Towards Integration of Discriminability and Robustness For Document-Level Relation Extraction



Jia Guo^{1,2}, Stanley Kok¹, Lidong Bing²

¹School of Computing, National University of Singapore ²DAMO Academy, Alibaba Group



Introduction		Method
Sentence-level relation	extraction	 Our PEMSCL model
 Single sentence input 	Input sentence: General Director <u>Placido Domingo</u> will return to sing and conduct	I. <u>Pairwise moving-threshold loss with Entropy Minimization</u>
 Single entity mention 	on the Washington stage.	II. <u>Supervised Contrastive Learning for multi-labels and long-tailed relations</u>
 Single entity pair 	Given entity pair: < <u>Placido Domingo</u> , <u>Washington stage</u> >	III. Negative label sampling strategy
\circ Single label output		Small difference

<u>Single</u> laber output

Output relation: "org:top_members/employees"

- **Document-level relation extraction** An example from a sentence-level RE dataset (TACRED)
- <u>Multiple</u> sentence inputs Ο
- <u>Multiple</u> entity mentions Ο
- Multiple entity pairs Ο
- Multiple label outputs Ο
- Challenges:
- Inadequate in effectively **a**. distinguishing relations
- Lack of sufficient learning b. for long-tailed relations
- Vulnerable to annotation С. errors or missing annotations

Input document: (1) The culture of Los Angeles is rich with arts and ethnically diverse. (2) The greater Los Angeles metro area has several notable art museums including the Los Angeles County Museum of Art_[2] (LACMA_[2]), the J. Paul Getty Museum on the Santa Monica mountains overlooking the Pacific, the Museum of Contemporary Art_[6] (MOCA[6]), the Hammer Museum[7]</sub> and the Norton Simon Museum₁₈₁. (3) In the 1920s and 1930s Given entity pair:

<J. Paul Getty Museum, Los Angeles> Output relations: "located in the administrative territorial entity", "headquarters location".

An example from a document-level RE dataset (DocRED)

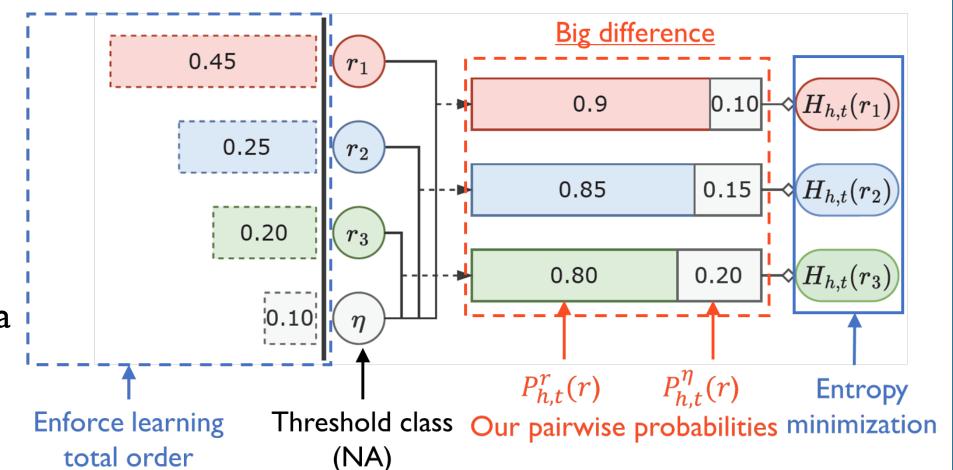
- **Our contributions:**
- A novel model based on Pairwise moving-threshold loss, Entropy Minimization, and Ο Supervised Contrastive Learning (PEMSCL)
 - Boost the <u>discriminability</u> of both probabilistic scores and internal embeddings

Problem definition:

□ <u>Inputs</u>:

*

- A document: $D = \{w_l\}_{l=1}^{L}$ containing L words
- A set of entities $\mathcal{E}_D = \{e_i\}_{i=1}^{|\mathcal{E}_D|}$
- Each entity e_i is associated with a set of mentions $\mathcal{M}_{e_i} = \{m_j^i\}_{j=1}^{|\mathcal{M}_{e_i}|}$ **Outputs:**



- For each entity pair, $(e_h, e_t)_{h,t=1,...,|\mathcal{E}_D|,h\neq t}$, the model predicts a subset of pre-defined relations $\mathcal{R} = \{r_k\}_{k=1}^{|\mathcal{R}|}$
- If an entity pair does not express any relation, it is labeled as NA
- Pairwise moving-threshold loss with entropy minimization
- \Box Split the predefined relation set into two mutually exclusive sets for (h, t): $\mathcal{R} = \mathcal{P}_{h,t} \cup \mathcal{N}_{h,t}$
 - Positive relations (i.e., the labels): $\mathcal{P}_{h,t}$
 - Negative relations: $\mathcal{N}_{h,t}$
- \Box Compare each r to the threshold class (NA), define their pairwise probabilities: $P_{h,t}(C = r | C = \{r, NA\}) = \frac{P_{h,t}^r(r)}{\exp(f_r)} = \frac{\exp(f_r)}{\exp(f_r) + \exp(f_\eta)}$ $P_{h,t}(C = \mathsf{NA}|C = \{r, \mathsf{NA}\}) = P_{h,t}^{\eta}(r) = 1 - P_{h,t}^{r}(r) = \frac{\exp(f_{\eta})}{\exp(f_{r}) + \exp(f_{n})}$
- Adapt supervised contrastive learning for <u>long-tailed relations</u> b.
- Improve model robustness by a novel negative label sampling strategy
- Validate the effectiveness of our model in various settings Ο

Negative Label Sampling Strategy

Annotation error problem:

Some entity pairs labeled as NA class should have at least one relation label [Tan et al., EMNLP'22].

 \rightarrow The negative relations of NA examples may be wrong!

 \Box We <u>randomly sample a small set</u> of negative relations for NA-labeled entity pairs. For $\mathcal{B}_{\mathcal{N}}, \mathcal{L}_1$ is

modified as:
$$\mathcal{L}' = \sum_{(h,t)\in\mathcal{B}_{\mathcal{N}}} \sum_{r\in\mathcal{N}'_{h,t}} -\log P^{\eta}_{h,t}(r) + \frac{1}{\gamma_2} \sum_{r\in\mathcal{N}'_{h,t}} H_{h,t}(r)$$

 \Box 1st loss function with negative label sampling :

$$\mathcal{L}_{1}^{\mathrm{NA}} = \mathcal{L}_{1}' + \sum_{\substack{(h,t) \in \mathcal{B}_{\mathcal{P}} \\ \text{Improve robustness}}} \mathcal{L}_{pmt}^{h,t} + \mathcal{L}_{em}^{h,t}}$$

$$\texttt{Our PEMCL with negative label sampling:} \qquad \texttt{Improve discriminability}$$

$$\mathcal{L}^{\mathrm{NA}} = \mathcal{L}_{1}^{\mathrm{NA}} + \lambda \mathcal{L}_{2}$$

Our pairwise moving-threshold loss:
 ↓ \$\mathcal{L}_{pmt}^{h,t} = -\log\left(\prod_{r \in \mathcal{P}_{h,t}} P_{h,t}^{r}(r) \prod_{r \in \mathcal{N}_{h,t}} (1 - P_{h,t}^{r}(r))\right)
 □ The definition
$$H_{h,t}(r) = -P_{h,t}^{r}$$
 □ The regularing $H_{h,t}(r) = -P_{h,t}^{r}$
 □ The regularing $L_{em}^{h,t}(r) = -P_{h,t}^{r}$
 □ The regularing $L_{em}^{h,t} = \frac{1}{\gamma_{1}} \sum_{r \in \mathcal{P}_{h,t}} L_{em} = \frac{1}{\gamma_{1}} \sum_$

on of information entropy: $P_{h,t}^{r}(r) \log P_{h,t}^{r}(r) - P_{h,t}^{\eta}(r) \log P_{h,t}^{\eta}(r)$ rization of entropy minimization:

$$\mathcal{L}_{em}^{h,t} = \frac{1}{\gamma_1} \sum_{r \in \mathcal{P}_{h,t}} H_{h,t}(r) + \frac{1}{\gamma_2} \sum_{r \in \mathcal{N}_{h,t}} H_{h,t}(r)$$

nction:

- $\sum \mathcal{L}_{pmt}^{h,t} + \mathcal{L}_{em}^{h,t}$ $(h,t)\in\mathcal{B}$
- Supervised contrastive learning for multi-labels & long-tailed relations
- "Pull" the embeddings of similar examples together, and "push" dissimilar examples apart:

$$\mathcal{L}_{scl}^{h,t} = -\log\left\{\frac{1}{|\mathcal{S}_{h,t}|}\sum_{p\in\mathcal{S}_{h,t}}\frac{\exp(|\mathbf{x}_{h,t}\cdot\mathbf{x}_p/\tau)|}{\sum_{d\in\mathcal{B},d\neq(h,t)}\exp(|\mathbf{x}_{h,t}\cdot\mathbf{x}_d/\tau)|}\right\} \text{ pull push }$$

 \Box Handling long-tailed relations \rightarrow for entity pairs with empty positive examples in a batch:

 $\mathcal{L}^{h,t}_{l\prime}$

$$\mathcal{L}_{lt}^{h,t} = \log \sum_{d \in \mathcal{B}, d \neq (h,t)} \exp(\mathbf{x}_{h,t} \cdot \mathbf{x}_d / \tau) \downarrow \quad \text{push}$$

2nd Loss function:

$$\mathcal{L}_{2} = \sum_{(h,t)\in\mathcal{B}_{\mathcal{P}}} \mathbb{I}_{\{|\mathcal{S}_{h,t}|\neq 0\}} \mathcal{L}_{scl}^{h,t} + \mathbb{I}_{\{|\mathcal{S}_{h,t}|=0\}}$$

Final loss function:

 $\mathcal{L} = \mathcal{L}_1 + \lambda \mathcal{L}_2$

Experiments and Analysis

Benchmarks:

• Our PEMSCL outperforms previous strong baselines on the original Our negative label sampling strategy is effective and robust in the noisy settings

• DocRED [Yao et al., ACL'19] & Re-DocRED [Tan et al., EMNLP'22] (both $|\mathcal{R}| = 96$).

- Two new data regimes
 - OOG-DocRE / OGG-DocRE:
 - Original labels for the train set
 - Original labels / Gold labels for the dev set
 - Gold labels for the test set

	DocRE	ED Dev	DocRI	ED Test	
Model	Ign F_1	F_1	Ign F_1	F_1	On OOG-Doci ATLOP (Zhou
Implemented on DeBERTa _{Large}					NCRL (Zhou a
ATLOP (Zhou et al., 2021)	62.16 ± 0.15	64.01 ± 0.12	62.12	64.08	PEMSCL (Our
ATLOP + BCE (Zhou and Lee, 2022)	$61.92 {\pm} 0.13$	$63.96 {\pm} 0.15$	61.83	63.92	PEMSCL [†] (Ou

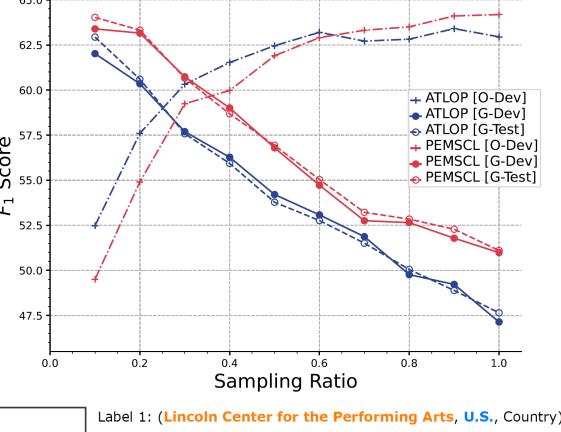
DocRED dataset and its cleaned version, the Re-DocRED dataset.

NCRL (Zhou and Lee, 2022)

PEMSCL (Ours)

62.98±0.18 64.79±0.13 63.03 64.96 63.25±0.09 65.15±0.10 63.40 65.41 **Re-DocRED** Dev **Re-DocRED** Test Ign F_1 F_1 Ign F_1

	Orig-Dev		Gold	-Dev	Gold-Test	
	Ign F_1	F_1	Ign F_1	F_1	Ign F_1	F_1
On OOG-DocRE Regime						
ATLOP (Zhou et al., 2021)	60.94	62.95	46.99	47.14	47.52	47.65
NCRL (Zhou and Lee, 2022)	61.42	63.52	49.06	49.21	48.41	48.53
PEMSCL (Ours)	62.05	64.19	<u>50.82</u>	<u>50.99</u>	<u>50.92</u>	<u>51.10</u>
PEMSCL [†] (Ours)	46.07	49.51	62.05	63.39	62.76	64.03
Dn OGG-DocRE Regime						
ATLOP (Zhou et al., 2021)	-	-	48.23	48.54	48.50	48.77
NCRL (Zhou and Lee, 2022)	-	-	49.92	50.08	50.10	50.25
PEMSCL (Ours)	-	-	<u>50.43</u>	<u>50.62</u>	<u>51.09</u>	<u>51.25</u>
PEMSCL [†] (Ours)	-	-	62.40	63.72	62.47	63.73



Deteret	Train	Dev	Test	Implemented on RoBERTa _{Large}	(0.12)	70.00		70.05	• The logit difference between	1. The Avery Fisher Career Crant, established by Avery	Label 1: (Lincoln Center for	• the Performing Arts, U.S., Count
Dataset	#Doc / #Example	#Doc / #Example	#Doc / #Example	JEREX (Eberts and Ulges, 2021)	69.12	70.33	68.97	70.25	the relation and the threshold	1. The Avery Fisher Career Grant, established by Avery Fisher, is an award given to up to five outstanding	PEMSCL's logits	ATLOP's logits
		•	<u> </u>	ATLOP + BCE* (Zhou and Lee, 2022)	75.86 ± 0.13	75.25 ± 0.11	75.91	75.36	alaga in any DEMCCL is much	instrumentalists each year	26.4	
DocRED	3,053 / 1,198,650	1,000 / 396,790	1,000 / 392,158	ATLOP (Zhou et al., 2021)	76.88	77.63	76.94	77.73	class in our PEMSCL is much	2. The Career Grants are a part of the Avery Fisher	9.2	16.8
Re-DocRED	3,053 / 1,193,092	500 / 193,232	500 / 198,670	DocuNet (Zhang et al., 2021)	77.53	78.16	77.27	77.92	lager than that of the ATLOP.	Artist Program, along with the Avery Fisher Prize and	Country NA	Country NA
				KD-DocRE (Tan et al., 2022a)	77.92	78.65	77.63	78.35	S	Special Awards.	´	
Our new data	regimes								 Our PEMSCL can correctly 	3. They are administered by the Lincoln Center for the	Label 2: (Avery Fisher Artis	st Program, U.S. , Country)
OOG-DocRE	3,053 / 1,198,650	$500^{3}/195,682$	500 / 198,670	NCRL* (Zhou and Lee, 2022)	78.41 ± 0.21	79.15 ± 0.20	<u>78.45</u>	79.19		Performing Arts.	PEMSCL's logits	ATLOP's logits
	3,053 / 1,198,650	500 / 193,232	,	PEMSCL (Ours)	79.02±0.20	79.89±0.17	79.01	79.86	predict relation that the	 5. Only U.S. citizens or permanent residents are eligible.		9.3 9.9
									- ATLOP model fails to identify.		Country NA	Country NA