

Classifying medical relations in clinical text via convolutional neural networks

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

Deep learning research on relation classification has achieved solid performance in the general domain. This study proposes a convolutional neural network (CNN) architecture with a multi-pooling operation for medical relation classification on clinical records and explores a loss function with a category-level constraint matrix. Experiments using the 2010 i2b2/VA relation corpus demonstrate these models, which do not depend on any external features, outperform previous single-model methods and our best model is competitive with the existing ensemble-based method.

Keywords: Relation classification; Clinical text; Convolutional neural network; Multi-pooling

1. Introduction

Relation classification, a natural language processing (NLP) task identifying the relation between two entities in a sentence, is an important technique used in many subsequent NLP applications such as question answering and knowledge base completion. This task has been widely studied in the general domain due to the large number of accessible datasets such as the SemEval-2010 task 8 dataset [1], which aims to classify the relation between two nominals in the same sentence.

In the clinical domain, Informatics for Integrating Biology and the Bedside (i2b2) released an annotated relation corpus on clinical records, attracting considerable attention [2]. Identifying relations in clinical records is more challenging than relations in the general domain because one sentence from clinical records may contain more than two medical concepts and concepts may be comprised of several words. For example, the sentence “at that time , she also had *cat scratch fever* and she had *resection of an abscess in the left lower extremity*” contains three concepts, two of which contain more than two words. The

annotated information given in the 2010 i2b2/VA relation corpus thus differs from that in the SemEval-2010 task 8 dataset. In the former, the category to which a concept pair belongs is given, and the classification objective is to identify the subcategory, also known as the *relation type*.

Deep neural networks have become a research trend in recent years due to powerful learning ability features without manual feature engineering. Various neural architectures have been proposed for classifying relations in general [3–6], biomedical [7–12] and clinical text [13, 14]. Conventional convolutional neural network (CNN) models use max-pooling operations to extract the most significant feature in a convolutional filter; however, information regarding feature positioning relative to the concepts cannot be captured. Responding to this issue, Chen et al. [15] designed a dynamic multi-pooling method to extract features from each part of a sentence in the argument classification task. A chunk-based max-pooling algorithm, proposed by [16], splits each sentence into a fixed number of segments to retain more semantics from the sentence for the statistical machine translation model. The position of features relative to concepts is vital for medical relation classification on clinical records. Based on the above studies, this study proposed a CNN-based method (without any external features) for

*Corresponding author

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

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

recognizing medical concept relations in clinical records. Its contributions are as follows:

- A multi-pooling operation was introduced into the proposed CNN architecture, which aims to capture the position information of local features relative to the concept pair.
- A novel loss function with a category-level constraint matrix was explored.
- The proposed models achieved improved performance compared to previous single-model methods, and the best model is competitive with the ensemble-based method for classifying relations between medical concepts.

2. Corpus and data preprocessing

The relation corpus¹ used in this study was released in the 2010 i2b2/VA challenge, and is comprised of 426 discharge summaries. Of these, 170 were used for training, and the remaining 256 for testing². This dataset contains three types of concepts (*medical problem*, *treatment*, and *test*), and each concept pair in the same sentence was assigned one relation type. Medical concept relations in this corpus can be grouped into 3 categories: *medical problem-treatment*, *medical problem-test*, and *medical problem-medical problem* relations. Table 1 describes the definitions³ and statistics of these relation types.

Although words within sentences were already separated by spaces, additional splits were required for some specific strings. This study employed the Natural Language Toolkit⁴ (NLTK) to tokenize sentence strings in clinical records, then realigned concept boundaries to avoid concept information errors. Tokens⁵ were lowercase, and numbers were replaced by zero.

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

²This follows the official data split in the 2010 i2b2/VA challenge.

³2010 i2b2/VA Challenge Evaluation Relation Annotation Guidelines: <http://www.i2b2.org/NLP/Relations/assets/Relation%20Annotation%20Guideline.pdf>.

⁴Natural Language Toolkit: <http://www.nltk.org/>.

⁵Definition of *token*: <http://nlp.stanford.edu/IR-book/html/htmledition/tokenization-1.html>.

Table 1
Relation type statistics.

Relation	Definition	Train	Test
<i>Medical problem-Treatment relations</i>			
TrIP	Treatment improves medical problem	51	152
TrWP	Treatment worsens medical problem	24	109
TrCP	Treatment causes medical problem	184	342
TrAP	Treatment is administered for medical problem	885	1732
TrNAP	Treatment is not administered because of medical problem	62	112
NTrP	No relation between treatment and problem	1702	2759
<i>Medical problem-Test relations</i>			
TeRP	Test reveals medical problem	993	2060
TeCP	Test conducted to investigate medical problem	166	338
NTeP	No relation between test and problem	993	1974
<i>Medical problem-Medical problem relations</i>			
PIP	Medical problem indicates medical problem	755	1448
NPP	No relation between two medical problems	4418	8089

Eight positive relation types were annotated in this relation corpus, and samples of three negative relation types (starting with “N” in this table) were extracted for model training to ensure each concept pair within a sentence could be classified into a certain relation type.

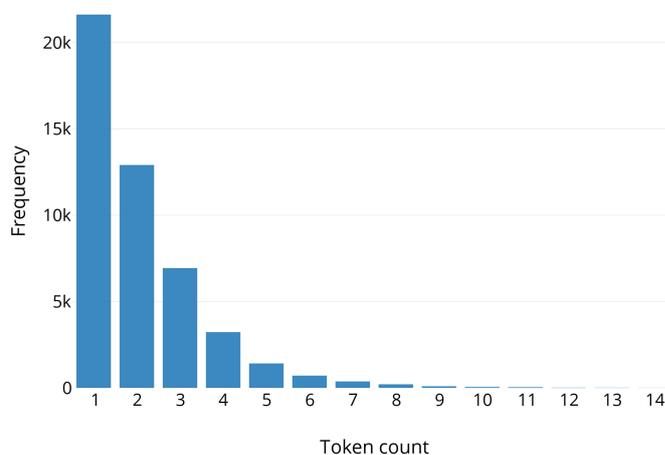


Fig. 1. Frequency distribution of token count in medical concepts. Concept lengths appearing less than five times were filtered.

Fig. 1 lists the frequency distribution of token count in medical concepts, showing over half the concepts contained more than one token. According to these statistics, the average token count for all concepts in the corpus was 2.09 (*medical problem* 2.42, *treatment* 1.85, and *test* 1.86, respectively). Considering

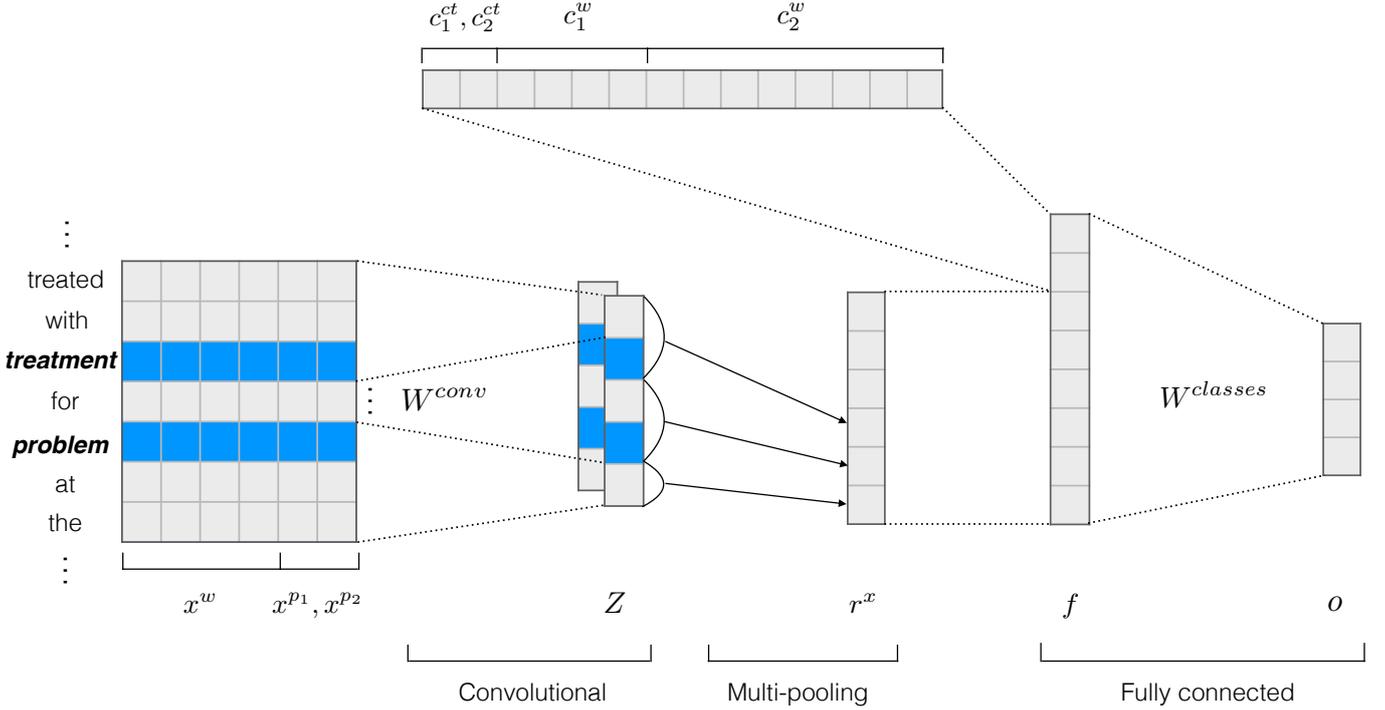


Fig. 2. Architecture of CNN-based model for medical relation classification. Concept contents in the sample sentence “she was treated with [steroids]_{treatment} for [this swelling]_{medical problem} at the outside hospital , and these were continued .” are replaced by their concept types.

concepts containing more than one token were detrimental to n-gram feature extraction in the proposed model, the concept pair contents were replaced by their concept types to make local features more apparent. The details of the replacement are shown in Fig. 2.

3. Methodology

Fig. 2 describes the architecture of the proposed CNN-based model for medical relation classification on clinical records. This model learns a distributed representation for each relation sample. A feature vector is generated to represent each sentence sample x containing two concepts, and final scores are calculated with relation type representations. Further details are discussed in the following subsections.

3.1. Word representation

In this relation classification task, the following information was given for each sample sentence $x = (x_1, x_2, \dots, x_n)$: (1) concept position indexes in the sentence c_1^{index} and c_2^{index} ; (2) concept contents $c_1 = (c_{11}, \dots, c_{1l_1})$ and $c_2 = (c_{21}, \dots, c_{2l_2})$;

(3) concept types c_1^{type} and c_2^{type} ; and (4) y is the sample’s relation type.

Previous studies on relation classification [4–6, 13] utilized word position features to capture information on the proximity of words to target concepts. Therefore, an word embedding matrix $W^w \in \mathbb{R}^{d^w \times |V^w|}$ and an word position embedding matrix $W^{wp} \in \mathbb{R}^{d^p \times |V^p|}$ were utilized in this study, where V^w represented the vocabulary, V^p represented the word position set, and d^w and d^p were pre-set embedding sizes. For each sample sentence, every word was mapped to a column vector $x_i^w \in \mathbb{R}^{d^w}$ representing the word feature. Additionally, the relative distances between the current word and concepts were mapped to the position vectors $x_i^{p1} \in \mathbb{R}^{d^p}$ and $x_i^{p2} \in \mathbb{R}^{d^p}$. Based on these features, each word could be represented as $x_i^x = [(x_i^w)^T, (x_i^{p1})^T, (x_i^{p2})^T]^T \in \mathbb{R}^{d^x}$, where d^x was the word vector size and $d^x = d^w + 2d^p$.

3.2. Convolutional multi-pooling

Semantic representations of n-grams are valuable features in relation classification tasks, and convolution operation can cap-

ture this information by combining word representations in a fixed window. Given a sample sentence $x = (x_1, x_2, \dots, x_n)$ and a context window size k , the concatenation of successive words in this window size could be defined as:

$$x_{i:i+k-1} = [(x'_i)^T, \dots, (x'_{i+k-1})^T]^T,$$

and $x_{i:i+k-1} \in \mathbb{R}^{d^x k}$. The representation of this sentence could be reformatted as $X = (x_{1:k}, \dots, x_{n-k+1:n})$ and $X \in \mathbb{R}^{d^x k \times (n-k+1)}$. The input X would then be fed into the convolutional layer to generate local features. Given W^{conv} as the weight matrix of the convolutional filters and a linear bias B_1 , a linear transformation followed by a non-linear function are calculated:

$$Z = f(W^{conv} \cdot X + B_1),$$

where $W^{conv} \in \mathbb{R}^{d^c \times d^x k}$, $B_1 \in \mathbb{R}^{d^c}$, f is the relu function, and the convolutional result is $Z \in \mathbb{R}^{d^c \times (n-k+1)}$.

Max-pooling operations are generally used to extract the most significant feature in a convolutional filter [4]; however, these are insufficient for relation classification on clinical records. The dataset used here contains ~ 3.3 concepts in one sentence, making the position of features relative to the concept pair necessary for relation classification. A multi-pooling operation was introduced in the proposed method to achieve more local features in each sentence. Although word position information was included in word representations, multi-pooling strengthened the significance of the relative position information.

Given the concept position index of a concept pair described in Section 3.1, the convolutional result Z can be split into three parts: $Z^1 = Z_{1:(c_1^{index}-1)}$, $Z^2 = Z_{c_1^{index}:(c_2^{index}-1)}$, and $Z^3 = Z_{c_2^{index}:(n-k+1)}$, where $Z_{p:q} = [Z_p, \dots, Z_q] \in \mathbb{R}^{d^c \times (q-p+1)}$. Max-pooling operations are then performed on each part to extract the most valuable features, defined as $r_j^l = \max[Z^l]_j$, $r^l \in \mathbb{R}^{d^c}$, and $l \in \{1, 2, 3\}$. These three vectors can be concatenate into the single vector

$$r^x = [(r^1)^T, (r^2)^T, (r^3)^T]^T \in \mathbb{R}^{3d^c},$$

creating an informative semantic representation of the sentence.

3.3. Concept feature representation

As described in Section 2, the concept types in a concept pair are given, allowing their relation category to be known directly. In response to this situation, concept type information is typically used in two ways: to train multiple independent models, or to train one model by adding the concept types as features. Both methods have distinct advantages and disadvantages: the former cannot maintain unified word representation, and each model loses some samples to update the word embedding matrix; the latter may produce misclassifications across categories.

In order to maintain unified word representation and tend to model simplicity, the latter method was selected here for model building, and two vectors were used to represent two concept types mapped from a concept type embedding matrix $W^{ct} \in \mathbb{R}^{d^{ct} \times |V^{ct}|}$. In the matrix, V^{ct} represents concept type set and d^{ct} represents a pre-set concept type embedding size. This concept type feature representation was formalized as $x^{ct} = [(c_1^{ct})^T, (c_2^{ct})^T]^T \in \mathbb{R}^{2d^{ct}}$. Concept content, in addition to the n-gram and concept type features described above, is also necessary for the relation classification model. Word embeddings of the concept contents were added to supplement concept feature representation, which can be formalized as $c^{fx} = [(x^{ct})^T, (c_1^w)^T, (c_2^w)^T]^T \in \mathbb{R}^{d^{cf}}$, where $c_i^w = [(c_{i1}^w)^T, \dots, (c_{ij}^w)^T]^T$ is the concept content representation, $i \in \{1, 2\}$, c_{ij}^w is the word representation of the j th word in the i th concept, and $d^{cf} = 2d^{ct} + d^w \times (l_1 + l_2)$.

3.4. Class embeddings and scoring

The n-gram feature representation and concept feature representation were concatenated into the single vector $rc = [(r^x)^T, (c^{fx})^T]^T$, and the confidence of each relation type with a class embedding matrix $W^{classes} \in \mathbb{R}^{m \times (3d^c + d^{cf})}$ was computed as

$$s = W^{classes} \cdot rc,$$

where each row vector $W_l^{classes}$ can be viewed as the representation of relation type l and m equals the number of relation types.

Training with logsoftmax. After obtaining relation type scores, a softmax operation was applied to obtain the probability of each relation type:

$$p(y|x, \theta) = \frac{e^{s_y}}{\sum_{l \in \mathcal{Y}} e^{s_l}},$$

where s_y is the score for the relation type y , \mathcal{Y} is the relation type set, and $\theta = (W^w, W^{wp}, W^{conv}, B_1, W^{ct}, W^{classes})$. Based on this probability, the loss function could be defined as

$$\mathcal{L} = -\log p(y|x, \theta) + \beta(\|W^w\|^2 + \|W^{wp}\|^2 + \|W^{conv}\|^2 + \|W^{ct}\|^2 + \|W^{classes}\|^2),$$

and β was the L_2 regularization parameter.

Category-level constraint. Training one model to cover all categories may cause cross-category misclassifications. It would also be inappropriate to regard samples in other categories as negative samples. Therefore, a loss function with a category-level constraint matrix was proposed:

$$\mathcal{L}_C = \log\left(\sum_{l \in \mathcal{Y}} C_{Vl}^x \cdot e^{s_l}\right) - s_y + \beta(\|C^x \cdot W^{classes}\|^2 + \|W^w\|^2 + \|W^{wp}\|^2 + \|W^{conv}\|^2 + \|W^{ct}\|^2),$$

$$C_{ij}^x = \begin{cases} 1, & \text{if } i = j \text{ and } i \in \text{Category}_x; \\ 0, & \text{otherwise.} \end{cases}$$

Here, C^x represents the constraint matrix of relation type indexes, l' represents the index number of relation type l , and Category_x represents the relation type index set for the category that sample x belongs to. After using this loss function during the training of sample x , only the class vectors $W_l^{classes} (l' \in \text{Category}_x)$ will be updated, and the other class vectors remain unchanged. This avoids treating samples in other categories as negative.

4. Experiments

4.1. Experimental setup

Evaluation metric. As shown in Table 1, there are eight positive relation types and three negative relation types. Precision,

recall, and F1-measure were used to evaluate the performance of each positive relation type. Simultaneously, as stipulated in the official evaluation metric [2], model performance was defined based on the micro-averaged F1 score across all positive relation types.

Parameter settings. Initial word representations were trained using the word2vec tool [17] and de-identified notes from the MIMIC-III database [18]. The other matrices in the proposed method were randomly initialized by normalized initialization [19]. The word embedding size was set to 50 and the concept type embedding size to 5, equal to those in [13]. The dropout technique [20] was used in the concatenated representation rc to avoid overfitting, and this value was set to 0.5. One fifth of the training set was randomly selected as the development set during experiments, and the model hyperparameters were tuned using a grid search: word position embedding size (5, 10, 20, 30); convolutional filter size (100, 200, 300, 400); learning rate (0.01, 0.025, 0.05, 0.075, 0.1); L_2 regularization parameter (0.00005, 0.0001, 0.0005, 0.001). The selected hyperparameter values were 10, 200, 0.075, and 0.0005, respectively.

4.2. Experimental results

Three method comparisons were designed: (1) CNN-Max, the CNN-based model using max-pooling in the convolutional layer; (2) CNN-Multi, the CNN-based model using multi-pooling in the convolutional layer; and (3) CNN-Multi-C, where the CNN-Multi model was trained with category-level constraint.

4.2.1. Filter window sizes and word embedding initializations

The efficacy of different filter window sizes and word embedding initializations was investigated using the CNN-Multi model. For each filter window size, model performance was evaluated under two word embedding initializations: (1) pre-trained, where word embeddings are initialized by pre-trained word embeddings as described in Section 4.1; and (2) randomly initialized, where word embeddings are randomly initialized by

normalized initialization [19]. Table 2 shows the system performance by measures of precision, recall, and F1 score.

Table 2

CNN-Multi model performance using various convolutional window sizes and different word embedding initializations.

Window size	Pre-trained			Randomly initialized		
	P	R	F1	P	R	F1
[3]	73.7	64.1	68.5	73.2	65.8	69.3
[4]	72.2	63.6	67.6	73.1	66.7	69.7
[5]	72.7	62.7	67.3	74.5	62.3	67.9
[3,4]	73.9	61.6	67.2	71.0	67.5	69.2
[3,5]	73.2	62.1	67.2	71.2	67.1	69.1
[4,5]	71.7	65.6	68.5	75.3	61.7	67.8
[3,4,5]	70.7	67.0	68.8	70.8	67.7	69.2

Pre-trained word embeddings demonstrated lower F1 scores in most window sizes. The highest F1 score was achieved using a window size 4 and randomly initialized word embeddings. Therefore, all proposed models were trained using a filter window size 4 and randomly initialized word embeddings.

4.2.2. Comparison with baselines

Previous methods [13, 14, 21] followed inconsistent data split schemes. To compare these to the proposed methods, all methods were reimplemented and evaluated using the official data split of the 2010 i2b2/VA relation corpus [2], as shown in Table 1. All model hyperparameters remained unchanged during the reimplementation. To maintain a fair comparison, the part-of-speech and chunk features used in [13] were removed, and word position embeddings were added to Raj et al. [14]’s models. The performance results are displayed in Table 3, including 95% confidence intervals for each performance metric derived via bootstrapping [22]. The same bootstrapping method described in [23] was used.

System performance. Rink et al. [24] presented a support vector machine (SVM) method and achieved the best result in the 2010 i2b2/VA challenge. As the relation corpus available to the research community is a subset of that used during the 2010 i2b2/VA challenge, Souza and Ng [21] re-implemented

this method using the accessible dataset and obtains a F1 score of 62.1. They also proposed an improved single-model method and an ensemble-based method within an integer linear programming (ILP) framework, which became the new single-model and ensemble-based state-of-the-art methods, respectively. CNN achieves a slightly lower F1 score than SVM. CRNN-Max improved upon SVM, but still lags behind the single-model state-of-the-art method SVM+ILP. The models proposed in this paper outperformed SVM+ILP. CNN-Multi achieved a very similar result to Ensemble+ILP without depending on external features, and was significantly less complex to implement.

Table 3

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

Classifier	External features	P	R	F1
<i>Single-model methods</i>				
SVM [24]	Set1	66.7	58.1	62.1
SVM+ILP [21]	Set2	58.9	75.0	66.0
CNN [13]	None	72.2	54.1	61.8
		(70.9, 73.5)	(52.8, 55.3)	(60.7, 62.9)
CRNN-Max [14]	None	62.0	64.6	63.3
		(60.8, 63.1)	(63.4, 65.7)	(62.2, 64.3)
CRNN-Att [14]	None	64.7	56.5	60.3
		(63.4, 65.9)	(55.3, 57.7)	(59.2, 61.4)
CNN-Max	None	73.4	62.4	67.5
		(72.2, 74.6)	(61.2, 63.6)	(66.4, 68.5)
CNN-Multi	None	73.1	66.7	69.7
		(71.9, 74.2)	(65.4, 67.8)	(68.7, 70.6)
CNN-Multi-C	None	72.8	65.9	69.2
		(71.7, 74.0)	(64.7, 67.0)	(68.2, 70.2)
<i>Ensemble-based method</i>				
Ensemble+ILP [21]	set2	72.9	66.7	69.6

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.

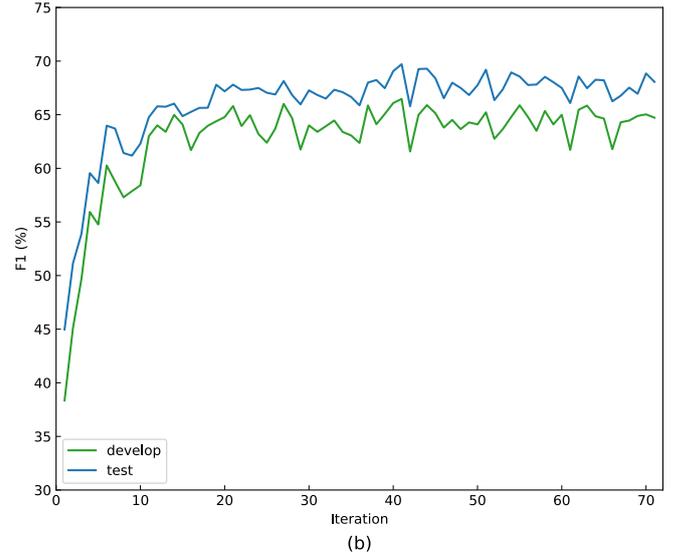
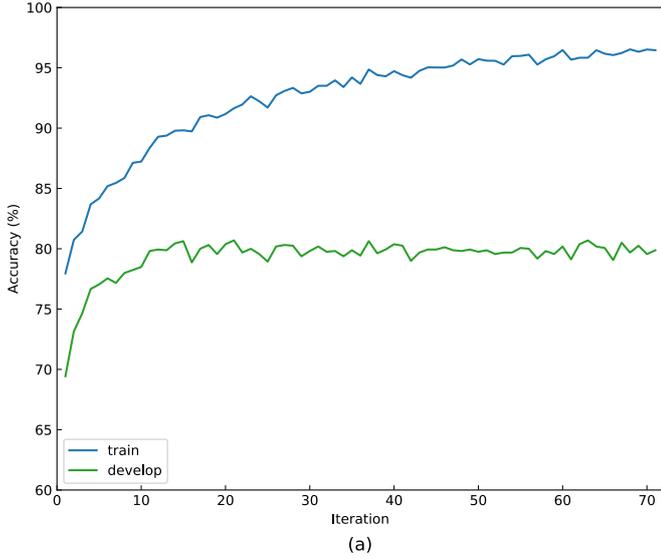


Fig. 3. Training progress of the CNN-Multi model.

Category-wise performance. Table 4 shows the performance of the neural network methods in the three relation categories. All methods demonstrated a better performance on *medical problem-test* relations, potentially due to two conditions: (1) the relation type number of *medical problem-treatment* relations is twice that of *medical problem-test* relations; and (2) as shown in Table 1, *medical problem-medical problem* relations have a high relation type imbalance, which is adverse for classification. Compared with CNN-Max, CNN-Multi obtained significantly higher F1 scores for both *medical problem-treatment* and *medical problem-medical problem* relations, but improved much less for *medical problem-test* relations. This may indicate the relative position information of features, extracted via multi-pooling operation, works well for relatively complex relation classifications whereas max-pooling is sufficient for simpler relation classifications. Among these neural network methods, CNN-Multi performed best for *medical problem-treatment* and *medical problem-medical problem* relations, whereas CNN-Multi-C performed best in *medical problem-test* relations.

Class-wise performance. Table 5 shows the performance of neural network methods for each positive relation type. In combination with Table 1, this demonstrates that relation types with

Table 4

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

Classifier	TrP relations			TeP relations			PP relations		
	P	R	F1	P	R	F1	P	R	F1
CNN [13]	64.5	47.5	54.7	79.5	68.6	73.7	70.6	41.0	51.9
CRNN-Max [14]	53.9	60.2	56.9	68.7	77.1	72.7	65.3	51.4	57.5
CRNN-Att [14]	62.8	46.7	53.6	66.0	75.7	70.5	64.3	41.2	50.2
CNN-Max	67.1	54.3	60.0	80.3	76.4	78.3	70.4	53.2	60.6
CNN-Multi	68.1	60.0	63.8	77.9	79.3	78.6	72.0	56.9	63.6
CNN-Multi-C	67.9	58.3	62.7	79.3	78.3	78.8	68.9	58.1	63.1

TrP, Medical problem-Treatment; TeP, Medical problem-Test; PP, Medical problem-Medical problem.

a small training size ($TrIP$, 51; $TrWP$, 24; $TrNAP$, 62) provided poor performance, and class-wise performance improved as the training size increased.

4.3. More analysis

Model training progress. Fig. 3(a) shows accuracy curves of the CNN-Multi model. The accuracy curve of the training set maintained high values, indicating the model fit the dataset well. The accuracy of the development set generally tends to stabilize after 15 iterations. In Fig. 3(b), the F1 score curves of the development and test sets show similar trends; these curves appear less smooth because the training set was shuffled for

Table 5

Class-wise performance comparison with other neural network models using the 2010:12b2/VA relation corpus.

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
CNN [13]	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	65.6	60.9	63.2	44.4	3.6	6.6	80.4	77.7	79.0	57.5	13.6	22.0	70.6	41.0	51.9
CRNN-Max [14]	35.8	12.5	18.5	0.0	0.0	0.0	43.9	41.8	42.8	55.7	75.6	64.1	100.0	0.9	1.8	74.4	83.0	78.5	35.6	41.4	38.3	65.3	51.4	57.5
CRNN-Att [14]	0.0	0.0	0.0	0.0	0.0	0.0	40.5	38.6	39.5	68.0	58.4	62.8	0.0	0.0	0.0	66.0	88.2	75.5	0.0	0.0	0.0	64.3	41.2	50.2
CNN-Max	0.0	0.0	0.0	0.0	0.0	0.0	54.3	43.9	48.5	69.2	67.9	68.5	100.0	1.8	3.5	80.7	85.0	82.8	72.3	24.0	36.0	70.4	53.2	60.6
CNN-Multi	50.0	2.0	3.8	0.0	0.0	0.0	56.7	41.8	48.1	69.6	76.2	72.8	100.0	2.7	5.2	77.8	88.9	83.0	81.6	21.0	33.4	72.0	56.9	63.6
CNN-Multi-C	100.0	2.0	3.9	0.0	0.0	0.0	54.9	42.4	47.9	69.7	73.6	71.6	57.1	3.6	6.7	79.8	86.5	83.0	70.6	28.4	40.5	68.9	58.1	63.1

each iteration. F1 score on the development set reached its optimal value at the 41st iteration, after which system parameters were maintained to evaluate system performance on the test set.

Errors. Table 6 contains no cross-category misclassification, and relation samples are evidently more often misclassified as the relation type whose training size is larger. This is due to the fact that during multi-class model training, models have more offset for classes with larger training sizes. Considering this situation, sampling methods can be considered as a strategy to improve model performance in future works.

Effect of category-level constraint. As shown in Table 3, system performance declined after adding a category-level constraint into the CNN-Multi model, whereas advantages of this constraint were not reflected. There are two potential reasons for this: (1) the CNN-Multi model was well trained and no cross-category misclassifications existed; and (2) the sample number of each relation type was too small for updating only the class vectors $W_l^{classes}(l' \in \text{Category}_x)$ during training to improve the generalization capability of the model. This constraint may be effective when experimenting with large datasets.

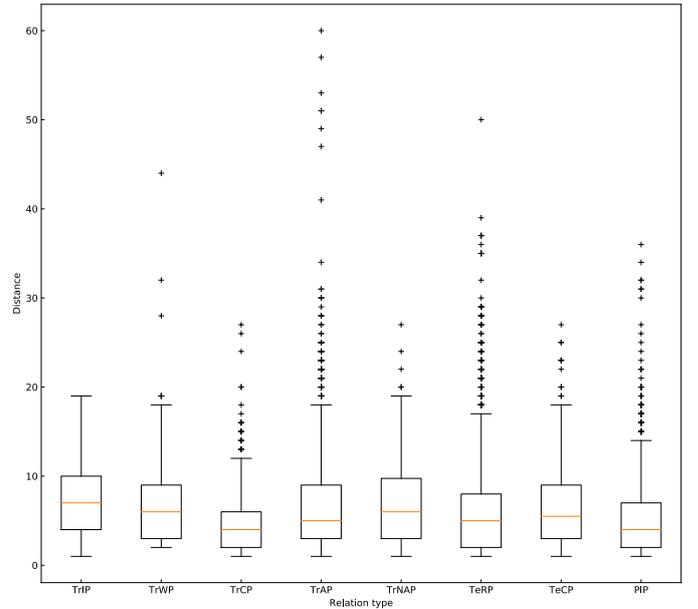


Fig. 4. Distance between medical concepts. Distance was calculated from the number of tokens between two medical concepts.

Table 6

Confusion matrix of the system output of the CNN-Multi model.

	System output										
	TrIP	TrWP	TrCP	TrAP	TrNAP	NTrP	TeRP	TeCP	NTeP	PIP	NPP
TrIP	3		<i>13</i>	<i>93</i>		<i>43</i>					
TrWP			<i>11</i>	<i>38</i>		<i>60</i>					
TrCP			143	<i>73</i>		<i>126</i>					
TrAP			<i>26</i>	1320		<i>386</i>					
TrNAP			<i>15</i>	<i>53</i>	3	<i>41</i>					
NTrP	<i>3</i>		<i>44</i>	<i>319</i>		2393					
TeRP							1831	<i>7</i>	<i>222</i>		
TeCP							<i>112</i>	71	<i>155</i>		
NTeP							<i>411</i>	<i>9</i>	1554		
PIP										824	<i>624</i>
NPP										<i>320</i>	7769

Zero items were removed from this table. Correctly classified items are bolded and the remainder are italicized.

Distance between medical concepts. The distance between two medical concepts was calculated from the number of tokens between the concepts. Fig. 4 illustrates the distance distribution of different relation types. In most samples, the distance between concepts was less than 20 tokens, however, there are still some long-distance relations, which are more challenging to be classified.

5. Related work

Before deep learning research became popular, most relation classification tasks used statistical machine learning methods. Many researchers in the general and medical domains focused on feature-based and kernel-based methods [21, 24–28], which are limited by conditions such as manual feature engineering and dependence on existing NLP toolkits.

More recently, researchers began investigating the performance of deep learning methods in relation classification tasks and achieved satisfactory results. Various deep architectures have been proposed for relation classification in the general domain, including the recurrent neural network (MV-RNN) [3], CNN with softmax classification [4], factor-based compositional embedding model (FCM) [29], and word embedding-based models [30]. Many RNN- and CNN-based variants ex-

ist. Because the max-pooling operation in CNN models experiences significant linguistic feature losses in sentences, some researchers introduced dependency trees for this application such as bidirectional long short-term memory networks (BLSTM) [31], dependency-based neural networks (DepNN) [32], shortest dependency path-based CNN [33], long short term memory networks along shortest dependency paths (SDP-LSTM) [34], deep recurrent neural networks (DRNN) [35], and jointed sequential and tree-structured LSTM-RNN [36]. Although the above studies achieved solid results, further research was devoted to eliminating the dependence on NLP parsers. dos Santos et al. [5] proposed a new pairwise ranking loss function where only two class representations were updated in every training round. Similarly, Wang et al. [6] introduced a pairwise margin-based loss function and multi-level attention mechanism, achieving new state-of-the-art results for relation classification. Some feature-free neural network methods also exist for relation classification on biomedical and clinical text. Liu et al. [7] employed CNN for drug-drug interaction (DDI) extraction, and Quan et al. [9] proposed a multichannel convolutional neural network (MCCNN) for this task. Furthermore, several attention-based methods [10–12, 14] were presented for automated biomedical relation extraction. This paper aimed is

to train a feature-free CNN-based model for classifying medical relations in clinical records.

6. Conclusion

This paper presented a novel CNN-based model for classifying relations between medical concepts in clinical records. The model performed well classifying relations in the 2010 i2b2/VA relation corpus. A multi-pooling operation helped to extract more precise and richer features in the convolutional layer, indicating feature extraction based on concept pair positioning can improve the efficacy of relation classification on clinical records. Although this model was applied to relations between medical concepts in clinical records, it could also be adapted to classify relations in other domains.

Acknowledgments

The authors are grateful to the Editors and the anonymous reviewers for their insightful comments. We would also like to thank the data support from the 2010 i2b2/VA challenge, as well as the useful discussions with Kehai Chen and Conghui Zhu.

References

References

- [1] 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. *Computational Linguistics* 2010;(June 2009):94–99.
- [2] 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.
- [3] 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.
- [4] Zeng, D., Liu, K., Lai, S., Zhou, G., Zhao, J.. Relation Classification via Convolutional Deep Neural Network. *COLING 2014;(2011):2335–2344.*
- [5] 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.
- [6] 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.
- [7] Liu, S., Tang, B., Chen, Q., Wang, X.. Drug-Drug Interaction Extraction via Convolutional Neural Networks. *Computational and Mathematical Methods in Medicine* 2016;2016.
- [8] Zhao, Z., Yang, Z., Luo, L., Lin, H., Wang, J.. Drug drug interaction extraction from biomedical literature using syntax convolutional neural network. *Bioinformatics* 2016;32(22):3444–3453.
- [9] Quan, C., Hua, L., Sun, X., Bai, W.. Multichannel convolutional neural network for biological relation extraction. *BioMed Research International* 2016;2016.
- [10] Asada, M., Miwa, M., Sasaki, Y.. Extracting Drug-Drug Interactions with Attention CNNs. In: *BioNLP 2017*. 2017, p. 9–18.
- [11] Sahu, S.K., Anand, A.. Drug-Drug Interaction Extraction from Biomedical Text Using Long Short Term Memory Network. *arXiv preprint arXiv:170108303* 2017;.
- [12] Zheng, W., Lin, H., Luo, L., Zhao, Z., Li, Z., Zhang, Y., et al. An attention-based effective neural model for drug-drug interactions extraction. *BMC Bioinformatics* 2017;18(1).
- [13] 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.
- [14] 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.
- [15] Chen, Y., Xu, L., Liu, K., Zeng, D., Zhao, J.. Event Extraction via Dynamic Multi-Pooling 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*. 2015, p. 167–176.
- [16] Zhang, J., Zhang, D., Hao, J.. Local translation prediction with global sentence representation. In: *International Joint Conference on Artificial Intelligence*; vol. 2015-Janua. 2015, p. 1398–1404.
- [17] Mikolov, T., Corrado, G., Chen, K., Dean, J.. Efficient Estimation of Word Representations in Vector Space. *Proceedings of the International Conference on Learning Representations* 2013;:1–12.
- [18] 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.

- [19] 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.
- [20] 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.
- [21] Souza, J.D., Ng, V.. Ensemble-Based Medical Relation Classification. In: Hajic, J., Tsujii, J., editors. COLING. ACL; 2014, p. 1682–1693.
- [22] DiCiccio, T.J., Efron, B.. Bootstrap confidence intervals. *Statistical Science* 1996;11(3):189–228.
- [23] 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;.
- [24] 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.
- [25] 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.
- [26] 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.
- [27] 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.
- [28] 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.
- [29] Yu, M., Gormley, M.R., Dredze, M.. Factor-based Compositional Embedding Models. *NIPS Workshop on Learning Semantics* 2014;:95–101.
- [30] 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.
- [31] Zhang, S., Zheng, D., Hu, X., Yang, M.. Bidirectional Long Short-Term Memory Networks for Relation Classification. In: *PACLIC*. 2015, p. 73–78.
- [32] 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.
- [33] 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.
- [34] 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.
- [35] 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.
- [36] 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.