dataset and experimental settings

Experiments are conducted on the 2010 i2b2/VA relation dataset¹ and the WI relation dataset². The former dataset comprises 426 English discharge summaries (170 for training and 256 for test), and the latter dataset contains 992 Chinese clinical records (521 for training and 471 for test). The relation types and their counts in these two datasets are listed in Table 1. As stipulated in the official evaluation metric in the 2010 i2b2/VA challenge, the model performance is based on the micro-averaged F1 score over all positive relation types.

In our methods, the initial word representations and the other matrices are randomly initialized by normalized initialization [25], and a 5-fold cross-validation is used on the training set to tune the model hyperparameters. The selected hyperparameter values are: word embedding size d^w , 100; word position embedding size d^p , 10; convolutional size d^c , 200; context window size k , 3; GRU dimension d^h , 100; learning rate, 0.01. Adam technique [26] is utilized to optimize our loss function. We use both l_2 regularization and dropout technique [27] to avoid overfitting, and the values are set to 0.0001 and 0.5, respectively.

3.2. Baselines

3.2.1. 2010 i2b2/VA relation dataset

When doing experiments on this dataset, the previous methods [12, 28, 29] followed inconsistent data split schemes. In

¹The relation dataset is available at <https://www.i2b2.org/NLP/Relations/>.

²https://github.com/WILAB-HIT/Resources/tree/master/entity_assertion_relation

Table 1
Relation type statistics.

2010 i2b2/VA relation dataset			WI Relation dataset		
Relation	Train	Test	Relation	Train	Test
TrIP	51	152	TrID	103	92
TrWP	24	109	TrWD	38	27
TrCP	184	342	TrAD	221	166
TrAP	885	1732	NTrD	675	656
TrNAP	62	112	TrIS	337	215
NTrP	1702	2759	TrWS	297	242
TeRP	993	2060	TrCS	125	176
TeCP	166	338	TrAS	334	238
NTeP	993	1974	NTrS	1062	901
PIP	755	1448	TeRD	301	227
NPP	4418	8089	NTeD	331	248
SID	969	620	TeRS	527	542
DCS	228	181	TeAS	313	564
NDS	777	635	NTeS	8628	7060

Positive relations were annotated in both relation datasets, and samples of negative relation types (starting with “N” in this table) were extracted to ensure each concept pair within a sentence could be assigned a certain relation type. For more details of these relation types, please refer to [1, 24].

order to compare these methods together, we choose the split scheme in [28], which is also the official data split.

SVM: due to the dataset available to the research community is only a subset of the dataset used in the 2010 i2b2/VA challenge, so Souza and Ng [28] reimplemented the state-of-the-art model in the challenge [30] and reevaluated this model on the relation dataset accessible. **SVM+ILP:** Souza and Ng [28] also proposed a better single-model method and an ensemble-based method within an integer linear programming (ILP) framework. In these feature-based state-of-the-art methods, a variety of external features sets are used, such as part-of-speech (POS) tagging and dependency parsing.

In this work, three previous neural network methods are reimplemented and reevaluated. **CNN:** a multiple-filter CNN with max-pooling proposed by Sahu et al. [29]. To evaluate the model performance independent of the external features, POS and chunk features used in this method are removed. **CRNN-Max and CRNN-Att:** a two-layer model comprising recurrent

and convolutional layers with max and attentive pooling [12]. However, only word embeddings were used in their work. In order to maintain a fair comparison, word position embeddings are added in our model reimplementation. In these three baseline reimplementations, we follow the selected hyperparameters used in the corresponding work and the word embeddings are pre-trained on the deidentified notes from the MIMIC-III database [31].

3.2.2. WI relation dataset

SVM: this model is implemented using scikit-learn³. And it involves the following features: entity e_1 , entity e_2 , entity type et_1 , entity type et_2 , distance between e_1 and e_2 , words in e_1 and e_2 , words between e_1 and e_2 , words behind e_2 , POS of words in e_1 and e_2 , POS of words between e_1 and e_2 , and POS of words behind e_2 .

CNN: the model version of C-BGRU-Max after removing the GRU layer, which is a CNN-based model.

3.3. Experimental results

3.3.1. System performance

The performance results are displayed in Table 2 and 3, including 95% confidence intervals for the models we implemented, which are derived using bootstrapping [32]. We use the same bootstrapping method described in [33]. We observe that our C-BGRU-Max model outperforms the previous single-model methods significantly in both datasets, without using any external features. After using attentive pooling, the model performance on the two datasets shows different changes: drops on the 2010 i2b2/VA relation dataset but increases on the WI relation dataset. The intuitive explanation is that descriptions in English discharge summaries tend to be more colloquial, making specific features more difficult to capture. More details of the category-wise and class-wise performance comparisons are listed in Table 4, 5, 6, and 7.

Table 2

System performance comparison with other models using the 2010 i2b2/VA relation dataset.

Classifier	External features	P	R	F1
<i>Single-model methods</i>				
SVM* [30]	Set1	58.1	66.7	62.1
SVM+ILP [28]	Set2	75.0	58.9	66.0
CNN [29]	None	68.0	55.1	60.9
		(67.4, 68.6)	(54.5, 55.7)	(60.4, 61.4)
CRNN-Max [12]	None	65.1	61.3	63.1
		(64.6, 65.6)	(60.7, 61.8)	(62.7, 63.6)
CRNN-Att [12]	None	63.2	58.5	60.7
		(62.6, 63.7)	(58.0, 59.0)	(60.3, 61.2)
C-BGRU-Max	None	69.3	66.3	67.8
		(68.8, 69.9)	(65.8, 66.8)	(67.3, 68.3)
C-BGRU-Att	None	69.6	63.7	66.5
		(69.0, 70.1)	(63.1, 64.2)	(66.0, 66.9)
<i>Ensemble-based method</i>				
Ensemble+ILP ^o [28]	Set2	66.7	72.9	69.6

The symbol * indicates that this model is reimplemented by [28] on the relation dataset available to the research community, due to the accessible dataset is only a subset of that used in the 2010 i2b2/VA challenge. The symbol ^o indicates that this classifier is the ensemble of 5 independent models. The bold item is the best result. Set1: POS, chunk, semantic role labeler, word lemma, dependency parse, assertion type, sentiment category, Wikipedia. Set2: POS, chunk, semantic role labeler, word lemma, dependency parse, assertion type, sentiment category, Wikipedia, manually labeled patterns. POS, part-of-speech; ILP, integer linear programming.

Table 3

System performance comparison using the WI relation dataset.

Classifier	P	R	F1
SVM	72.9	63.9	68.1
CNN	72.7	64.5	68.3
	(72.0, 73.4)	(63.7, 65.2)	(67.7, 69.0)
C-BGRU-Max	73.2	68.3	70.7
	(72.5, 73.9)	(67.6, 69.0)	(70.1, 71.3)
C-BGRU-Att	74.8	68.8	71.6
	(74.1, 75.5)	(68.1, 69.5)	(71.0, 72.3)

The bold item is the best result.

³<http://scikit-learn.org/stable/>.

Table 4

Category-wise performance comparison with other neural network models using the 2010 i2b2/VA relation dataset.

Classifier	TrP relations			TeP relations			PP relations		
	P	R	F1	P	R	F1	P	R	F1
	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)
CNN [29]	60.9	48.2	53.8	75.8	69.2	72.3	64.8	43.3	51.9
	1.0	0.9	0.9	0.9	0.8	0.7	1.3	1.1	1.1
CRNN-Max [12]	58.4	53.8	56.0	73.3	73.1	73.2	61.6	54.4	57.8
	0.9	0.9	0.8	0.8	0.8	0.7	1.1	1.2	1.0
CRNN-Att [12]	55.2	50.8	52.9	70.1	73.8	71.9	63.3	46.3	53.5
	0.9	0.9	0.8	0.8	0.8	0.7	1.3	1.2	1.1
C-BGRU-Max	62.7	59.7	61.2	78.4	77.5	77.9	64.8	58.9	61.7
	0.9	0.9	0.8	0.8	0.8	0.6	1.2	1.1	1.0
C-BGRU-Att	63.9	57.1	<u>60.4</u>	79.4	72.5	<u>75.8</u>	62.8	60.0	<u>61.4</u>
	0.9	0.9	0.8	0.7	0.8	0.7	1.1	1.2	1.0

TrP, Treatment-Problem; TeP, Test-Problem; PP, Problem-Problem. CI(±) is confidence interval for P, R, and F1. The bold item is the best result. Compared with previous models, the underlined item is statistically significant.

3.3.2. Discussion of attentive pooling

As show in Table 2, the F1 scores of CRNN-Att and C-BGRU-Att are lower than that of CRNN-Max and C-BGRU-Max, respectively. This indicates that the attention mechanism, which presents a positive effect in the general domain [9, 34], does not show any performance improvement on the 2010 i2b2/VA relation dataset. In this dataset, there exist ~3.3 entities in each sentence on average. Therefore, input representations of relation samples generated from the same sentence are quite similar, and the only difference is that some of the word position representations between these relation samples are different, which may not be able to show sufficient sample differentiation. In addition, attentive pooling does not extract the most significant features like max-pooling, which may lead to relative deficiencies in distinguishing model similar samples. We will try to analysis and validate these speculations in our future work.

3.3.3. F1 score vs. distance

Figure 3a and 4a show the frequency distribution of different distances in the two datasets, and Figure 3b and 4b depict the

Table 5

Class-wise performance comparison with other neural network models using the 2010 i2b2/VA relation dataset.

Classifier	TrIP			TrWP			TrCP			TrAP			TrNAP			TeRP			TeCP			PIP		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)
CNN [29]	18.7	5.9	9.0	0.0	0.0	0.0	30.1	39.3	63.0	61.1	62.0	35.9	9.1	14.5	77.7	77.8	77.7	45.0	16.9	24.5	64.8	43.3	51.9	
	5.0	1.7	2.5	0.0	0.0	0.0	2.2	2.4	1.0	1.0	0.9	7.5	2.3	3.5	0.8	0.8	0.6	3.9	1.7	2.3	1.3	1.1	1.1	
CRNN-Max [12]	33.0	4.6	8.1	0.0	0.0	0.0	28.1	33.6	61.3	69.8	65.3	11.7	3.2	5.0	76.9	80.1	78.5	41.8	30.4	35.2	61.6	54.4	57.8	
	8.9	1.5	2.6	0.0	0.0	0.0	2.1	2.3	1.0	1.0	0.8	5.1	1.4	2.2	0.8	0.8	0.6	2.8	2.2	2.2	1.1	1.2	1.0	
CRNN-Att [12]	0.0	0.0	0.0	0.0	0.0	0.0	10.8	16.4	56.4	69.6	62.3	0.0	0.0	0.0	71.5	83.0	76.8	45.7	17.8	25.6	63.3	46.3	53.5	
	0.0	0.0	0.0	0.0	0.0	0.0	3.7	1.4	1.0	1.0	0.8	0.0	0.0	0.0	0.8	0.7	0.6	3.9	1.9	2.4	1.3	1.2	1.1	
C-BGRU-Max	51.4	4.7	8.7	36.4	0.7	<u>1.4</u>	38.9	44.5	64.4	75.8	69.6	38.5	7.1	12.0	79.9	84.4	82.1	61.1	35.3	44.8	64.8	58.9	61.7	
	12.0	1.5	2.7	29.2	0.7	1.3	2.7	2.3	2.1	1.0	0.9	9.7	2.1	3.4	0.8	0.7	0.6	3.1	2.3	2.4	1.2	1.1	1.0	
C-BGRU-Att	43.8	11.1	17.6	26.7	2.9	5.3	48.5	41.1	44.5	67.7	70.9	69.3	30.1	8.8	13.6	81.2	79.3	80.3	58.7	30.8	40.4	62.8	60.0	61.4
	7.0	2.2	3.3	10.9	1.4	2.4	2.5	2.2	2.0	0.9	0.8	7.0	2.3	3.4	0.8	0.8	0.6	3.2	2.2	2.4	1.1	1.2	1.0	

CI(±) is confidence interval for P, R, and F1. The bold item is the best result. Compared with previous models, the underlined item is statistically significant.

Table 6

Category-wise performance of neural network models using the WI relation dataset.

Classifier	TrD relations			TrS relations			TeD relations			TeS relations			DS relations		
	P	R	F1	P	R	F1									
	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)	CI(±)									
CNN	59.5	56.3	57.8	58.1	50.5	54.1	90.1	90.7	90.4	86.9	59.0	70.3	72.5	82.7	77.3
	2.6	2.7	2.5	1.5	1.5	1.3	1.7	1.7	1.2	1.1	1.3	1.1	1.3	1.2	1.0
C-BGRU-Max	59.9	60.4	60.1	59.3	54.2	56.7	91.7	92.0	91.8	85.5	64.7	<u>73.7</u>	73.6	84.7	78.8
	2.4	2.5	2.3	1.5	1.5	1.3	1.6	1.5	1.1	1.1	1.3	1.0	1.2	1.1	1.0
C-BGRU-Att	61.8	58.7	60.2	64.4	55.5	<u>59.6</u>	91.6	93.6	92.5	82.6	68.0	<u>74.6</u>	75.0	80.8	77.8
	2.6	2.6	2.5	1.5	1.5	1.4	1.6	1.4	1.1	1.1	1.2	0.9	1.3	1.2	1.0

TrD, Treatment-Disease; TrS, Treatment-Symptom; TeD, Test-Disease; TeS, Test-Symptom; DS, Disease-Symptom. CI(±) is confidence interval for P, R, and F1. The bold item is the best result. Compared with CNN, the underlined item is statistically significant.

Table 7

Class-wise performance of neural network models using the WI relation dataset.

Classifier	TrID			TrWD			TrAD			TrIS			TrWS			TrCS		
	P	R	F1															
	CI(±)	CI(±)	CI(±)															
CNN	53.7	39.6	45.6	53.4	28.9	37.5	62.0	70.0	65.8	50.6	43.7	46.9	71.5	69.0	70.2	69.0	20.2	31.3
	5.4	4.7	4.5	11.7	7.6	8.4	3.1	3.2	2.6	3.0	3.0	2.6	2.5	2.7	2.0	5.5	2.7	3.5
C-BGRU-Max	53.7	39.1	45.3	56.6	44.4	49.8	62.3	74.7	67.9	54.3	42.0	47.3	75.3	70.5	72.8	57.5	27.7	37.4
	5.4	4.4	4.2	9.5	8.6	8.0	2.9	3.0	2.4	3.4	3.0	2.7	2.5	2.6	2.0	4.7	2.9	3.3
C-BGRU-Att	57.9	40.7	47.8	58.7	40.0	47.6	63.5	71.7	67.3	58.8	50.0	<u>54.1</u>	75.8	69.5	72.5	66.4	26.2	37.6
	5.5	4.6	4.4	10.4	8.4	8.1	3.1	3.1	2.6	3.2	3.0	2.7	2.5	2.6	2.1	5.0	2.9	3.4
Classifier	TrAS			TeRD			TeRS			TeAS			DCS			SID		
	P	R	F1															
	CI(±)	CI(±)	CI(±)															
CNN	50.2	60.3	54.8	90.1	90.7	90.4	88.8	84.1	86.4	82.9	35.0	49.2	56.6	63.2	59.7	77.0	88.4	82.3
	2.6	2.8	2.3	1.7	1.7	1.2	1.2	1.3	0.9	2.1	1.8	1.9	3.1	3.2	2.6	1.4	1.2	1.0
C-BGRU-Max	51.1	68.3	58.5	91.7	92.0	91.8	85.4	84.7	85.1	85.6	45.4	<u>59.4</u>	56.3	69.0	62.0	79.1	89.4	83.9
	2.4	2.7	2.1	1.6	1.5	1.1	1.4	1.4	1.0	1.7	1.9	1.7	2.8	3.0	2.4	1.3	1.1	0.9
C-BGRU-Att	58.6	67.7	<u>62.8</u>	91.6	93.6	92.5	80.6	86.0	83.2	86.2	50.7	<u>63.9</u>	58.7	61.7	60.2	79.6	86.4	82.9
	2.6	2.6	2.2	1.6	1.4	1.1	1.4	1.3	1.0	1.6	1.9	1.7	3.2	3.2	2.7	1.4	1.3	1.0

CI(±) is confidence interval for P, R, and F1. The bold item is the best result. Compared with CNN, the underlined item is statistically significant.

trend of the F1 score as the distance increases. The F1 score is the average value of the relation samples belonging to the distance window $[d - 2, d + 2]$. In order to ensure the reliability of the evaluation, the maximum distance value with a statistic greater than 20 is selected as the truncation of the distance value. On the 2010 i2b2/VA relation dataset, C-BGRU-Max and C-BGRU-Att outperform the baselines over all dis-

tances. On the WI relation dataset, C-BGRU-Max and CNN do not show significant differences when the distance is less than 20, but as the distance increases, the performance gap gradually expands. These results verify that our model has the ability to learn long-term dependencies and this information works in the relation classification task.

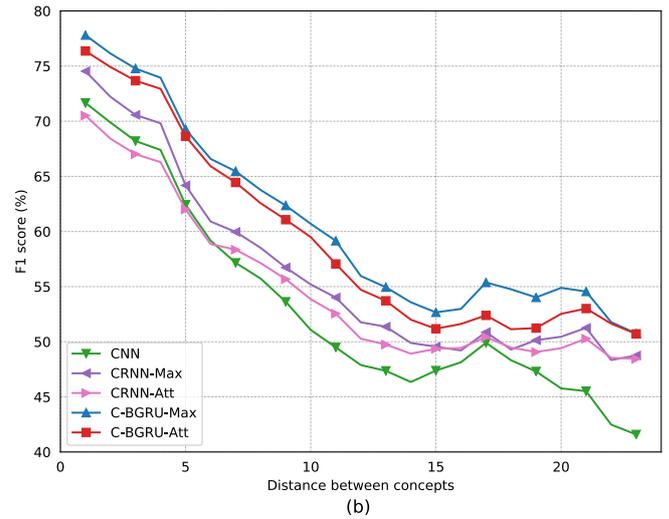
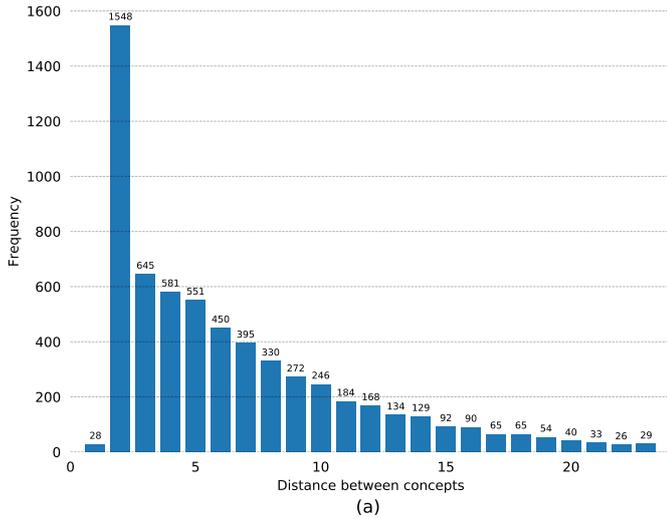


Fig. 3. The frequency distribution of the distance between concepts in the 2010 i2b2/VA relation dataset (a) and F1 score comparisons over different distances (b). The “distance” means the difference in word position between two concepts in the relation sample.

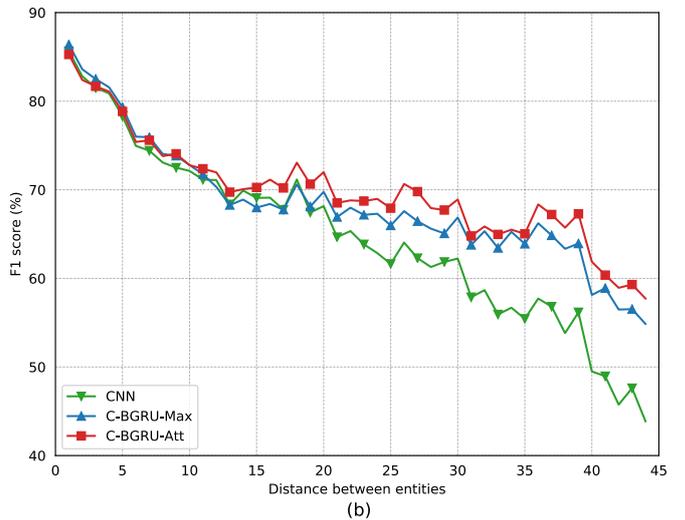
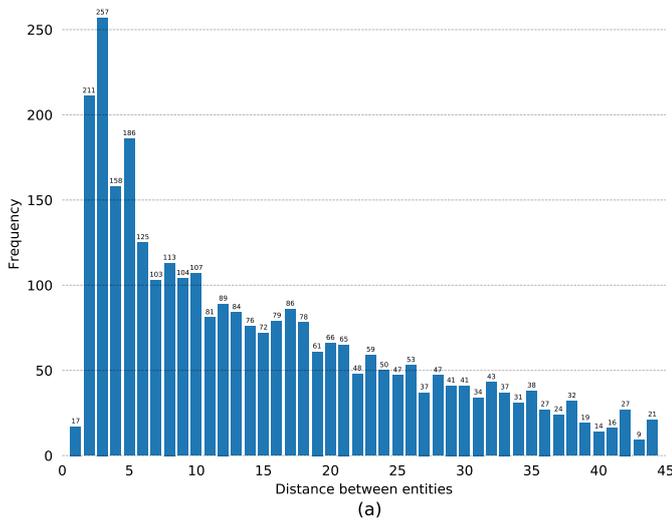


Fig. 4. The frequency distribution of the distance between entities in the WI relation dataset (a) and F1 score comparisons over different distances (b). The “distance” means the difference in word position between two entities in the relation sample.

4. Related Work

Before deep learning research became popular, statistical machine learning methods were the main approaches in the relation classification task. Most of the researchers in the general and clinical domain focused on feature-based and kernel-based methods [35–39].

In recent years, researchers have gradually tried the effect of deep learning methods in the relation classification task and achieved satisfactory results. A variety of deep architectures have been proposed to classify the relations, such as recurrent

neural network (MV-RNN) [17], CNN with softmax classification [21], factor-based compositional embedding model (FCM) [18], and word embedding-based models [40]. Next, there exist many RNN-based and CNN-based variants. Because the max-pooling operation in CNN models will lose significant linguistic features in a sentence, some researchers introduced dependency trees for this work, e.g., bidirectional long short-term memory networks (BLSTM) [2], dependency-based neural networks (DepNN) [3], shortest dependency path-based CNN [4], long short term memory networks along shortest dependency paths

(SDP-LSTM) [5], deep recurrent neural networks (DRNN) [6], and jointed sequential and tree-structured LSTM-RNN [7]. Although the above studies achieved solid results, further research was devoted to eliminating the dependence on the NLP parser because of its limited performance. dos Santos et al. [19] proposed a new pairwise ranking loss function, and only two class representations were updated in every training round. Similarly, Wang et al. [8] introduced a pairwise margin-based loss function and multi-level attention mechanism and achieved the new state-of-the-art results for relation classification.

More recently, neural network methods have show promising performance for relation classification on clinical records. Sahu et al. [29] proposed a multiple-filter CNN with some linguistic features, and experiments on the 2010 i2b2/VA relation dataset verified the effectiveness of the neural network model for medical relation classification. Raj et al. [12] trained a two-layer model by feeding short phrase features extracted by a bidirectional LSTM layer into CNN, and the model performed better than CNN on relation samples where the distance between the medical concepts are large. Different from Raj et al. [12]’s study, we think n-gram features and sequential correlations among them are the key to relation classification, so we explore another unified architecture that utilizes the strengths of CNN and RNN simultaneously.

5. Conclusion

In this paper, we present a unified architecture based on the combination of CNN and RNN to classify medical relations in English and Chinese clinical records. Our model captures long-term dependencies of phrase-level features through a bidirectional GRU layer and this information improves model performance. To the best of our knowledge, this is the first time that neural network methods have been used to classify relations in Chinese clinical text. Experiments show that the proposed model achieves a significant improvement over comparable methods on the 2010 i2b2/VA relation dataset and the WI relation dataset.

Acknowledgments

The authors would like to thank all the anonymous reviewers for their insightful comments. We would also like to thank the data support from the 2010 i2b2/VA challenge.

References

- [1] Uzuner, Ö., South, B.R., Shen, S., DuVall, S.L.. 2010 i2b2/VA challenge on concepts, assertions, and relations in clinical text. *Journal of the American Medical Informatics Association* 2011;18(5):552–556.
- [2] Zhang, S., Zheng, D., Hu, X., Yang, M.. Bidirectional Long Short-Term Memory Networks for Relation Classification. In: *PACLIC*. 2015, p. 73–78.
- [3] Liu, Y., Wei, F., Li, S., Ji, H., Zhou, M., WANG, H.. A Dependency-Based Neural Network for Relation Classification. In: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Beijing, China: Association for Computational Linguistics; 2015, p. 285–290.
- [4] Xu, K., Feng, Y., Huang, S., Zhao, D.. Semantic Relation Classification via Convolutional Neural Networks with Simple Negative Sampling. In: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Lisbon, Portugal: Association for Computational Linguistics; 2015, p. 536–540.
- [5] Xu, Y., Mou, L., Li, G., Chen, Y., Peng, H., Jin, Z.. Classifying Relations via Long Short Term Memory Networks along Shortest Dependency Paths. In: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 2015, p. 1785–1794.
- [6] Xu, Y., Jia, R., Mou, L., Li, G., Chen, Y., Lu, Y., et al. Improved relation classification by deep recurrent neural networks with data augmentation. In: *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. Osaka, Japan: The COLING 2016 Organizing Committee; 2016, p. 1461–1470.
- [7] Miwa, M., Bansal, M.. End-to-End Relation Extraction using LSTMs on Sequences and Tree Structures. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Berlin, Germany: Association for Computational Linguistics; 2016, p. 1105–1116.
- [8] Wang, L., Cao, Z., de Melo, G., Liu, Z.. Relation Classification via Multi-Level Attention CNNs. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2016, p. 1298–1307.
- [9] Zhou, P., Shi, W., Tian, J., Qi, Z., Li, B., Hao, H., et al. Attention-Based Bidirectional Long Short-Term Memory Networks for Relation Classification. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* 2016;:207–212.

- [10] Vu, N.T., Adel, H., Gupta, P., Schütze, H.. Combining Recurrent and Convolutional Neural Networks for Relation Classification. In: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. San Diego, California: Association for Computational Linguistics; 2016, p. 534–539.
- [11] Cai, R., Zhang, X., Wang, H.. Bidirectional Recurrent Convolutional Neural Network for Relation Classification. Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL 2016) 2016;:756–765.
- [12] Raj, D., Sahu, S.K., Anand, A.. Learning local and global contexts using a convolutional recurrent network model for relation classification in biomedical text. Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017) 2017;:311–321.
- [13] Tang, D., Qin, B., Liu, T.. Document Modeling with Gated Recurrent Neural Network for Sentiment Classification. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. 2015, p. 1422–1432.
- [14] Nguyen, T.H., Grishman, R.. Combining Neural Networks and Log-linear Models to Improve Relation Extraction. arXiv preprint arXiv:151105926 2015;.
- [15] Zhou, C., Sun, C., Liu, Z., Lau, F.. A C-LSTM neural network for text classification. arXiv preprint arXiv:151108630 2015;.
- [16] Choi, K., Fazekas, G., Sandler, M., Cho, K.. Convolutional recurrent neural networks for music classification. In: Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on. IEEE; 2017, p. 2392–2396.
- [17] Socher, R., Huval, B., Manning, C.D., Ng, A.Y.. Semantic Compositionality through Recursive Matrix-Vector Spaces. Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning 2012;(Mv):1201–1211.
- [18] Yu, M., Gormley, M.R., Dredze, M.. Factor-based Compositional Embedding Models. NIPS Workshop on Learning Semantics 2014;:95–101.
- [19] dos Santos, C.N., Xiang, B., Zhou, B.. Classifying Relations by Ranking with Convolutional Neural Networks. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). 3; 2015, p. 626–634.
- [20] Bahdanau, D., Cho, K., Bengio, Y.. Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:14090473 2014;.
- [21] Zeng, D., Liu, K., Lai, S., Zhou, G., Zhao, J.. Relation Classification via Convolutional Deep Neural Network. COLING 2014;(2011):2335–2344.
- [22] Hochreiter, S., Schmidhuber, J.. Long short-term memory. Neural computation 1997;9(8):1735–1780.
- [23] Graves, A.. Supervised sequence labelling with recurrent neural networks; vol. 385. Springer; 2012.
- [24] He, B., Dong, B., Guan, Y., Yang, J., Jiang, Z., Yu, Q., et al. Building a comprehensive syntactic and semantic corpus of Chinese clinical texts. Journal of Biomedical Informatics 2017;69:203–217.
- [25] Glorot, X., Bengio, Y.. Understanding the difficulty of training deep feedforward neural networks. In: Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS); vol. 9. 2010, p. 249–256.
- [26] Kingma, D., Ba, J.. Adam: A method for stochastic optimization. arXiv preprint arXiv:14126980 2014;.
- [27] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. Journal of Machine Learning Research 2014;15:1929–1958.
- [28] Souza, J.D., Ng, V.. Ensemble-Based Medical Relation Classification. In: Hajic, J., Tsujii, J., editors. COLING. ACL; 2014, p. 1682–1693.
- [29] Sahu, S.K., Anand, A., Oruganty, K., Gattu, M.. Relation extraction from clinical texts using domain invariant convolutional neural network. Proceedings of the 15th Workshop on Biomedical Natural Language Processing 2016;:71.
- [30] Rink, B., Harabagiu, S., Roberts, K.. Automatic extraction of relations between medical concepts in clinical texts. Journal of the American Medical Informatics Association 2011;18(5):594–600.
- [31] Johnson, A.E., Pollard, T.J., Shen, L., Lehman, L.W.H., Feng, M., Ghassemi, M., et al. MIMIC-III, a freely accessible critical care database. Scientific Data 2016;3.
- [32] DiCiccio, T.J., Efron, B.. Bootstrap confidence intervals. Statistical Science 1996;11(3):189–228.
- [33] Gao, S., Young, M.T., Qiu, J.X., Yoon, H.J., Christian, J.B., Fearn, P.A., et al. Hierarchical attention networks for information extraction from cancer pathology reports. Journal of the American Medical Informatics Association 2017;.
- [34] Kim, J., Lee, J.H.. Multiple Range-Restricted Bidirectional Gated Recurrent Units with Attention for Relation Classification. arXiv preprint arXiv:170701265 2017;.
- [35] Bunescu, R., Mooney, R.. A shortest path dependency kernel for relation extraction. In: Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing. 2005, p. 724–731.
- [36] Hendrickx, I., Kim, S.N., Kozareva, Z., Nakov, P., Ó Séaghdha, D., Padó, S., et al. Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. In: Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions. Association for Computational Linguistics; 2009, p. 94–99.
- [37] Rink, B., Harabagiu, S.. UTD: Classifying Semantic Relations by Combining Lexical and Semantic Resources. Proceedings of the 5th International Workshop on Semantic Evaluation 2010;(July):256–259.
- [38] Zhu, X., Cherry, C., Kiritchenko, S., Martin, J., de Bruijn, B.. Detecting concept relations in clinical text: Insights from a state-of-the-art

model. *Journal of Biomedical Informatics* 2013;46(2):275–285.

- [39] Kim, J., Choe, Y., Mueller, K.. Extracting clinical relations in electronic health records using enriched parse trees. In: *Procedia Computer Science*; vol. 53. 2015, p. 274–283.
- [40] Hashimoto, K., Stenetorp, P., Miwa, M., Tsuruoka, Y.. Task-Oriented Learning of Word Embeddings for Semantic Relation Classification. In: *Proceedings of the Nineteenth Conference on Computational Natural Language Learning*. Beijing, China: Association for Computational Linguistics; 2015, p. 268–278.