A Joint Neural Model for Information Extraction with Global Features

Ying Lin, Heng Ji, Fei Huang, Lingfei Wu

ACL2020

Tasks

- This work jointly performs three tasks of Information Extraction at sentence-level:
 - + Entity Extraction.
 - + Relation Extraction.
 - + Event Extraction.

Model

- Their model performs all the tasks in four stages:
 - + Encoding tokens.
 - + Identifying nodes (i.e., triggers or entity mentions).
 - + Scoring nodes and edges (i.e., relations or argument roles).
 - + Searching for best graph.

Model



Model: Encoding Tokens

- Input: a setence of L words.
- Uses BERT as the encoder.
- Word representations are the averaged vector of their wordpiece representations.



Model: Identifying Nodes

• BERT outputs are fed into a Feed-Forward Net to obtain score vectors.

 $\hat{\boldsymbol{y}}_i = \operatorname{FFN}(\boldsymbol{x}_i)$

• Node identification is formalized as a sequence labeling task (e.g., B-Life:Marry, B-GPE) with a CRF layer.



Model: Scoring Nodes and edges

- At this layer, the model computes node and edge representations.
 - + Node: average sum over its component words.
 - + Edge: concatenation of node representations.
- Then, score vectors for nodes ($\hat{y}_i^t = FFN^t(v_i)$) & edges ($\hat{y}_k^t = FFN^t(v_i, v_j)$) are computed via softmax layers.
- Node that, the model does not make predictions here.



Model: Searching for Best Graph

• With the score vectors obtained from the previous step. They use Beam search to efficiently find the configuration with the highest score.



• Score of a graph is computed by:

$$s(G) = s'(G) + \boldsymbol{u}\boldsymbol{f}_G$$

Where:

+
$$s'(G)$$
 is local score: $s'(G) = \sum_{t \in T} \sum_{i=1}^{N^t} \max \hat{y}_i^t$

+ uf_G is global score where $f_G = \{f_1(G), ..., f_M(G)\}$ global feature vector

Model: Global Features

• Global features are introduced to capture cross-subtask and cross-instance dependencies.

Categary	Description				
Role	1. The number of entities that act as $< role_i > $ and $< role_j > $ arguments at the same time.				
	2. The number of <event_type;> events with <number> <role;> arguments.</role;></number></event_type;>				
	3. The number of occurrences of $$, $$, and $$ combination.				
	4. The number of events that have multiple <rolei> arguments.</rolei>				
	5. The number of entities that act as a <role<sub>i> argument of an <event_type<sub>j> event and a <role<sub>k> argument of an <event_type<sub>1> event at the same time.</event_type<sub></role<sub></event_type<sub></role<sub>				
Relation	6. The number of occurrences of $< entity_type_i >$, $< entity_type_j >$, and $< relation_type_k >$ combination.				
	7. The number of occurrences of $< entity_type_i > and < relation_type_j > combination.$				
	8. The number of occurrences of a <relation_type<sub>i> relation between a <role<sub>j> argument and a <role<sub>k> argument of the same event.</role<sub></role<sub></relation_type<sub>				
	9. The number of entities that have a <relation_typei> relation with multiple entities.</relation_typei>				
	10. The number of entities involving in <relation_type<sub>i> and <relation_type<sub>j> relations simultaneously.</relation_type<sub></relation_type<sub>				
Trigger	11. Whether a graph contains more than one $$ event.				

Training

• Identification loss: negative log-likelihood

$$\mathcal{L}^{\mathrm{I}} = -\log p(\boldsymbol{z}|\boldsymbol{X}) = -s(\boldsymbol{X}, \boldsymbol{z}) + \log \sum_{\hat{\boldsymbol{z}} \in Z} e^{s(\boldsymbol{X}, \hat{\boldsymbol{z}})}$$

Classification loss: cross-entropy

$$\mathcal{L}^{\mathrm{t}} = -\frac{1}{N^{t}} \sum_{i=1}^{N^{t}} \boldsymbol{y}_{i}^{t} \log \boldsymbol{\hat{y}}_{i}^{t}$$

• Global feature constraint: the ground-truth graph G should be the one with the highest score. Mimize this:

$$\mathcal{L}^{\mathbf{G}} = s(\hat{G}) - s(G)$$

• Overall loss: $\mathcal{L} = \mathcal{L}^I + \sum_{t \in T} \mathcal{L}^t + \mathcal{L}^G$

Experiments

• Datasets: ACE, ERE

Dataset	Split	#Sents	#Entities	#Rels	#Events
	Train	10,051	26,473	4,788	-
ACE05-R	Dev	2,424	6,362	1,131	-
	Test	2,050	5,476	1,151	-
	Train	17,172	29.006	4,664	4,202
ACE05-E	Dev	923	2,451	560	450
	Test	832	3,017	636	403
	Train	6,841	29,657	7,934	2,926
ACE05-CN	Dev	526	2,250	596	217
	Test	547	2,388	672	190
	Train	19,240	47,525	7,152	4,419
ACE05-E ⁺	Dev	902	3,422	728	468
	Test	676	3,673	802	424
	Train	14,219	38,864	5,045	6,419
ERE-EN	Dev	1,162	3,320	424	552
	Test	1,129	3,291	477	559
	Train	7,067	11,839	1,698	3,272
ERE-ES	Dev	556	886	120	210
	Test	546	811	108	269

Experiments

• Monolingual performance on English language:

Dataset	Task	DYGIE++	BASELINE	ONEIE
ACE05_P	Entity	88.6	-	88.8
ACE03-K	Relation	63.4	-	67.5
	Entity	89.7	90.2	90.2
	Trig-I	-	76.6	78.2
ACE05 E	Trig-C	69.7	73.5	74.7
ACE05-E	Arg-I	53.0	56.4	59.2
	Arg-C	48.8	53.9	56.8

Experiments

• Multilingual performance (with additional English data) on Chinese and Spanish.

Dataset	Training	Entity	Relation	Trig-C	Arg-C
ACE05 CN	CN	88.5	62.4	65.6	52.0
ACE05-CN	CN+EN	89.8	62.9	67.7	53.2
EDE ES	ES	81.3	48.1	56.8	40.3
ERE-ES	ES+EN	81.8	52.9	59.1	42.3